

Washington University in St. Louis

## Washington University Open Scholarship

---

Arts & Sciences Electronic Theses and  
Dissertations

Arts & Sciences

---

Winter 12-15-2022

### Genetic and Transcriptomic Aspects of Major Depressive Disorder: In Vivo Functional Assays of Risk-Associated Variation, Candidate Disease Cell Types, and Their Pharmacologic and Sex Interactions

Bernard Mulvey

*Washington University in St. Louis*

Follow this and additional works at: [https://openscholarship.wustl.edu/art\\_sci\\_etds](https://openscholarship.wustl.edu/art_sci_etds)



Part of the [Genetics Commons](#), [Neuroscience and Neurobiology Commons](#), and the [Psychiatric and Mental Health Commons](#)

---

#### Recommended Citation

Mulvey, Bernard, "Genetic and Transcriptomic Aspects of Major Depressive Disorder: In Vivo Functional Assays of Risk-Associated Variation, Candidate Disease Cell Types, and Their Pharmacologic and Sex Interactions" (2022). *Arts & Sciences Electronic Theses and Dissertations*. 2806.  
[https://openscholarship.wustl.edu/art\\_sci\\_etds/2806](https://openscholarship.wustl.edu/art_sci_etds/2806)

This Dissertation is brought to you for free and open access by the Arts & Sciences at Washington University Open Scholarship. It has been accepted for inclusion in Arts & Sciences Electronic Theses and Dissertations by an authorized administrator of Washington University Open Scholarship. For more information, please contact [digital@wumail.wustl.edu](mailto:digital@wumail.wustl.edu).

WASHINGTON UNIVERSITY IN ST. LOUIS

Division of Biology and Biomedical Sciences  
Neurosciences

Dissertation Examination Committee:

Joseph Dougherty, Chair

Barak Cohen

Christina Gurnett

Celeste Karch

Joshua Rubin

Genetic and Transcriptomic Aspects of Major Depressive Disorder: *In Vivo* Functional Assays of  
Risk-Associated Variation, Candidate Disease Cell Types, and Their Pharmacologic and Sex  
Interactions

by

Bernard John Mulvey

A dissertation presented to  
Washington University in St. Louis  
in partial fulfillment of the  
requirements for the degree  
of Doctor of Philosophy

December, 2022  
St. Louis, Missouri



© 2022, Bernard Mulvey

# Table of Contents

List of Figures .....	vi
List of Tables.....	viii
Acknowledgments .....	ix
Abstract.....	xiii
<b>Chapter 1: Introduction .....</b>	<b>1</b>
<b>1.1 Abstract.....</b>	<b>1</b>
<b>1.2 Major depressive disorder (MDD).....</b>	<b>2</b>
1.2.1 Epidemiology, clinical presentation, and treatment of MDD.....	3
1.2.2 Candidate molecular, cellular, and brain mechanisms of MDD.....	5
1.2.3 MDD heritability and genetics .....	8
<b>1.3 Common noncoding variant discovery and association in psychiatry.....</b>	<b>8</b>
1.3.1 MPRA for identification of functional regulatory elements and variants.....	11
1.3.2 MPRA identify functional elements in specific cellular contexts.....	13
1.3.3 MPRA, assay context, and functional <i>variants</i> : MPRA can be designed not only to identify functional elements, but to assay and compare genetic variants in contexts known—or predicted—to mediate disease.....	16
1.3.4 MPRA: limitations and design considerations.....	18
1.3.5 Complementary methods in high-throughput study of DNA and RNA regulatory elements .....	22
1.3.6 The utility of MPRA for parsing linked variation.....	29
1.3.7 MPRA as an avenue to dissect multiallelic and polygenic mechanisms of neuropsychiatric traits.....	30
<b>1.4 Approach.....</b>	<b>34</b>
<b>1.5 Acknowledgments for Sections 1.3 and 1.6.....</b>	<b>35</b>
<b>1.6 Supplementary Material.....</b>	<b>37</b>
<b>1.7 References .....</b>	<b>39</b>
<b>Chapter 2: Molecular and functional sex differences of noradrenergic neurons in the mouse locus coeruleus.....</b>	<b>55</b>
<b>2.1 Introduction .....</b>	<b>55</b>
<b>2.2 Results .....</b>	<b>58</b>
2.2.1 Generation and validation of reagents for transcriptional profiling of	

noradrenergic neurons.....	58
2.2.2 Transcriptional response of NE neurons to LPS can be identified with <i>Slc6a2</i> TRAP .....	62
2.2.3 Sex-differential LC genes and putative <i>cis</i> -elements underlying differential regulation....	65
2.2.4 Molecular differences predict functional differences between sexes .....	66
<b>2.3 Discussion .....</b>	<b>69</b>
<b>2.4 Experimental Procedures .....</b>	<b>73</b>
<b>2.5 Acknowledgements .....</b>	<b>76</b>
<b>2.6 Supplementary Material.....</b>	<b>77</b>
<b>2.7 References .....</b>	<b>83</b>
<b>Chapter 3: From sex differences in gene expression to sex differences in gene regulation..</b>	<b>86</b>
<b>Chapter 4: Transcriptional-regulatory convergence across functional MDD risk variants identified by massively parallel reporter assays .....</b>	<b>88</b>
<b>4.1 Introduction .....</b>	<b>89</b>
<b>4.2 Methods.....</b>	<b>94</b>
4.2.1 Identifying candidate psychiatric GWAS regulatory variants.....	94
4.2.2 Massively parallel reporter assays .....	95
4.2.3 Targeted sequencing of RNA and input plasmid.....	96
4.2.4 Analysis.....	96
4.2.5 MotifbreakR analysis and functional SNP enrichment for perturbed motifs .....	98
4.2.6 Analysis of functional SNP-enriched TF expression in human brain and chromatin immunoprecipitation (ChIP)-seq.....	98
4.2.7 Code availability.....	99
<b>4.3 Results .....</b>	<b>99</b>
4.3.1 Many MDD loci contain more than one functional SNP .....	99
4.3.2 Shared regulatory architecture across distinct loci .....	105
4.3.3 Retinoids unmask additional functional SNPs in MDD loci.....	107
4.3.4 Retinoids reveal additional axes of convergent regulation at functional MDD-associated SNPs at levels of TF and cell type .....	110
4.3.5 Integrative analysis of TF sets at functional variants: TF binding, spatiotemporal Brain enrichment, and putative target genes .....	111
<b>4.4 Discussion .....</b>	<b>114</b>
<b>4.5 Acknowledgments .....</b>	<b>118</b>
<b>4.6 Supplementary Materials .....</b>	<b>119</b>

4.6.1	Supplementary methods.....	119
4.6.2	Supplementary figures and tables.....	132
<b>4.7</b>	<b>References .....</b>	<b>138</b>
<b>Chapter 5: Sex significantly impacts the function of major depression-linked variants <i>in vivo</i> .....</b>		
<b>5.1</b>	<b>Introduction .....</b>	<b>149</b>
<b>5.2</b>	<b>Results .....</b>	<b>151</b>
5.2.1	Combining Translating Ribosome Affinity Purification (TRAP) with MPRA.....	151
5.2.2	Identification of rSNPs and their sex-allele interactions in total hippocampus and excitatory neurons.....	153
5.2.3	rSNPs in the adult hippocampus and hippocampal <i>Vglut1</i> <sup>+</sup> neurons.....	156
5.2.4	Sex-interacting rSNPs in hippocampus.....	157
5.2.5	Transcriptional-regulatory systems shared across hippocampal rSNPs.....	159
5.2.6	Transcriptional-regulatory systems implicated in SxG interactions at MDD rSNPs in hippocampus .....	160
5.2.7	Identification of rSNPs in developing whole mouse brain.....	162
5.2.8	Sex-allele interactions are widespread neonatally but absent during hormonal quiescence.....	165
5.2.9	Transcriptional-regulatory systems implicated in brain-wide rSNP function in postnatal development .....	165
5.2.10	Landscape of functional variation differs broadly across age, sex, and brain region/cell type .....	168
<b>5.3</b>	<b>Discussion and Conclusion .....</b>	<b>169</b>
<b>5.4</b>	<b>Acknowledgments and additional data.....</b>	<b>173</b>
<b>5.5</b>	<b>Supplementary Materials .....</b>	<b>174</b>
5.5.1	Methods.....	174
5.5.2	Additional highlighted findings from TF analysis of rSNPs in P0 and P10 brain.....	201
5.5.3	Annotation datasets and other outside datasets .....	202
5.5.4	Supplementary figures and tables.....	203
<b>5.6</b>	<b>References .....</b>	<b>210</b>
<b>Chapter 6: Discussion .....</b>		
<b>6.1</b>	<b>Overview of findings.....</b>	<b>219</b>
<b>6.2</b>	<b>LC sex differences: human validity and transcriptional-regulatory mechanisms.....</b>	<b>220</b>

<b>6.3</b>	<b>Further dissection of functional MDD risk variants .....</b>	<b>222</b>
6.3.1	Defining variant-perturbed enhancers and repressors; minimal promoter selection in MPRA.....	223
6.3.2	Validation/demonstration of predicted regulatory factor mechanisms .....	224
6.3.3	Confirming the role of sex hormones in SxG variants from the adult hippocampus.....	224
6.3.4	Leveraging <i>in vivo</i> MPRA to explore consequences of other environmental factors on transcriptional-regulatory variants.....	226
6.3.5	Single-cell MPRA: opportunities ahead.....	227
<b>6.4</b>	<b>Conclusions .....</b>	<b>230</b>
<b>6.5</b>	<b>References .....</b>	<b>231</b>

# List of Figures

Figure 1.1: Example Allele-Differential Phenomenon in Common MPRA Approaches, and Analysis of MPRA Data.....	12
Figure 1.2: Regulatory Assays are Influenced by a Range of Conditions, from Environment to Sequence Context .....	15
Supplementary Figure 1.1: Common Designs and Library Delivery in MPRA .....	37
Supplementary Figure 1.2: Example of a Hypothetical MPRA .....	38
Figure 2.1: Characterization of noradrenergic bacTRAP lines .....	60
Figure 2.2: Transcript and protein expression in LC neurons.....	61
Figure 2.3: Differentially-expressed genes by sex .....	64
Figure 2.4: Motifs discovered in conserved, noncoding regions near sex-DEGs .....	65
Figure 2.5: Sex differences in <i>Ptger3</i> expression can be reflected in LC-mediated behavior .....	68
Supplementary Figure 2.1: Examples of Allen Brain Atlas ISH scoring .....	78
Supplementary Figure 2.2: Pathway analysis of transcripts altered in hindbrain by LPS and of sex-differentially expressed transcripts in noradrenergic neurons .....	79
Supplementary Figure 2.3: Map of LC cannula placements from post-mortem brain tissue .....	80
Supplementary Figure 2.4: OFT data by estrous stage from one cohort of sulprostone-behavior mice.....	81
Supplementary Figure 2.5: Multidimensional scaling (MDS) and heatmaps illustrate the clustering of samples and gene sets .....	82
Figure 4.1: Design of an MPRA library to identify candidate functional SNPs in MDD loci ...	101
Figure 4.2: MPRA defines SNPs with a functional effect on gene expression.....	104
Figure 4.3: MPRA signal at SNPs disrupting putative retinoid TF motifs.....	106
Figure 4.4: Retinoid treatment alters transcriptional regulation and unmask additional functional variants .....	109
Figure 4.5: Distinct TFs underlying retinoid-dependent functional SNPs and implication of serotonergic neurons .....	112
Supplementary Figure 4.1: Cross-replicate correlations from initial N2A MPRA .....	132
Supplementary Figure 4.2: Extended motifbreakR results from the first MPRA .....	133
Supplementary Figure 4.3: Additional functional SNPs corresponding to an enriched TF motif group (NR3C1) from the first assay .....	134
Supplementary Figure 4.4: Extended motifbreakR results from the ATRA treatment MPRA ..	135

Supplementary Figure 4.5: Expression of TFs <i>GATA2</i> , <i>GATA3</i> , and <i>FEV</i> in a variety of human neuroblastoma cell lines at baseline, with retinoic acid treatment, and in melanoma cells and neural precursors for comparison .....	136
Figure 5.1: Proof-of-principle for cell type-specific MPRA <i>in vivo</i> .....	153
Figure 5.2: Experimental design, analysis plan, and quality control: adult mouse hippocampus and its excitatory neurons.....	155
Figure 5.3: Adult hippocampus rSNPs and complex context-dependent, polygenic architecture of the <i>RSRC1</i> locus .....	158
Figure 5.4: Shared regulatory architecture of rSNPs by sex, cell type, and sex-interacting SNP type.....	162
Figure 5.5: Validating the <i>in utero</i> MPRA delivery method, and identification of rSNPs and sex-interacting rSNPs in the developing brain .....	164
Figure 5.6: Regulatory architecture of rSNPs at P0 and P10, permutation analysis evaluating the number of detected SxG rSNPs, and the context-dependent landscape of MDD loci .....	167
Supplementary Figure 5.1: Template MPRA oligonucleotide.....	203
Supplementary Figure 5.2: Barcode random effect coefficients are consistent within sex across ages/cell types.....	203
Supplementary Figure 5.3: Barcode random effect coefficients are consistent between sexes..	204
Supplementary Figure 5.4: Example of barcode random effect fitting on expression values ....	204
Supplementary Figure 5.5: Comparison of Benjamini-Hochberg corrected p-values and empirical p-values of allelic differences in each single-sex analysis, and of SxG interactions in each interaction analysis .....	205
Supplementary Figure 5.6: Log2FC between results from main adult female hippocampal TRAP experiment and a validation set of three additional female hippocampal-TRAP sample sets ....	205
Supplementary Figure 5.7: Basal (minimal promoter alone) barcode expression values do not vary by sex in total hippocampus or <i>Vglut1</i> <sup>+</sup> TRAP.....	206
Supplementary Figure 5.8: IF negative controls for P0 and P10 AAV9 delivery.....	207
Supplementary Figure 5.9: Basal (minimal promoter alone) barcode expression values do not vary by sex in P0 or P10 whole brain .....	207
Supplementary Figure 5.10: Absolute MPRA SNP effects correlate more strongly to MDD GWAS summary statistic effects for SNPs directly genotyped than for those imputed in the Howard 2019 meta-GWAS of MDD .....	208
Figure 6.1: Example MPA oligo product for a hypothetical MPRA library with co-delivered shRNAs .....	229

# **List of Tables**

Table 1.1: Strengths and Limitations of Functional-Regulatory Assays in Terms of Throughput and Sequence and Cellular Contexts .....	25
Supplementary Table 2.1: Antibodies used for immunofluorescence and ISH .....	77
Supplementary Table 4.1: Power analysis for allelic effect size detection.....	136
Supplementary Table 4.2: Number of functional SNPs, by effect type(s), overlapping at least 1 ChIP peak .....	137
Supplementary Table 5.1: Sequencing preparation steps and QC.....	209
Supplementary Table 5.2: qPCR primers used for <i>Vglut1</i> <sup>+</sup> TRAP validation.....	209



# Acknowledgments

This dissertation and my arrival at its defense are the product of mentorship, guidance, and support, my gratitude for which could not be fit into a dissertation all its own. Despite my tendencies otherwise, I will try and *briefly* highlight and heap praise on the many people without whom this work would not have been possible.

Joe Dougherty is an exemplary mentor and scientist who took me in from an undergraduate research career on the outer edges of rodent behavior and genetic methodology, fiercely interested but inexperienced in molecular neuropsychiatry. Thanks to Joe, my intellectual interests then (and ones yet unrecognized) have become my areas of expertise. It has been amazing being in his lab the past six years, and seeing how despite its astonishing growth, he still has time to meet with all his trainees at least once a week, still turns around his feedback on drafted manuscripts in 2-3 days, keeps up with everyone's off-the-cuff data dumps and messages, all while getting into the wet lab now and again, *and* getting the grants to make it all happen. Joe's unbridled enthusiasm for science in and beyond his lab kept me not only motivated, but excited about my work through the years, even in those long spells worth nothing seems to work. His contagiously unflappable good nature shaped his entire lab of supportive, high-spirited scientists, and it has been a pleasure to learn and grow in this environment through the years. To Joe, my deepest thanks for playing such a formative role in my development into a scientist, academic, and future mentor.

My entire thesis committee—Josh Rubin, Barak Cohen, Celeste Karch, and Christina Gurnett—is also deserving of great thanks for their involvement in shaping and directing this work. Barak and I met very early into my rotation through Joe’s lab, when I sat in on a meeting on experiments trying to leverage massively parallel reporter assays in the mouse brain (which I do here) and was at a total loss for the entire hour. Little did I know that Barak and his lab would soon enough teach me everything I needed to know to use these same assays to study psychiatric genetics and provide a key source of feedback on my use of their innovations in this area. Barak and his lab, especially PhD alums Brett Maricque and Dana King, and current PhD candidate Ryan Friedman, are all to thank for guidance through the design, preparation, and troubleshooting stages (i.e., the majority of time spent on) the work presented. Celeste Karch joined the committee as I began preparing to pursue *in vitro* studies of sex differences using the assays developed under Barak’s guidance. Toward this end, Celeste and her lab provided substantial hands-on training of myself and labmates over a year or more, helping us to learn how to differentiate male and female neural progenitors derived from induced pluripotent stem cells. While COVID-19 and graduations from our lab prevented us from seeing this goal to fruition, Dr. Karch’s material and time investments expanded my neuroscience skill set to include competence in a domain otherwise foreign to me: human cell culture—a skill set I hope to advance further in my future research. Christina Gurnett was on the committee from the start to guide the *in vitro* neuroscience and computational/disease genetics aspects of my thesis aims, and to provide me and the committee perspectives of a physician scientist on both my own career path and research timeline. Though the computational genetics aspect of the work hardly took form until the last couple of committee meetings, Dr. Gurnett stayed the course and helped shape some of the statistical approaches used in the work here, and helped me and the committee keep an eye on the bigger picture of my MD/PhD training path. Finally, this

committee was chaired by Josh Rubin, who I met at the suggestion of Dr. Gurnett based on my interest in sex differences. Josh is one of the only scientists I've worked with whose jolliness may actually exceed that of Joe, and was a pleasure to work with meeting to meeting. Over the years, he has helped me weigh the career paths ahead of me based on his own—very similar—experiences from his training, easing the weight of my worries about the future (weight which would likely have cumulatively prevented at least one of the chapters here from reaching completion). To all of my committee, thank you.

I am also greatly indebted to several labmates and lab alumni for their intellectual and material support over the years. The entire reporter assay team in Joe's lab provided a great platform and desk area to constantly reconsider, innovate, collaborate on, and run these experiments, especially thanks to Tomas Lagunas, Jr., Ph.D., Michael Rieger, Ph.D., Tony Fischer, Ph.D., Stephen Plassmeyer, and Din Selmanovic. Deserving of an omni-thanks for providing one-on-one and lab meeting-based insights into mouse behavior research, sex differences, and for helping me to keep my time and materials organized is Susan Maloney, Ph.D. Likewise, thanks to Katie McCullough, who led the way with the behavior experiments and subsequent study of the brain tissue presented in Chapter 2. My earliest mentor in the lab was Kristina Sakers, Ph.D., who taught me to design and perform *in situ* hybridizations of mouse brain tissue and use a confocal laser microscope, bringing me full-bore into the world of neurogenetics. She later gave me her ultra-accurate pipettes, which I would argue are the only reason that experiments in Chapters 4-5 worked out, so for all this, thank you, Krissy. I also must thank alumna programmer Allie Lake—now my fiancée, so thanks once more, Joe—I first forayed into programming and computational genetics, interests which visibly grow and mature throughout this work. My last individual thanks are to the two

undergraduate/post-baccalaureate researchers in the lab who tirelessly helped with experiments trying to utilize human brain tissue for isolation of locus coeruleus RNA: Selma Avdagic and Haley Crosby.

Finally, there are four mentors before my time at Wash U whom I would be remiss not to thank here. Tom Jeffries, Ph.D., at UW Madison brought me into molecular genetics through a side door of yeast bioengineering, where I first learned and came to love many of the base techniques I now use in the setting of mammalian genetics. Tom Sutula, M.D., Ph.D., also at UW, mentored me in functional and pharmacologic epilepsy research throughout my entire four years, introducing me to the world of neuroscience and indeed to the concept of dual-track training. So to Tom and Tom, thank you. Dating even further back, I have to thank two high school science teachers, Mimi Verhoeven and Julie Fangmann, for encouraging me to embrace and invest in my affinity and enthusiasm for science.

I would like to dedicate this work to all those who have or do struggle with major depressive disorder. It is for you that I do this research: in hopes that the field of psychiatry can one day better identify, understand, and treat the underlying biological processes of MDD; it is for a world free of hopelessness beyond hope.

Bernie Mulvey

*Washington University in St. Louis*

*December 2022*

ABSTRACT OF THE DISSERTATION

Genetic and Transcriptomic Aspects of Major Depressive Disorder: *In Vivo* Functional Assays of Risk-Associated Variation, Candidate Disease Cell Types, and Their Pharmacologic and Sex Interactions

by

Bernard Mulvey

Doctor of Philosophy in Biology and Biomedical Sciences

Neurosciences

Washington University in St. Louis, 2022

Professor Joseph D. Dougherty, Chair

Major depressive disorder (MDD) is a debilitating illness that affects hundreds of millions globally, with substantial personal, medical, economic, and societal consequences. While depression occurs more commonly in females, the biology of the brain and sex underlying this skewed prevalence remains unclarified. This body of work explores two aspects of how biological sex may influence the brain at the level of gene expression: through intrinsic sex differences and through sex-mediated effects of depression risk genetics.

The monoamine hypothesis of depression suggests that modulatory neurotransmitters including serotonin and norepinephrine constitute a key axis in development of MDD. Large-scale studies of MDD treatment response have found that women respond better to serotonergic agents, while males respond better to mixed serotonergic-noradrenergic agents, suggesting one or both of these cell types may play a role in sex-differentiated MDD risk biology. Using translating ribosome

affinity purification (TRAP), gene expression in norepinephrine neurons of mouse locus coeruleus (LC) was profiled and compared across sexes, revealing over 100 genes with both sex-differential and LC-enriched expression. Three female-upregulated genes of interest emerged: *SLC6A15* and *LIN28B*, implicated in MDD, and prostaglandin receptor *PTGER3*. Pharmacologic activation of *PTGER3* had female-specific effects on LC electrophysiology and behavior, confirming that genetic sex differences in noradrenergic neurons have functional consequences on these neurons and behavior.

The role of genetic variation in MDD has recently come to be appreciated as an underlying cause of MDD, though whether sex interacts with genetic risk factors remains unknown. The primary work in this thesis focused on over 1,000 noncoding, putatively transcription-regulating common variants from 31 MDD-associated genomic regions—including those near *LIN28B* and *SLC6A15*—using functional assays in mouse brain *in vivo* to examine sex-by-genotype interactions. This work identified extensive sex-by-allele effects in mature hippocampus and, using TRAP, its excitatory neurons in particular. Unbiased informatics approaches indicated a role for nuclear hormone receptors, further supported by comparative analysis of analogous experiments in neonates during the masculinizing testosterone surge and in older, hormonally quiescent juveniles. This study provides novel insights into MDD genetics as influenced by age, biological sex, and cell type, and provides a framework for *in vivo* parallel assays to functionally define interactions between disease-linked genetic variation and complex biological or environmental variables.

# Chapter 1: Introduction

Sections 1.3, 1.5, and 1.6 of this chapter were previously published:

Mulvey, B., Lagunas, T. & Dougherty, J. D. Massively Parallel Reporter Assays: Defining Functional Psychiatric Genetic Variants across Biological Contexts. *Biol Psychiat* (2020) doi:10.1016/j.biopsych.2020.06.011.

## 1.1 Abstract

Sex differences in numerous psychiatric disorders have been recognized for decades at the level of patient populations, including differences in their incidence/prevalence, natural history, progression, and treatment responsiveness. The molecular bases of such differences, however, remain largely unknown across psychiatry. Rapid advances in genetic and functional neuroscience techniques have provided glimpses into possible cellular and molecular mechanisms of sex differences, especially at levels of gene expression, regulation, and function. These technologies have been applied in human postmortem and model organism brain, affording a new depth of resolution in the neurogenetics of health, disease, and their interplay with sex. The vast majority of this work, however, has targeted single brain regions and often single genes, limiting the scope of genetic-molecular investigations of how sex and the brain interact. To contribute broader insights into the interplay of sex and genetics within the brain, I first employed existing molecular-genetic techniques to characterize innate sex differences in an undercharacterized cell population implicated by pharmacotherapeutic strategies in depression: the noradrenergic locus coeruleus. Subsequently, I adapted high-throughput gene-regulatory assays to examine genetic risk loci for depression both for pharmacologic interactions *in vitro* and biological sex interactions in the mouse brain *in vivo*. These experiments illustrated that most depression-associated loci contain several

functional variants, and that sex interactions with variant regulatory effects are only observable during periods of sex hormonal circulation. Altogether, this body of work illustrates that sex differences are widespread in the molecular brain, with impacts spanning cellular physiology to behavior and disease risk.

## **1.2 Major depressive disorder (MDD)**

Major depressive disorder (MDD) is a widespread psychiatric disorder, affecting hundreds of millions of people worldwide during the course of their lifetime. Studies of the disorder—spanning from levels of populations to single genomic nucleotides—point to evident sex differences in the prevalence, severity, recurrence, and in the underlying tissue, cellular, and molecular phenomenon observed in the disorder. Nonetheless, mechanisms for the influence of biological sex on MDD remain a vexing problem in biological and molecular psychiatry. Human and model organism research has implicated several brain regions and/or cell types in both the general and sex-differentiated pathophysiology of depression, as have pharmacologic approaches to MDD treatment. Meanwhile, studies of genomic variation in large cohorts of MDD cases and population-matched controls have discovered over one hundred genomic regions associated with having the disorder, though these studies cannot identify the precise variants that are causal for disease association. To characterize how sex, gene expression, and genetic risk loci may influence the molecular brain with respect to MDD, I utilized and developed molecular techniques: 1) to characterize sex differences in gene expression of norepinephrine-releasing neurons of the brain; 2) to experimentally identify variants from MDD-associated loci which influence gene regulation and their interplay with retinoid signaling *in vitro*; and 3) to identify variants from these same MDD-associated loci affecting gene regulation in a sex-and-cell-type specific manner in the mouse



brain *in vivo*. This introduction first describes clinical, putative biological, and genetic aspects of MDD, emphasizing sex differences throughout. Subsequently, a primer in genetic association studies, their limitations, and biological questions they raise about psychiatric disease heritability is provided as motivation for the experiments presented.

### **1.2.1 Epidemiology, clinical presentation, and treatment of MDD**

Major depressive disorder (MDD) affects between to 5-15% of people in the United States<sup>1,2</sup>, with approximately 2/3 of lifetime diagnoses occurring in women<sup>3-5</sup>. The disorder is responsible for a substantial proportion of total life years effectively lost<sup>6</sup> due to health-impaired social and occupational function, extending the toll of the disease on patients to families, friends, and communities. Moreover, MDD has been repeatedly associated with up to two-fold increases in cardiovascular disease risk, morbidity, and mortality<sup>7-11</sup>, compounding patient- and society-level burdens of the disease.

MDD is characterized by nine primary features: feelings or manifestations of sadness/hopelessness, suicidal ideation or behavior, loss of interest (anhedonia), fatigue, excessive feelings of guilt, impaired concentration, and more heterogenous disruptions to psychomotor function (agitation or retardation), appetite/weight change (gain or loss), and sleep (insomnia or hypersomnia). Any five of these symptoms present on a near-daily basis are indicative of depression, resulting in hundreds of observable diagnostic MDD phenotypes<sup>12</sup>. MDD generally occurs in an episodic manner, once or more in a lifetime, with symptomatic period(s) of months to years followed by remission. Earlier age of onset and recurring episodes over the lifetime are indicative of more severe disease in both clinical<sup>13,14</sup> and biological<sup>15-17</sup> terms.

Two commonly identified patterns of depressive symptoms include “atypical” and “melancholic” depression, the former featuring hyperphagia, hypersomnia, and psychomotor retardation, while the latter features appetite loss and insomnia. Atypical depression is more common in females, while melancholic depression is more common in males<sup>18-21</sup>, illustrating a role for biological sex underlying disease manifestations. A third form of the disorder, postpartum MDD, provides an even stronger case for the role of biological sex in disease, as illustrated by a novel treatment mechanism, discussed below.

One of the first pharmacotherapeutic approaches to major depression was monoamine oxidase inhibitors (MAOIs), which were incidentally discovered in the process of developing treatments for tuberculosis<sup>22</sup>. The efficacy of these drugs—which increase brain monoamine neurotransmitters dopamine, serotonin, and norepinephrine—was a major contributor to the “monoamine hypothesis” of depression<sup>23,24</sup> and the development of (most) other pharmacotherapeutics for the disorder. The three main classes of antidepressants commonly used in treatment today are, in order of their development: tricyclic antidepressants (TCAs), which elevate extracellular serotonin and norepinephrine and modulate levels of several additional neurotransmitters; selective serotonin reuptake inhibitors (SSRIs), which selectively increase extracellular serotonin; and selective serotonin-norepinephrine reuptake inhibitors (SNRIs), which increase extracellular serotonin and norepinephrine without the myriad off-target effects of their TCA predecessors. Monoamine-releasing (or -receptive) areas of the brain likely play a role in sex differences in depression, as males show better response to TCAs, while females respond better to SSRIs<sup>25</sup>. Interestingly, the intersection of sex differences in depression and cardiovascular

comorbidities has been hypothesized to originate from dysregulated autonomic (adrenergic, noradrenergic, and cortisolic) functions and their modulation by the hippocampus<sup>10</sup>.

Novel treatments for depression have also emerged in the last decade, including the rapidly-acting antidepressant ketamine, which transiently blocks glutamate receptors, resulting in days or weeks of symptom relief near-immediately after each dose<sup>26</sup>. More notably for our purposes, brexanolone, approved specifically for postpartum MDD, is a progesterone metabolite administered intravenously in the weeks or months after delivery. The drug modulates inhibitory signaling in the brain via GABA receptors<sup>27</sup>, and likely plays a role in modulating sex hormonal signaling given the dysregulation of female sex hormones in postpartum MDD<sup>28</sup>.

### **1.2.2 Candidate molecular, cellular, and brain mechanisms of MDD**

Studies in depression patients have identified a number of potential features of the disease at the structural level of the brain. Chief among these is a loss of hippocampal volume, an observation replicated in several cohorts to date<sup>15,29–31</sup>. The hippocampus plays numerous roles in cognition and behavior, including memory formation and dampening of fear and stress circuit activity. As is the case for most psychiatric disorders, changes in subregions of the frontal cortex have also been observed and reviewed<sup>31,32</sup>. Sex differences are manifest in gross-anatomic and molecular-pathology brain measures in MDD, including surface area and volume of the prefrontal cortex (PFC), with opposite disease effect directions between sexes<sup>33</sup>; female-specific loss of functional connectivity between frontal and parietal cortices in depression with inattentive symptoms<sup>34</sup>; and greater degree of volume loss in hippocampal subregions in females<sup>29</sup>.

The treatment approaches discussed above highlight several candidate cell types of the brain which may play direct roles in MDD, or whose modulation may compensate the (unknown) primary dysregulation. The primary monoamine targets of MDD drugs—serotonin and norepinephrine—are released from small populations of neurons with wide-ranging projections throughout the brain, respectively the raphe nuclei and the locus coeruleus. Given the sex differences in antidepressant response discussed above, these two populations constitute candidate cell types involved in sex-differential risk and presentation of MDD.

Studies of postmortem brain gene expression have, principally, identified extensive sex differences across cortical and subcortical regions in MDD, including the hippocampus and frontal cortices<sup>35–37</sup>. Work in healthy and MDD human tissue and in human cell types *in vitro* have nominated several additional cell types on the basis of gene expression and regulation, discussed more below. In brief, cell types implicated by genetic assays across model organisms and human tissues include several excitatory neuron layers of the PFC<sup>38–41</sup>, inhibitory and granule **neurons of the hippocampus**<sup>42–44</sup>, **and fetal excitatory and inhibitory neurons**<sup>45,46</sup>. Importantly, excitatory neurons as a broad cell type have been implicated by several forms of genetic and intersectional analysis (*see Chapter 5 for more discussion of this cell type*); moreover, the glutamatergic system shows sex differences in human brain levels of the transmitter and expression of its receptors, hormone-mediated levels of pathway activity, and sex-specific behavioral responses to its modulation (reviewed in<sup>47</sup>).

Rodent models of depression generally revolve around the use of chronic stressors and/or genetic perturbations to induce depression-like behavior. In these paradigms, females are more susceptible to stress induction of depressive phenotypes, showing behavioral changes with several fewer days

of stress than males<sup>35,48–50</sup>, reinforcing the roles for both the environment (stress) and biological sex in MDD. Consistent with human postmortem studies, multiple brain regions of the mouse including cortex and nucleus accumbens exhibit broad sex differences after chronic stress<sup>51</sup>. While monogenic forms of MDD have not been described in the literature, regional perturbations of candidate genes in rodent **hippocampal excitatory neurons**, of *Crebl1* and *Ppara* in **total hippocampus**, and myriad other genes in **PFC**, nucleus accumbens, and raphe have been able to recapitulate depressive phenotypes<sup>35,52–57</sup>, including sex-specific behavioral effects<sup>35,55</sup>.

Finally, given the sex differences in MDD, several research groups have directly investigated the roles of sex hormones in human MDD and rodent models thereof. Treatment of postpartum MDD with a single three-day course of brexanolone (a progesterone metabolite, as noted above) indicates sex hormones as a proximal cause of this subtype of disease. For MDD more broadly, low serum testosterone in males has been associated with a 1.5-fold increase in five-year odds of developing MDD<sup>58</sup>. Rodent model studies of hormonal signaling in the brain have identified sex differences in hippocampal neurite outgrowth in response to estrogen<sup>59</sup>, protective effects of estrogen (locally converted from androgen) against stress susceptibility to depressive phenotypes in male mice<sup>60</sup>, and changes in behavior and hippocampal gene expression specific to stages of the estrus (female hormonal) cycle<sup>61–63</sup> consistent with robust epigenomic sensitivity of the female hippocampus to estrogens<sup>64</sup>.

Likewise, perturbations of rodent hormones using gonadectomy/ovariectomy with or without hormone replacement have shown consequent effects on candidate cell types in MDD, including in the norepinephrine nexus of the brain, the locus coeruleus (LC). There, the population of neurons

is larger in females or with developmental estrogen treatment, while it is smaller in males and androgen-treated females<sup>65,66</sup>; estrogens regulate expression of *Th* and *Dbh*<sup>67</sup>, the two key enzymes for the conversion of tyrosine to norepinephrine; the sexes exhibit opposite, sustained electrophysiologic responses to adolescent stress<sup>68</sup>; and females have greater receptor availability and response to the endocrine stress signal, corticotrophin releasing factor (CRF) (reviewed in<sup>69,70</sup>).

### **1.2.3 MDD heritability and genetics**

Concentration of MDD in families has been consistently observed, with approximately 35% heritability based on twin studies<sup>71-74</sup>. These findings imply that part of MDD risk is conferred through the genome, which has spurred much research into the disorder using genome-wide association studies (GWASes; see **Section 1.3** for further description). GWAS for MDD across multiple studies involving hundreds of thousands of cases<sup>75-79</sup> has identified population-level genetic associations explaining a similar degree of heritable risk, ranging 10-33%. Intriguingly, GWAS of clinically-ascertained or self-reported MDD required tens of thousands of cases to identify the first associated loci, except for an early GWAS of a *single-sex* (female) clinical cohort exclusively with recurrent, melancholic MDD, wherein two significant loci were identified using only five thousand cases<sup>80</sup>. That mixed-sex and -subtype GWAS required a substantially greater number of subjects to be powered for such detections suggests that the sexes, MDD subtypes, or both may have non-shared genetic risk factors. Nonetheless, these studies have cumulatively and fundamentally confirmed that there is an inherited, genetic basis for MDD.

### 1.3 Common noncoding variant discovery and association in psychiatry

Psychiatric diseases are genetically influenced by both heritable variation (common and rare) and non-inherited, *de novo* mutations. Estimated common variant (frequency  $\geq 1\%$ ) influence on disease liability ranges from 10-33% for major depressive disorder (MDD)<sup>74,78,81</sup> and schizophrenia (SCZ)<sup>74,82,83</sup> to over 50% in autism spectrum disorders (ASD)<sup>84,85</sup>. The remaining familial heritability of psychiatric—especially neurodevelopmental and psychotic<sup>86,87</sup>—diseases is largely conferred by rare variants<sup>85</sup>. Two major hurdles have prevented variant data from illuminating disease mechanisms: the volume of variant discoveries/associations to test for functionality and causality, and imperfect methods of predicting variant consequences.

Variant-disease associations arise from correlational methodologies. Genome-wide association studies (GWAS) identify overrepresented single nucleotide polymorphisms (SNPs), tagging hundreds of mostly non-protein coding, linked SNPs<sup>85</sup>. Similarly, family studies identify proband-specific (*de novo*) or -enriched (rare, inherited) variants, though few of these are causal for disease at the individual level. However, these statistical association-based approaches alone are incapable of specifying which variants have biological function.

Predicting whether and how noncoding variants are functional is a nontrivial enterprise. The majority of GWAS loci bear indirect indication(s) of transcriptional regulatory function, including expression quantitative trait locus (eQTL) associations, chromatin accessibility, or histone marks<sup>88-90</sup>. As others have noted, these data alone cannot define functional regulatory elements/variants<sup>91,92</sup>. However, even within one cell type, such data are often mutually discordant: an emerging (*i.e.*, preprinted) study examining six epigenomic datasets from K562

cells showed 49% of functional regulators did not overlap *any* epigenomic annotations; another 40% only overlapped one of the six<sup>93</sup>. Similarly, MPRA of chromatin-based K562 enhancer predictions found only 30% regulated transcription<sup>94</sup>. Unsurprisingly, these discrepancies apparently extend to disease variant interpretation: only a minority of GWAS variants (except for blood traits) overlap tissue-specific regulatory predictions<sup>95</sup> from histone marks<sup>96</sup>. Such findings collectively suggest that heritable, disease tissue-specific regulatory phenomenon are both missed and mislabeled when relying solely on chromatin states.

Despite the clear excess of *de novo* variation in coding sequences in ASD and other neurodevelopmental disorders, and though coding variant consequences can often be predicted (*e.g.*, nonsense mutations), this constitutes the minority of heritable risk for several psychiatric diagnoses<sup>97</sup>. The remaining burden falls within putative transcriptional<sup>97</sup> and translational regulatory elements (*e.g.*, promoters, UTRs)<sup>98,99</sup>. ASDs provide a representative case: among 1,902 subjects, over 255,106 *de novo* variants were recently identified, with thousands each in upstream/promoter sequences and untranslated regions (UTRs)<sup>100</sup>. UTRs regulate transcript stability and miRNA interactions<sup>101</sup>; emerging work further implicates UTRs in nuclear transcript trafficking in the brain<sup>102</sup>. The occurrence of most disease-linked variation in the least-well understood features of the genome/transcriptome thus obstructs understanding of disease biology. Collectively, these two problems necessitate **high-throughput assays with functional readout** for putative regulatory elements and variants. Such assays enable identification of functional variants and the biological contexts in which they act. This knowledge can shape hypotheses regarding shared mechanisms by which disparate genetic factors converge on shared phenotypic endpoints.



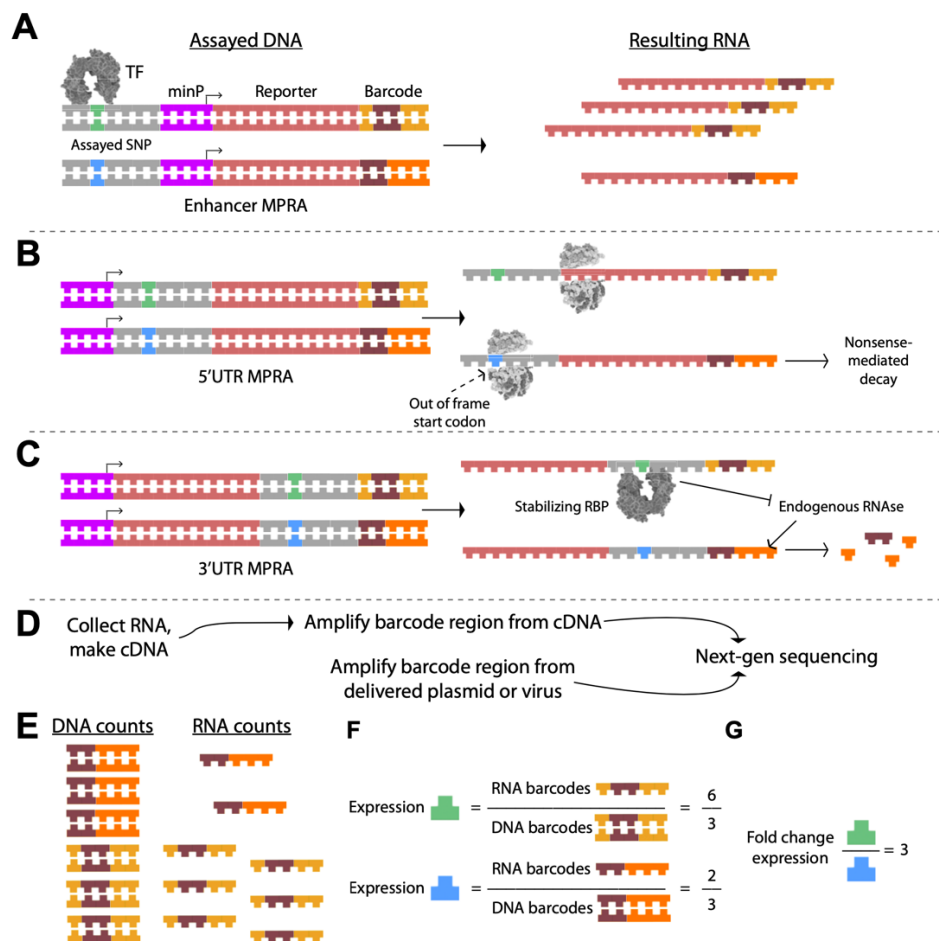
Here, I will primarily discuss MPRA for high-throughput parcellation of genetic discoveries. MPRA technology pairs genomic features (*e.g.*, each allele of a genomic sequence) to a reporter gene bearing unique, transcribed barcodes, allowing multiplexed RNA-level readout of element activity<sup>103,104</sup>. Critically, there is substantial potential for MPRA to identify functional variants from neuropsychiatric-associated loci. I first discuss uses of MPRA in functional identification of gene regulatory elements and variants, design/interpretation considerations for MPRA, and methods to complement/follow-up findings. Subsequently, I discuss potential applications of MPRA to identify mechanistic convergence across polygenic risk space.

### **1.3.1 MPRA for identification of functional regulatory elements and variants**

MPRA offer a flexible framework to study elements regulating transcription (*e.g.*, enhancers, promoters), splicing, protein translation, and post-transcriptional phenomenon. Though too numerous to review deeply here, I point readers to published and emerging applications of MPRA to splicing<sup>105–107</sup>, RNA editing<sup>108</sup>, and protein translation<sup>109</sup>. MPRA have been most broadly applied to explore and computationally model transcriptional “regulatory grammar”—how sequence features such as binding motifs, their abundance, and arrangement affect regulatory capacity<sup>94,110–116</sup>. More recently, these approaches have been applied to characterize UTR functions in RNA stability and translation<sup>117–121</sup>, and to identify SNPs and rare variants influencing transcription<sup>122–128</sup>.

As shown in **Supplementary Figure 1.1A-B**, a canonical ‘enhancer’ MPRA utilizes a promoter with candidate elements either upstream or in a 3’UTR (*i.e.*, for the STARRseq variant of MPRA)<sup>129</sup>. Each element is paired with unique barcodes in the transcribed UTRs, which are

sequenced as quantitative readouts. Expression—representing transcription or RNA stability—is measured as the RNA barcode counts per DNA barcode count (**Supplementary Figure 1.1E**). This measure can be leveraged to define active or differentially active enhancer elements. Functional elements have been defined by either comparing to minimal-promoter only barcodes<sup>94,112,115,116,130,131</sup>, or individual sequences against their shuffled counterparts<sup>110,116</sup>; MPRA have also successfully compared activity between alleles<sup>122,124–128,128,132</sup>.



**Figure 1.1. Example Allele-Differential Phenomenon in Common MPRA Approaches, and Analysis of MPRA Data.** **A)** In a transcriptional-regulatory assay, a putative regulatory SNP may create, ablate, strengthen, or weaken a TF binding site. As a result, one allele drives more transcription (detected via its 3'UTR barcodes) per encoding DNA than the other allele. **B)** In a 5'UTR assay, a functional SNP may sequence features controlling translation initiation. For example, a variant allele may introduce an upstream start codon out of frame with the reporter

gene, resulting in nonsense mediated decay, and thus, decreased detection of the barcodes paired to that UTR allele. **C)** In a 3'UTR assay, a variant may alter an RBP binding site; in this example, an RBP site specific to one allele increases the stability of the reporter transcript, and thus of the barcode paired to it. **D)** After transfection/transduction, RNA is collected from specimens and prepared along with DNA (often the delivered DNA, though sometimes this is recovered from the specimens as well) to generate sequencing libraries to quantify expression of the delivered elements in the RNA, compared to starting abundance in the DNA. **E)** Example read counts, presented visually, for the DNA and RNA barcode counts of one barcode paired to each allele. **F)** MPRA analysis centers on taking the ratio of RNA/DNA counts (or counts per million), represented by the sequence fragments at top left, as a measure of expression—*i.e.*, approximating the number of transcripts generated per encoding DNA. These can be compared relative to the expression of all elements to find the strongest features (e.g., strongest enhancers and repressors, or most stabilizing and destabilizing UTR elements), or **G)** compared on a variant-wise basis to determine significant allelic regulatory effects.

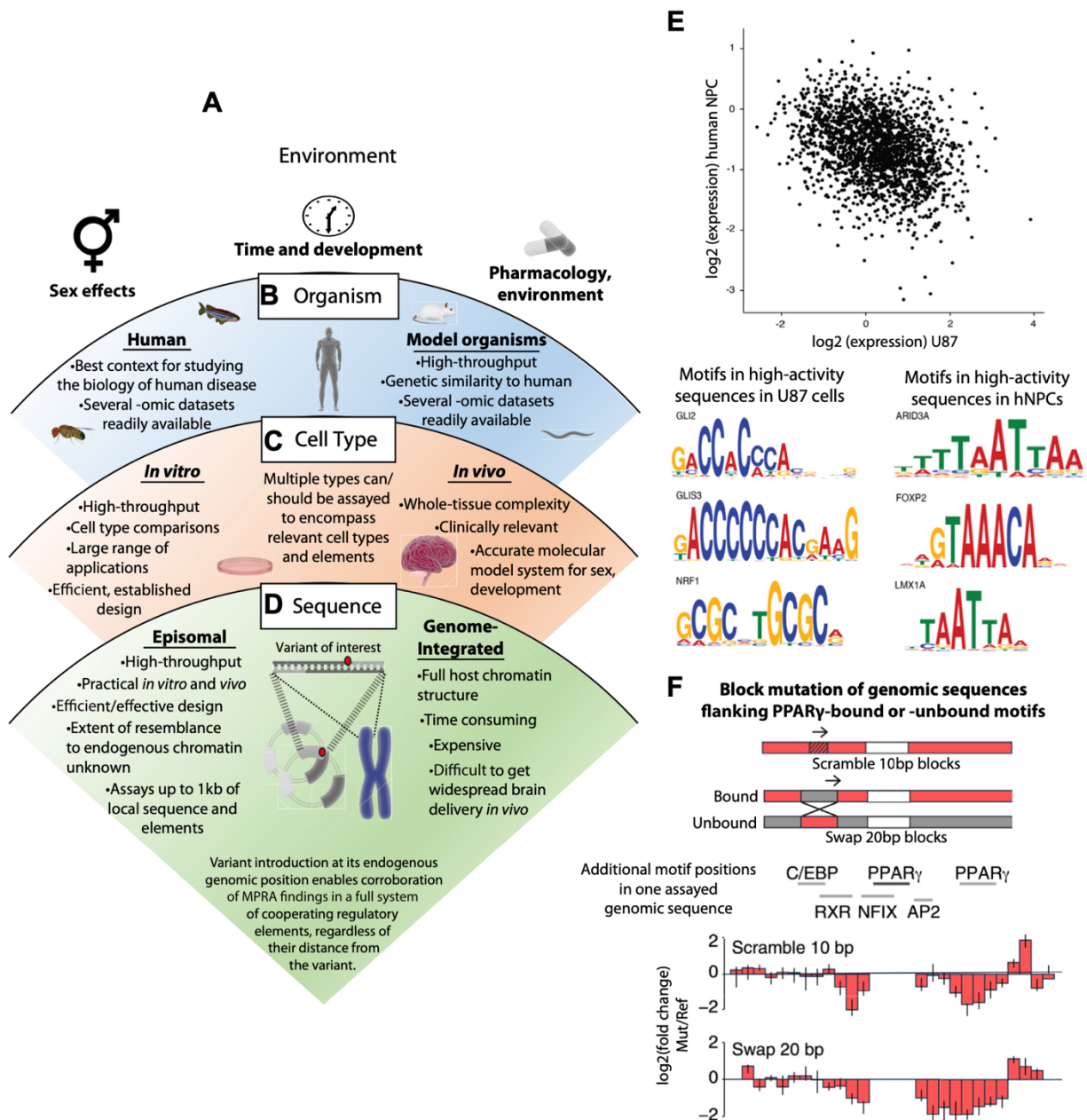
### 1.3.2 MPRA identifies functional elements in specific cellular contexts

Perhaps the most exciting—if underappreciated—property of MPRA is the ability to assay elements using disease-relevant cells and conditions. Functional elements are defined by each cell type's unique milieu of expressed TFs, chromatin modifiers, miRNAs, and RBPs, which mediate regulatory element activity. The breadth of published and emerging tissue/cell type differences in gene expression<sup>133,134</sup>, chromatin marks<sup>96,135</sup>, and chromatin interactions<sup>136–138</sup> all illustrate the magnitude of these regulatory differences. The importance of cell type was illustrated by an MPRA of the same elements in U87 glioblastoma and neural progenitor cells (NPCs): the most active enhancers in each cell type contained entirely different motifs and sequence features<sup>115</sup>. Recent<sup>139,140</sup> and emerging<sup>141</sup> approaches have identified highly cell type-specific brain enhancers using adeno-associated viral (AAV) vectors alongside traditional (*e.g.*, immunofluorescent) readouts. Moreover, a novel, AAV-based MPRA (*i.e.*, using barcodes) identified novel functional enhancers for somatostatin interneurons<sup>142</sup>. Aside from these examples, MPRA in neural cells have been limited to date. Several early MPRA utilized explanted retina to explore influences of TF binding sites and their arrangements on activity<sup>116,125,131,143</sup>. One novel study, relevant to functional contexts (discussed below), assayed mouse neuron enhancers for activity changes

following KCl depolarization<sup>144</sup>. Other studies include an MPRA characterizing temporal patterns of *cis*-regulatory element activity across seven timepoints in human NPC differentiation into neurons<sup>145</sup>. This delineation of regulatory element function illustrates the power of regulatory assays to reveal timepoints and cellular states wherein gene regulation—especially for neurodevelopmental disorders—may exert its causal effects.

*In vivo* regulatory assays—including in the brain—have more recently been demonstrated, generally at smaller scales than *in vitro* MPRA. Osterwalder, et. al<sup>146</sup> singly or multiply knocked out putative limb development enhancers in mice, illustrating enhancer redundancy—that is, limb development disruption only with perturbation of multiple elements. McClymont et. al<sup>147</sup> identified 2,000 candidate embryonic mouse enhancers in purified midbrain dopamine neurons, and validated the developmental and regional specificity of a subset using transgenic reporter mice. The scale of these assays has been expanded by groundbreaking implementation of MPRA in the brain *in vivo*<sup>125,142</sup> to query functional effects in native cell contexts.

This transition to *in vivo* MPRA is beneficial because, while cell type overwhelmingly influences regulatory assays, additional conditions may equally alter outcomes (**Figure 1.2A-C**). Age, sex, pharmacology, and environment (*e.g.*, stress)—all can shape gene expression. For example, MPRA have identified elements responsive to hormonal contexts such as steroid-responsive glucocorticoid receptor binding<sup>148</sup>. Altogether, MPRA enable identification of functional regulatory elements across varied internal and external environments.



**Figure 1.2. Regulatory Assays are Influenced by a Range of Conditions, from Environment to Sequence Context.** The range of conditions that influence regulatory assays (from top to bottom) starts when considering **A**) the environment, *e.g.*, sex, time, and pharmacology. These parameters have the potential to affect various –omic profiles in a given system. **B**) The next level of consideration is the organism, which can include human-derived tissue or one of the many model organisms. Human genomic context is ideal for studying the biology of human disease – though a comparatively limited scope of techniques for human-derived tissues exists. **C**) Next, one should consider the selected cell type(s) and whether to assay *in vitro* or *in vivo*. Each of these provides a unique set of benefits, and one approach can be used to validate findings from the other. In the case of modeling the brain and psychiatric genetic variants, cell type-specific/enriched MPRA *in vivo* would constitute the highest-fidelity model of variant effects by accounting for

regulatory effects of endogenously interacting cell types. **D)** Lastly, sequence (or chromatin) context is dictated by delivery method, yielding extragenomic or intragenomic MPRA DNAs. In either case, only limited length of sequence surrounding a feature of interest is preserved (*e.g.*, in ~120bp tiles of genomic sequence in custom oligonucleotide cloning, or  $\leq$  1kb in clone-and-capture methods), preventing assessment of any interactive effects from elements further away. (A recent study suggests that size of a tile negatively correlates with reproducibility of expression driven compared to that driven by ~120bp tiles, emphasizing the importance of this consideration (93)). While AAV-transduced episomes gain histones (105) and chromosome-like nucleosome spacing (106), it is unknown whether gene-regulatory histone marks on these episomes mirror those of endogenous regulatory chromatin. I suggest corroboration of MPRA findings in native genomic settings, by, for example, introducing the variant to the genome of a cell line using CRISPR methods to account for local sequence and chromatin structure effects on MPRA results. **E)** The consequences of cell type in a previous MPRA testing the enhancer potential of random sequences in human NPCs (hNPCs) and U87 glioma cells. Notably, activity of sequences in one cell type was negatively correlated to activity in the other cell type (left), and sequence motifs corresponding to high-activity enhancers for the cell types were strikingly distinct (right) (reproduced from<sup>115</sup>). **F)** Previously, Grossman, et. al<sup>110</sup> took a series of mouse genomic sequences with motif matches to the TF PPAR $\gamma$ , only some of which were actually bound by PPAR $\gamma$  in ChIP-seq of mouse adipose. They performed an MPRA on several such motif-and-binding genomic sequences, centered on the binding site, but shuffling the bases in 10bp windows surrounding the site, or swapping 20bp windows between bound and unbound, motif-centered sequences. Shown are results for one genomic sequence, illustrating differential MPRA expression between each tile mutant relative to the reference sequence.

### **1.3.3 MPRA, assay context, and functional *variants*: MPRA can be designed not only to identify functional elements, but to assay and compare genetic *variants* in contexts known—or predicted—to mediate disease.**

As transcriptomic and epigenomic studies highlight an enormous role for cell type, it is unsurprising that this influence extends to regulatory variants. For example, variants exert cell type-dependent effects on chromatin structure even within a neurodevelopmental lineage: an emerging study discovered chromatin accessibility QTLs in human NPCs and neurons, with ~80% of QTL SNPs being specific to one of the cell types<sup>149</sup>. Cell type roles in putatively functional variation are also implicated by GWAS SNPs enrichment in tissue-specific eQTLs<sup>134,150</sup>, neural cell type-specific chromatin interactions<sup>136</sup>, and eQTLs that evade detection in bulk brain (*i.e.*, multi-cell type) tissue but are evident in purified populations like dentate granule neurons<sup>151</sup>. MPRA has likewise demonstrated the essentiality of cell type in defining functional variants: the

Critical Assessment of Genome Interpretation 5 (CAGI5) consortium performed an MPRA on saturation-mutagenized human enhancers and disease-associated promoters in numerous cell lines, challenging analysts to computationally predict functionality and effect size for held-out variants. The most predictive annotations for a given cell line were often from the same cells across several top-performing analyses<sup>152</sup>. Thus, experimental study of putative disease-associated variants requires firm hypotheses on where (tissue/cell type), when (development/differentiation), and how elements are expressed/active and biologically relevant. Careful consideration needs to be given to the appropriate cellular context when designing assays for psychiatric genetics: key variant-interacting TFs and RBPs expressed in neurons may not be present in convenient cell lines (e.g., K562), potentially rendering functional neural elements/variants apparently silent.

Despite their potential, few MPRA have examined disease-associated variants while considering both cell type and -omic predictions. Tewhey, et. al<sup>126</sup> screened 30,000 eQTL SNPs from human lymphoblastoid cell lines (LCLs) using MPRA in LCLs, maintaining the discovery context in their assay. Over 3,400 active regulatory sequences were identified, including 850 activity-modulating variants (24%), consistent with functional (expression-modulating) SNP associations tagging linked, non-functional SNPs, akin to GWAS. Illustrating MPRA's sensitivity, significant allelic differences in activity were detectable at effect sizes <2-fold. Emerging work by Choi, et. al<sup>153</sup> prioritized over 800 SNPs—guided by fine-mapping and epigenomics—from 16 melanoma GWAS loci, to assay for transcriptional-regulatory activity in cultured melanocytes. Candidate variants with concordant eQTL signal in independent melanocyte data were further investigated, ultimately enabling experimental demonstration of biophysical (TF binding), molecular (target gene expression), cellular (growth rate), and *in vivo* (melanoma rate in transgenic zebrafish) variant

mechanisms. Finally, a recent MPRA of autoimmune GWAS loci yielded replicable findings across 12 donor lines of CD4+ T-cells, which were discordant with the more easily accessible—but leukemic—Jurkat cell line<sup>154</sup>. These experiments exquisitely illustrate MPRA’s capacity for sensitivity, context specificity, and high discovery rates, especially when integrating both association data and multi-omic annotations.

As with functional element assays, functional variant assays have recently moved *in vivo*, again including the brain. Kvon, et. al<sup>155</sup> utilized a novel knock-in system to generate transgenic mice expressing a LacZ reporter expressed under putatively regulatory elements containing rare, polydactyly-associated variants; subsequent LacZ staining clarified which variants were functional based on alterations of limb bud LacZ patterns. Excitingly, a small-scale MPRA has recently emerged using *in vivo* tissue: after prioritizing a single SNP from a bipolar disorder GWAS locus using epigenomic annotations, the two alleles of this sequence region were paired to 20 barcodes each and electroporated into embryonic mouse brains to confirm variant function<sup>125</sup>.

### **1.3.4 MPRA: limitations and design considerations**

With the powerful opportunities of MPRA come limitations. A major caveat lies in gathering candidate variants to assay. For example, a prominent and functionally characterized schizophrenia GWAS locus in the major histocompatibility complex (MHC) region<sup>82</sup>—containing hundreds of linked SNPs—turns out to mark heritable copy number variations in the complement C4a gene<sup>156</sup>; assaying only SNPs from this locus would not reveal the primary causal variant. Likewise, an MDD-associated SNP tags the absence of a transposon with regulatory effects on a noncoding RNA<sup>157</sup>. In other words, MPRA’s usefulness is contingent on investigation (and size—see below) of sequences to be assayed.



Further considerations include appropriateness of biological ‘contexts’ (**Figure 1.2**). At the level of ‘sequence context,’ MPRAs generally use multiplex oligonucleotide synthesis to custom-design sequences and variations by the thousands. However, such approaches are size-limited to ~300bp, which precludes assay of large or spaced regulatory sequences. Oligonucleotide synthesis also is error-prone; tagging each element with multiple barcodes safeguards against error-driven false-positives. Bulk capture-and-clone strategies circumvent these issues by utilizing larger, genomic DNA fragments directly<sup>124,158–160</sup> at the expense of precision assay design. Finally, element positioning can substantially influence results and replicability. While STARR-seq is favorable for one-step cloning (putative enhancers doubling as 3’UTR barcodes), emerging works illustrate that enhancer-like sequence placement in 3’UTRs yields results which cluster separately from other MPRA designs testing the same sequences<sup>161</sup>, and that such sequence placement can exert RNA stability effects that, without correction, may confound interpretation<sup>162</sup>.

Reporter gene features are also important in regulatory assays. Previous enhancer MPRA have demonstrated replicability by testing the same elements with a second promoter, with element activities highly correlated between the two<sup>122,163,164</sup>. However, these cross-promoter correlations (Pearson  $r$  0.7-0.8) have been weaker than often reported for biological replicates in MPRA ( $r > 0.9$ ). Promoter choice thus can influence assay results, via, for example, absence—or species differences in—promoter elements a *cis*-regulator requires. Likewise, UTR regulatory elements may be sensitive to the stoichiometry of transcripts and RBPs or miRNAs in the cell; excess transcript production by a strong promoter could potentially render effects of interacting regulators undetectable. In brief, rigorous MPRA or follow-up assays should use both a minimal promoter

and either a strong exogenous (*e.g.*, CMV) or a genomic promoter from the pertinent cell/tissue type (*e.g.*, a constitutively expressed, neuron-specific gene).

Importantly, the ability to test candidate sequences in their endogenous locus is not a feature of MPRA. Thus, ‘genomic context’—that is, episomal (AAV, plasmid) vs. genome-integrated (lentiviral) approaches—require consideration. Emerging comparisons find these approaches correlate well<sup>161</sup>, though certain applications may require a specific approach (*e.g.* MPRA of chromatin conformation<sup>165</sup>). The comparative throughput for a fixed number of cells is greater for plasmid transfection—thousands of plasmids per neural cell *in vitro*<sup>166</sup> compared to viral transduction ( $\leq$  tens of sequences/cell). These limitations and alternative methods are further considered in **Table 1.1**.

Other considerations include determining an appropriate ‘cellular’ and ‘organismal’ context (**Figure 1.2B-C**). Common strategies for choosing cellular contexts include using pathology (*e.g.*, substantia nigra in Parkinson’s disease), expression patterns of disease-associated genes (*e.g.*, cortical excitatory neurons in SCZ<sup>167</sup>), or GWAS-eQTL overlaps (*e.g.*, neurogenic niches of mid-fetal brain in ASD and SCZ<sup>168</sup>). (Cell type prioritization was further covered elsewhere in the Special Issue<sup>169</sup> of *Biological Psychiatry* wherein this section was previously published.)

A notable opportunity is utilization of MPRA in human iPSC-derived neural cell types, which offer the ability to conduct cell type-specific assays in a human genetic context. Very recent<sup>145</sup> and emerging<sup>125,170</sup> MPRA are proof-of-principle for this approach, supporting advancement to assaying *variants* in the setting of iPSC derivatives. Moreover, while cell type-specific MPRA have

been restricted to *in vitro* settings, where reproducing tissue physiology (*e.g.*, inter-cell type interactions, hormones, stress) is difficult, barcoded multiplex AAV regulatory assays<sup>142</sup> indicate *in vivo*, cell type-specific MRPA is possible. Nonetheless, negative MPRA results should be interpreted cautiously; absence of function in one context may not extend to *all* contexts.

Statistical considerations in MPRA include appropriate library size (number of elements and paired barcodes) for the cell type to be tested. Generally, library size should be downsized for rarer, hard-to-maintain, or hard-to-transfect/transduce cell types to ensure robust barcode recovery and measurement. MPRA has tested  $\sim 10^7$  sequences simultaneously in easily transfected cancer cell lines<sup>159,160</sup>, though in physiological cell types, like NPCs, this capacity is  $10^4$ - $10^5$ <sup>115,144,171</sup>, with emerging work approaching  $10^6$ <sup>170,172</sup>. Library size is further constrained by element-per-cell (*i.e.*, lentiviral) approaches. In other words, the fidelity of the model system and the MPRA library size it can support are generally anticorrelated. A consensus on the depth of barcodes-per-element is, to date, absent, ranging from 1 (STARR-seq<sup>129</sup>) to several hundred in previous<sup>145</sup> and emerging<sup>172,173</sup> work, with highly correlated replicate measurements across this range. Tewhey, et al estimated that statistical benefits for small-effect transcriptional-regulatory variants accrue by 5 barcodes, and asymptote around 25-50<sup>126</sup>; another finds  $> 10$  barcodes consistently yield inter-replicate  $r > 0.97$  in several cell types<sup>161,172</sup>. Whether these guidelines apply to UTR assays remains unclear. Overall, MPRA power guidelines would benefit substantially from deep assessment by modelers and statisticians.

Finally, given a finite number of elements that can be simultaneously assayed, one can choose whether to prioritize candidate variants using epigenomic data, or simply include all linked SNPs

(**Supplementary Figure 1.2**). An assay's 'hit rate' may be improved by prioritizing variants with *indirect* evidence of function, with the caveats of relying on epigenomic data discussed previously. However, foregoing such prioritization enables analysis of how well such features actually predict measured expression. Thus, the decision of prioritization must balance the value of 'hits' vs. identifying functionally predictive indirect measures (epigenetics) for the target cell type or disease.

### **1.3.5 Complementary methods in high-throughput study of DNA and RNA regulatory elements**

There are a variety of other approaches that complement MPRA (**Table 1.1**). Of course, lower throughput enhancer assays allow screening of the same elements or variants across a variety of contexts, even *in vivo*. Whether conducted using AAV (*e.g.*,<sup>174</sup>), or transgenesis (*e.g.*,<sup>175</sup>), these should remain gold standard approaches for validation and deep characterization of small numbers of elements and variants, including those identified by MPRA.

A primary limitation of MPRA is the inability to test regulators in their endogenous genomic position and sequence context. Sequence-specific targeting using CRISPR/Cas9 has enabled several additional techniques for probing molecular and cellular effects of regulatory variation, with the caveat that, unlike MPRA, these techniques do not currently allow for the multiplexed study of post-transcriptional/translational regulatory variants. Nonetheless, these techniques enable study of putative disease gene roles in gene expression networks and cellular phenotypes. *Perturb-seq*<sup>176</sup> combines genewise perturbation by CRISPR with single-cell RNA-seq to identify gene sets dysregulated by loss of function of each candidate gene. These have, for example, been used to discover co-transcribed gene networks involved in neuronal remodeling<sup>177</sup> and for *in vivo*

assessment of genes harboring *de novo* loss of function mutations in ASDs<sup>178</sup>. Likewise, CRISPR screens can be used to define functional elements influencing selectable traits (*e.g.*, proliferation), as in an emerging study perturbing both genes and *cis*-regulatory elements to define their roles in human neural stem cell proliferation<sup>179</sup>. Finally, CRISPR editing has been used *in vitro* to assess single-transcript noncoding variant effects by comparing allelic RNA and genomic DNA abundances in edited cultures<sup>180</sup>, a potential means of single-variant validation/follow-up of UTR MPRA findings. To my knowledge, such assays have not been conducted at a genome-wide scale in psychiatric disease, but have been used to identify genes that alter expression of the Parkinson's-associated *PARKIN*<sup>181</sup>.

*Cis*-regulatory MPRA cannot identify the endogenous target gene(s) of functional elements. Fortunately, CRISPR-derived methods using a mutagenically-‘dead’ Cas9 (dCas9) conjugated to a transcriptional activator or repressor allow targeted potentiation or repression of endogenous genomic regulatory elements (CRISPRa and CRISPRi, respectively) to assess altered gene expression and other outcomes. These technologies are already online in state-of-the-art human neuroscience models: a recent CRISPRi study knocked down over 2000 genes by targeting their promoters in iPSC-derived excitatory neurons, defining their context-specific roles in their survival, differentiation, and proliferation—including gene effects altered by co-culture with astrocytes<sup>182</sup>. Emerging work has further leveraged CRISPRi’s cell type specificity to study ASD-associated gene knockdown effects in an etiologically relevant cellular context (NPCs)<sup>183</sup>. A recently introduced extension of CRISPRi (‘CRISPRi-FlowFISH’) targets intergenic regulators, identifying their target gene by concurrent fluorescent *in situ* hybridization against genes from the same chromodomain. Fluorescence-intensity sorting into bins and subsequent RNA-seq can then

pair regulators (via guide RNA sequence) and target genes (altered FISH signal in a guide RNA's presence)<sup>184</sup>. While this assay was performed in K562 cells, it is not hard to envision its extension to neural cell types *in vitro* or *in vivo*. Altogether, CRISPR-based follow-up of MPRA candidates to define target genes and verify of genomic activity of regulators/variants will be key to developing insights in psychiatric genomics.

**Table 1.1. Strengths and Limitations of Functional-Regulatory Assays in Terms of Throughput and Sequence and Cellular Contexts.** *Method family:* An umbrella term covering multiple adaptations of an assay. *Technique:* The particular adaptation of the family's assay. "CRISPR editing" signifies precision replacement of an endogenous genomic sequence with a desired sequence (as opposed to CRISPR *mutagenesis*). *Table begins on next page.*

Method Family	Method	Can assay variants (e.g. SNPs) for function?	Can assay elements (e.g., TFBS, enhancers) for function?	Largest sequence / target	Simultaneous throughput for variants / perturbations per sample	Can assay cellular phenotype?	Genome-integrated?	Assays at the endogenous genomic sequence?	Demonstrated in model organism CNS <i>in vivo</i> ?	Demonstrated in human primary or iPSC-derived neural stem cells, NPCs or neurons?	Can identify target gene of endogenous <i>cis</i> -regulator?
MPRA	Multiplex (Barcoded) AAV Transcription Regulatory Assays	Yes	Yes	3-5 kb	100s-1000s	No	No	No	Yes <sup>142</sup>	No	No
	Enhancer MPRA and STARR-seq			~150-200 (custom oligos); ~700 (capture-and-clone)	10,000-10 <sup>6</sup>				Yes <sup>125</sup>	Yes <sup>125,170,172</sup>	No
	3'UTR MPRA/PTR E-seq			Not demonstrated but see above	No				N/A		
	5'UTR MPRA			"	No				N/A		
	RNA Splicing MPRA			"	No				N/A		
	Protein Translation MPRA			"	No				N/A		
CRISPR Regulatory Disruption Assays	CRISPRi	No	Yes	~50 bp	•Max demonstrated in CNS <i>in vivo</i> : 5 targets (2 sgRNAs each)	Yes	Yes	Yes	Yes <sup>185</sup>	Yes <sup>182,186,187</sup>	Yes



					•Max demonstrated in neural cell types <i>in vitro</i> : ~12,000						
	CRISPRa				•Max demonstrated in CNS <i>in vivo</i> : 10 targets (5 sgRNAs each) •Max demonstrated in neural cell types <i>in vitro</i> : 3	Yes	Yes		Yes <sup>188-190</sup>	Yes <sup>187,191,192</sup>	Yes
	CRISPR Mutagenesis of regulatory elements				•Max demonstrated in neural cell types <i>in vitro</i> : 26,000 targets (2 sgRNAs per target)	Yes	Yes		No	Yes <sup>179</sup>	For an <i>a priori</i> defined gene <sup>193,194</sup>
	CRISPRi-FlowFISH				~900	Not demonstrated	Yes		No	No	Yes
Low/single throughput	CRISPR Editing	Yes	Yes	Several kb	1-2	Yes	Yes	Yes	Yes	Yes	Yes
	AAV Transcription Regulatory Assays with traditional readouts			3-5 kb	1	No	No	No	Yes (esp. tacitly via cell-type targeted optogenetic, chemogenetic, and circuit-labeling)	Yes <sup>174</sup>	No

	(fluorescence , LacZ, etc.)								techniques, as in <sup>174</sup> )		
	Luciferase Reporter Assay			3-5 kb	1	No	No	No	Yes <sup>195</sup>	Yes	No

### 1.3.6 The utility of MPRA for parsing linked variation

One minimally-explored challenge in parsing loci implicated in common variant association studies is that multiple variants in the region are near-equally statistically associated, and thus near-equally plausible functional mediators of risk. Indeed, in statistical genetics, repeating a GWAS or eQTL association analysis after conditioning on the lead SNP within a block often reveals one or more additional independent variants associated with traits or gene expression, respectively (e.g.,<sup>196-198</sup>); in fact, nearly half (~8,000) of brain expressed genes have  $\geq 1$  conditional eQTL<sup>199</sup>. An exemplary reporter assays systematically evaluated such linked sets of variants, using epigenomic data to identify 16 enhancers near the *RET* gene, known to be downregulated in Hirschsprung's disease (failure of terminal colonic nerves to form *in utero*). These putative regulators were assayed for allele-differential activity in a model of the disease-relevant cell type, neural crest cells. These were likewise assayed with transient expression in mouse embryos driving a LacZ reporter to identify the crest-relevant regulatory sequences. They also validated roles for regulator-binding TFs via siRNA knockdown. In all, this identified three functional SNPs in linkage disequilibrium with synergistic effects on Hirschsprung-like deficits in colonic nerve development<sup>200</sup>. It is easy to imagine from this small-scale example how MPRA could be used to simultaneously dissect several linked blocks of disease-associated common variants.

While the *RET* SNPs were several hundred kilobases apart, another challenge is with variants in very close proximity, potentially within the same regulatory element. While oligonucleotide synthesis is limited in length, MPRA-based assay of such variants spanning up to 700bp is now possible with use of PCR amplification of oligos with uniquely complimentary ends<sup>161</sup>. Likewise, 'capture-and-clone' MPRA designs, often using STARR-seq architecture to simplify cloning, fragment genomic DNA or capture e.g. ChIP-seq DNA fragments to test larger fragments than

attainable via oligonucleotide synthesis<sup>124,159,160,201</sup>. Altogether, MPRAs allow for discovery of multiple functional variants per linkage region, as well as close-proximity discovery and disentanglement of multi-variant regulatory effects. Such efforts at a consortium scale could likewise characterize functionality and properties of linked untranscribed and untranslated variants across myriad cell types.

### **1.3.7 MPRA as an avenue to dissect multiallelic and polygenic mechanisms of neuropsychiatric traits**

While MPRA cannot intrinsically scale up to functional demonstration of cell-, tissue-, or behavior-level phenotypes, they have the potential to provide key information to guide molecular hypotheses for how these higher-order phenotypes emerge from large sets of regulators and/or their target genes. I focus here on examination of variants across disease loci—that is, defining shared and recurrent features among MPRA-nominated functional variants across the genome that may collectively underlie large portions of polygenic disease risk.

The most vexing question that remains after individual functional variant mechanisms are elucidated is how variants *collectively* contribute to phenotypic risk. MPRA provides several ways to begin addressing this question: 1) identifying regulatory features shared by across several functional risk variants; 2) identifying functional modules enriched for variant-impacted genes; 3) providing functional annotations to variants for computational genomic approaches; and finally, 4) by enabling study of variant-by-environment interactions contrasting MPRA across conditions.

Firstly, MPRA experiments running the gamut from basic regulatory genomics to human trait-associated variation have defined ‘regulatory grammars’ of assayed contexts. Identification of

functional variants in the MPRA setting enables similar establishment of the ‘regulatory grammar’ of a trait or disease. Functional variants identified by MPRA across several UTRs may feature a specific RBP’s binding site, for example, or could be used to deliberately define functional activity of a disease-associated miRNA, like miR-137<sup>82</sup>. Likewise, variants associated with a trait could be more likely to fall in particular TF binding sites or be enriched in cell type-specific marks of genomic regulation. Evidence of this convergence is seen in *de novo* variants associated with ASD: several distinct variants disrupt binding sites for a single TF, *NFIX*<sup>202</sup>. Similarly, putative gene targets of schizophrenia-associated variants are also putative—biases aside<sup>203</sup>—*Fmrp* targets<sup>204</sup>. MPRA has also identified such regulatory convergence by, for example, intersecting identified functional SNPs with TF ChIP-seq datasets in pertinent cell types to discover recurrently disrupted TF binding sites<sup>173</sup>. Assays of downstream consequences of variation also confirm biological convergence across association loci. A four-element-target CRISPRi/a assay revealed that schizophrenia risk genes act synergistically via shared influence on synaptic activity, and concurrent alteration of expression of all four genes results in a cellular transcriptome more accurately reflective of postmortem schizophrenia brain tissue<sup>191</sup>. For both rare and common variants, identifying common regulators among risk genes provides information which can refine predictions of disease-related cell types based on TF, RBP, or epigenomic mark expression.

Secondly, genes and gene networks affected by statistically associated variation are often predicted using MAGMA<sup>205</sup>, which in essence scores genes based on proximity to an associated variant and its linkage partners. Resulting gene sets are subjected to analyses such as Gene Ontology enrichment or are examined for enrichment in WGCNA (coexpression) networks from candidate tissue types to identify pathways and mechanisms on which these genes converge. While its use is

ubiquitous in genomic studies, standard MAGMA gene association statistics for psychiatric disorders only modestly correlate to those from a tissue-specific, chromatin configuration-aware modification of MAGMA<sup>206</sup>, suggesting that biological hypotheses from MAGMA gene sets may miss disease-associated genes in brain. Being able to refine implicated genes by functional validation using—or in follow-up to—MPRA will help to benchmark such approaches and refine prediction convergence with ‘truly’ dysregulated candidate genes.

Thirdly, epigenomic data alone is not comprehensively predictive of active regulators. However, well-informed analyses of human genetic findings rely heavily on such annotations to convert associations into biological hypotheses. Critically, these epigenomic data—unlike MPRA data—can be collected from postmortem human tissue. MPRA focused on neuropsychiatric disorder associated variation stand to benefit *from* high-information datasets by aiding variant prioritization for assay inclusion. Several recent datasets on synthetic UTRs<sup>117,121</sup>, RNA binding proteins<sup>207,208</sup>, and postmortem human brain multi-omics<sup>135,197,209–217</sup> are worth noting for readers investigating disease-associated variation. Integrative computational analyses have brought these datasets together predict functional variation in SCZ, bipolar disorder, and ASDs<sup>218,219</sup>. These constitute high-priority candidates for experimental validation by MPRA. Furthermore, emerging work reveals a symbiotic relationship developing between epigenomics and functional assays: functional element/variant information from MPRA has been used alongside epigenomic annotations to improve machine learning predictions of functional variants<sup>220</sup>. Predictions from these refined algorithms are another low-hanging fruit for candidates to assay by MPRA; those results could then constitute new training data. Such refinement of epigenomic data interpretation coupled with functionally-demonstrated regulatory variation would mutually benefit one another

and myriad downstream analyses, such as variant enrichment in genomic features and disease gene identification. For example, TWAS<sup>221</sup> and Predixcan<sup>222</sup> intersect gene expression QTLs (eQTL) with trait-associated variants to predict expression differences between cases and controls, thus identifying dysregulated gene sets. MPRA data can disentangle which eQTL SNPs are truly functional from those associated only due to LD, which could thus refine variant-gene pairings used in these analyses. Altogether, MPRA can serve to refine both epigenomic and genic definitions of truly causal disease features.

Finally, the context-specificity of MPRA represents a newfound ability to assess variant effects on gene regulation *en masse* under different biological and environmental contexts, including with *in vivo* models. While issues of convergent disease effects across genes and regulators are indeed complex, environmental effects—perhaps most canonically, stress—on these regulators are questions at the forefront of understanding polygenic risk in neuropsychiatric disorders. Pharmacologic variables have been successfully tested in MPRA, namely in the identification of glucocorticoid-responsive<sup>148</sup> and p53-responsive<sup>201</sup> regulatory elements. MPRA could further be layered with concurrent gene perturbations (*e.g.*, knockdown of a putative regulator), or cell culture conditions for *in vitro* identification of variant-environment interactions, exemplified by MPRA identification of neuronal activity-dependent enhancers<sup>144</sup>. As mouse and human brain cell types and their gene expression patterns are largely (though not entirely) conserved both in development<sup>223</sup> and adulthood<sup>224</sup>, the extension of MPRA to the *in vivo* context will enable study of broader endogenous and exogenous disease-associated factors, such as sex or stress. Identifying variants with environment-dependent functions would be a start toward identifying convergent molecular mechanisms behind conditional disease risk in disorders such as MDD.

Overall, MPRA presents unique opportunities to dissect polygenicity of psychiatric disorders via simultaneous identification of functional variants across identified risk space. Beyond the primary benefits of identifying ‘true positive’ functional variants in specific biological and environmental contexts, MPRA stands to rapidly broaden, deepen, and refine hypotheses and mechanisms of both noncoding disease risk and of gene-regulatory architecture itself.

## **1.4 Approach**

To characterize sex differences in MDD-relevant cell types and in genetic risk, I leveraged multiple techniques, including MPRA as described above. Chapter 2 of this dissertation describes my work interrogating whether sex differences in gene expression and/or regulation are present in the noradrenergic locus coeruleus (LC) by using Translating Ribosome Affinity Purification (TRAP) to measure sex-differential expression, followed by pharmacology and behavior experiments demonstrating functional consequences of the observed sex differences. Pattern searching in the mouse genome near sex-differentially expressed genes of the LC revealed enrichment of several putative regulatory sequences. Chapter 3 then briefly outlines a project I designed to assay these features by MPRA to functionally demonstrate sex-differential regulatory activity of small regions of mouse genome containing them, and the unfortunate demise of these experiments before they began. In Chapter 4, I begin using MPRA to identify functional variation from MDD-associated loci in mouse neuroblastoma cells *in vitro*. Enrichment analysis of sequence features disrupted by the functional variants identified recurrent roles across the loci for retinoic acid-responsive transcription factors. As neuroblastoma cells are robustly retinoid sensitive, this system was readymade for a subsequent drug-variant interaction MPRA, identifying additional retinoid-



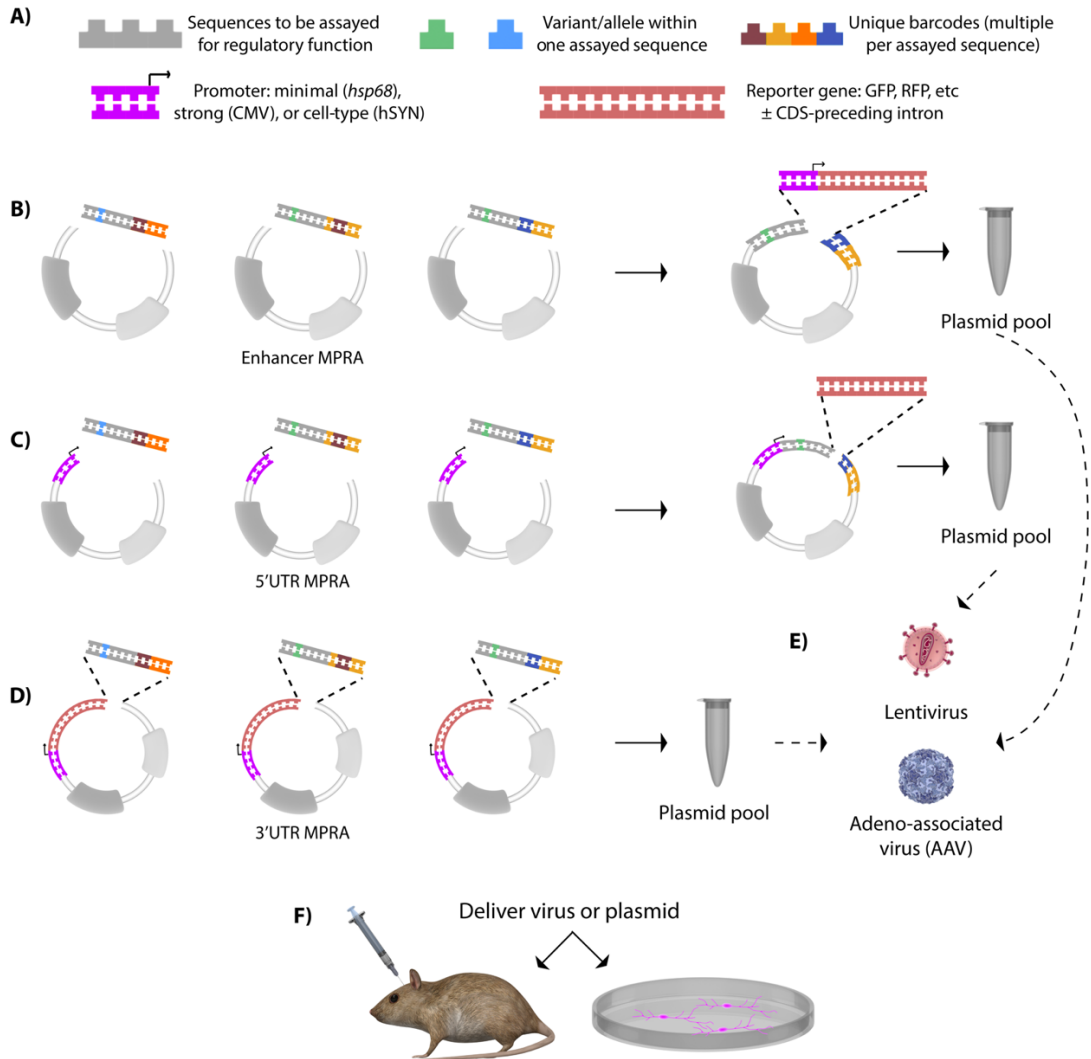
mediated or -altered variants from across MDD loci. Finally, in Chapter 5, I describe the results of successfully adapting highly-multiplexed MPRA to the mouse brain, enabling me to measure variant function in a high-fidelity model of both the brain as an organ and its organismal environment dictated by sex and age. I further complex this with TRAP to specifically identify functional and sex-interacting variation within hippocampal excitatory neurons, honing the resolution of the assay to a specific candidate cell type underlying molecular perturbations in MDD. Identification of shared regulatory features for functional and sex-interacting variation from these first *in vivo* MPRA indicated enrichment of transcription factors that co-operate with sex hormone receptors. To verify this, I then performed additional MPRA in the neonatal whole mouse brain at day 0 and day 10 of life in order to utilize a developmentally distinct, but naturally occurring, period of sex hormone release (perinatal) and quiescence (day 10, juvenile age). These experiments confirmed that sex-interacting variation is absent solely at day 10, supporting an “activational” (ligand-dependent) role for sex hormones in sex-differential activity of MDD risk variants.

## **1.5 Acknowledgements for Sections 1.3 and 1.6**

This work was funded by NIMH grants 1F30MH1116654 and 1R01MH116999, and Simons Foundation 571009. Previously published figures were reproduced in Figure 1.2, panels E and F, with permission respectively from Oxford University Press (Rightslink license # 4842680549544; generally licensed under Creative Commons 4.0 BY-NC (<https://creativecommons.org/licenses/by-nc/4.0/>)), and *a priori* permission for review article use from Proceedings of the National Academy of Sciences. I would like to thank Sergej Djuranovic, Ph.D., Barak Cohen, Ph.D., and Cohen lab alumni Dana King, Ph.D., and Brett Maricque, Ph.D.

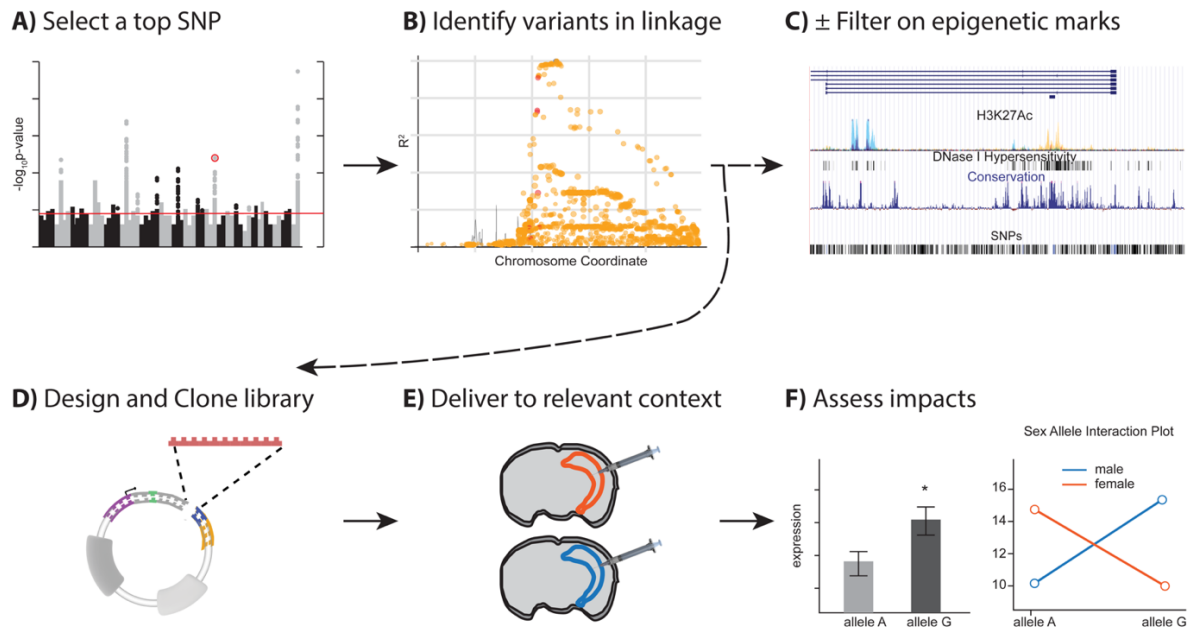
for their collaboration and guidance designing, adapting, and analyzing MPRAAs applied to neuropsychiatry. I would also like to thank Idoya Lahortiga, Ph.D. and Luk Cox, Ph.D. curators of Somersault1824 (<https://www.Somersault1824.com>), for their open-access, Creative Commons BY-NC-SA 4.0 licensed libraries of high-quality biomedicine graphics (especially those from Graphite Life Explorer, ePMV, and Eyewire), adapted for figures in this review. Finally, I would like to thank Lexi Harris for assistance with figure design, and Mike Vasek and Tony Fischer for their assistance in editing the text.

## 1.6 Supplementary Material



**Supplementary Figure 1.1. Common Designs and Library Delivery in MPRA.** **A)** Visual key for subsequent panels. **B)** Enhancer MPRA most commonly clone a pool of custom oligonucleotide pools containing sequences to be assayed, each paired to multiple unique barcode sequences, into a vector. Subsequently, a reporter gene driven by a minimal or core promoter—alone, driving only modest transcription—is added downstream of the element of interest, placing the barcode in a 3'UTR. In the case of STARR-seq (not shown), the paradigm in panel **D)** is instead used, with the *cis*-regulatory element being transcribed and thus acting as its own barcode. **C)** 5'UTR assays likewise use a two-step cloning assay, placing elements immediately *downstream* of a promoter. A reporter gene itself is inserted between the element and barcode. **D)** 3'UTR assays place the elements of interest immediately adjacent to barcodes, all downstream of a promoter-reporter in a single cloning step. **E)** A sequence pool from panel **B)**, **C)**, or **D)** can then be packaged into AAV or lentivirus (if in a compatible vector), or used for assay as plasmid directly. It is important to note that in all three of these scenarios, the library is constructed as a pool of sequences, which can result in variable levels of DNA for each element and barcode, hence the

normalization of RNA counts to DNA counts in downstream analysis (**Figure 1.1**). **F**) The plasmid or viral MPRA library can be delivered to cells in culture or *in vivo*.



**Supplementary Figure 1.2. Example of a Hypothetical MPRA.** **A)** Starting from common variant GWAS for a psychiatric disease (e.g. MDD), loci showing statistical association are selected for study. **B)** Within each locus, potential regulatory variants are identified based on showing sufficient linkage with the lead SNPs that they may be causal. Consideration should be given to whether such variants should be studied as candidate transcriptional or post-transcriptional regulatory elements. (Linkage plot generated using LDLink<sup>225</sup>. **C)** Variants can be prioritized based on epigenetic data from appropriate tissues or cell types, *or* all variants can be utilized (if such data are not available, or if one wants to test the predictive power of such data) (Data shown is from emerging neuronal chromatin contact data<sup>137</sup> using the WashU Epigenome Browser<sup>226</sup>. **D)** Elements with variants are cloned and prepared as described in **Figure 1.1**. **E)** MPRA library is delivered into the appropriate context for the disorder. In this hypothetical example, the library was packaged in AAV (to allow delivery to adult neurons *in vivo*) of the hippocampus (associated with MDD by imaging studies<sup>15</sup>, and powered to look for sex differences (since MDD has higher prevalence in females<sup>5</sup>. **F)** RNA is recovered and data is analyzed to define impact of variants on the expression of each element as described in **Figure 1.2**. Hypothetical results might discover a significant main effect of an allele (left panel) and/or sex specific interactions (right panel).

## 1.7 References

1. Kessler, R. C. *et al.* Lifetime Prevalence and Age-of-Onset Distributions of DSM-IV Disorders in the National Comorbidity Survey Replication. *Arch Gen Psychiat* 62, 593 (2005).
2. Hasin, D. S., Goodwin, R. D., Stinson, F. S. & Grant, B. F. Epidemiology of Major Depressive Disorder. *Arch Gen Psychiat* 62, 1097 (2005).
3. Salk, R. H., Hyde, J. S. & Abramson, L. Y. Gender Differences in Depression in Representative National Samples: Meta-Analyses of Diagnoses and Symptoms. *Psychol Bull* 143, 783–822 (2017).
4. Weissman, M. M. *et al.* Cross-National Epidemiology of Major Depression and Bipolar Disorder. *Jama J Am Medical Assoc* 276, 293 (1996).
5. Brody, D. J., Pratt, L. A. & Hughes, J. P. Prevalence of Depression Among Adults Aged 20 and Over: United States, 2013-2016. *Nchs Data Brief* 1–8 (2018).
6. Collaborators, G. 2019 D. and I. *et al.* Global burden of 369 diseases and injuries in 204 countries and territories, 1990–2019: a systematic analysis for the Global Burden of Disease Study 2019. *Lancet Lond Engl* 396, 1204–1222 (2020).
7. Barefoot, J. C. *et al.* Depression and Long-Term Mortality Risk in Patients With Coronary Artery Disease\*\*This study was supported in part by grants P01 HL36587, R01 HL43028, R01 HL44998, R01 HL45702, and R01 HL49572 from the National Heart, Lung and Blood Institute; and AG-09276, AG09663, 5P60 AG-11268, and P02 AG-12058 from the National Institute on Aging. *Am J Cardiol* 78, 613–617 (1996).
8. Glassman, A. H. & Shapiro, P. A. Depression and the Course of Coronary Artery Disease. *Am J Psychiat* 155, 4–11 (1998).
9. Jones, D. J., Bromberger, J. T., Sutton-Tyrrell, K. & Matthews, K. A. Lifetime History of Depression and Carotid Atherosclerosis in Middle-aged Women. *Arch Gen Psychiat* 60, 153–160 (2003).
10. Garcia, R. G. *et al.* Impact of sex and depressed mood on the central regulation of cardiac autonomic function. *Neuropsychopharmacol* 45, 1280–1288 (2020).
11. Penninx, B. W. J. H. *et al.* Depression and Cardiac Mortality. *Arch Gen Psychiat* 58, 221 (2001).
12. Zimmerman, M., Ellison, W., Young, D., Chelminski, I. & Dalrymple, K. How many different ways do patients meet the diagnostic criteria for major depressive disorder? *Compr Psychiat* 56, 29–34 (2015).
13. Zisook, S. *et al.* Effect of Age at Onset on the Course of Major Depressive Disorder. *Am J Psychiat* 164, 1539–1546 (2007).

14. Kessler, R. C. *et al.* The Epidemiology of Major Depressive Disorder. *Jama* 289, 3095 (2003).
15. Schmaal, L. *et al.* Subcortical brain alterations in major depressive disorder: findings from the ENIGMA Major Depressive Disorder working group. *Mol Psychiatr* 21, 806–812 (2015).
16. Mitchell, B. L. *et al.* The Australian Genetics of Depression Study: new risk loci and dissecting heterogeneity between subtypes. *Biol Psychiatr* (2021) doi:10.1016/j.biopsych.2021.10.021.
17. Shukla, R. *et al.* Molecular characterization of depression trait and state. *Mol Psychiatr* 1–12 (2021) doi:10.1038/s41380-021-01347-z.
18. Marcus, S. M. *et al.* Gender differences in depression: Findings from the STAR\*D study. *J Affect Disorders* 87, 141–150 (2005).
19. Matza, L. S., Revicki, D. A., Davidson, J. R. & Stewart, J. W. Depression With Atypical Features in the National Comorbidity Survey. *Arch Gen Psychiatr* 60, 817 (2003).
20. Posternak, M. A. & Zimmerman, M. The prevalence of atypical features across mood, anxiety, and personality disorders. *Compr Psychiatr* 43, 253–262 (2002).
21. Benazzi, F. Prevalence and clinical features of atypical depression in depressed outpatients: a 467-case study. *Psychiat Res* 86, 259–265 (1999).
22. COLE, C. E. *et al.* A Controlled Study of Efficacy of Iproniazid in Treatment of Depression. *M Archives Gen Psychiatry* 1, 513–518 (1959).
23. SCHILDKRAUT, J. J. THE CATECHOLAMINE HYPOTHESIS OF AFFECTIVE DISORDERS: A REVIEW OF SUPPORTING EVIDENCE. *Am J Psychiatr* 122, 509–522 (1965).
24. BUNNEY, W. E. & DAVIS, J. M. Norepinephrine in Depressive Reactions: A Review. *Arch Gen Psychiatr* 13, 483 (1965).
25. LeGates, T. A., Kvarita, M. D. & Thompson, S. M. Sex differences in antidepressant efficacy. *Neuropsychopharmacol Official Publ Am Coll Neuropsychopharmacol* 44, 140–154 (2018).
26. Daly, E. J. *et al.* Efficacy of Esketamine Nasal Spray Plus Oral Antidepressant Treatment for Relapse Prevention in Patients With Treatment-Resistant Depression. *Jama Psychiatr* 76, 893–903 (2019).
27. MacKenzie, G. & Maguire, J. The role of ovarian hormone-derived neurosteroids on the regulation of GABAA receptors in affective disorders. *Psychopharmacology* 231, 3333–3342 (2014).
28. Edinoff, A. N. *et al.* Brexanolone, a GABAA Modulator, in the Treatment of Postpartum Depression in Adults: A Comprehensive Review. *Frontiers Psychiatry* 12, 699740 (2021).

29. Salminen, L. E. *et al.* Hippocampal subfield volumes are uniquely affected in PTSD and depression: International analysis of 31 cohorts from the PGC-ENIGMA PTSD Working Group. *Biorxiv* 739094 (2019) doi:10.1101/739094.
30. Hilland, E. *et al.* Exploring the links between specific depression symptoms and brain structure: A network study. *Psychiat Clin Neuros* 74, 220–221 (2020).
31. Sha, Z. & Banihashemi, L. Integrative omics analysis identifies differential biological pathways that are associated with regional grey matter volume changes in major depressive disorder. *Psychol Med* 1–12 (2020) doi:10.1017/s0033291720002676.
32. Suh, J. S. *et al.* Cortical thickness in major depressive disorder: A systematic review and meta-analysis. *Prog Neuro-psychopharmacology Biological Psychiatry* 88, 287–302 (2018).
33. Hu, X. *et al.* Sex-specific alterations of cortical morphometry in treatment-naïve patients with major depressive disorder. *Neuropsychopharmacol* 1–8 (2022) doi:10.1038/s41386-021-01252-7.
34. Keller, A. S., Ball, T. M. & Williams, L. M. Deep phenotyping of attention impairments and the ‘Inattention Biotype’ in Major Depressive Disorder. *Psychol Med* 1–10 (2019) doi:10.1017/s0033291719002290.
35. Labonté, B. *et al.* Sex-specific transcriptional signatures in human depression. *Nat Med* 23, 1102–1111 (2017).
36. Seney, M. L. *et al.* Opposite Molecular Signatures of Depression in Men and Women. *Biol Psychiat* 84, 18–27 (2018).
37. Brivio, E., Lopez, J. P. & Chen, A. Sex Differences: Transcriptional Signatures of Stress Exposure in Male and Female Brains. *Genes Brain Behav* e12643 (2020) doi:10.1111/gbb.12643.
38. Kwon, D. Y. *et al.* Neuronal Yin Yang1 in the prefrontal cortex regulates transcriptional and behavioral responses to chronic stress in mice. *Nat Commun* 13, 55 (2022).
39. Duman, R. S., Sanacora, G. & Krystal, J. H. Altered Connectivity in Depression: GABA and Glutamate Neurotransmitter Deficits and Reversal by Novel Treatments. *Neuron* 102, 75–90 (2019).
40. Maynard, K. R. *et al.* Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat Neurosci* 24, 425–436 (2021).
41. Ruzicka, W. B. *et al.* Single-cell dissection of schizophrenia reveals neurodevelopmental-synaptic axis and transcriptional resilience. *Biorxiv* (2020) doi:10.1101/2020.11.06.20225342.
42. Ayhan, F. *et al.* Resolving cellular and molecular diversity along the hippocampal anterior-to-posterior axis in humans. *Neuron* (2021) doi:10.1016/j.neuron.2021.05.003.

43. Boldrini, M. *et al.* Resilience Is Associated With Larger Dentate Gyrus, While Suicide Decedents With Major Depressive Disorder Have Fewer Granule Neurons. *Biol Psychiat* 85, 850–862 (2019).
44. Bryois, J. *et al.* Genetic identification of cell types underlying brain complex traits yields insights into the etiology of Parkinson’s disease. *Nat Genet* 1–12 (2020) doi:10.1038/s41588-020-0610-9.
45. Zhang, K. *et al.* A single-cell atlas of chromatin accessibility in the human genome. *Cell* (2021) doi:10.1016/j.cell.2021.10.024.
46. Cruceanu, C. *et al.* Cell-Type-Specific Impact of Glucocorticoid Receptor Activation on the Developing Brain: A Cerebral Organoid Study. *Am J Psychiatry* appiajp202121010095 (2021) doi:10.1176/appi.ajp.2021.21010095.
47. Wickens, M. M., Bangasser, D. A. & Briand, L. A. Sex Differences in Psychiatric Disease: A Focus on the Glutamate System. *Front Mol Neurosci* 11, 197 (2018).
48. Williams, E. S. *et al.* Androgen-Dependent Excitability of Mouse Ventral Hippocampal Afferents to Nucleus Accumbens Underlies Sex-Specific Susceptibility to Stress. *Biol Psychiat* 87, 492–501 (2019).
49. Zhang, S. *et al.* Sex Differences in the Neuroadaptations of Reward-related Circuits in Response to Subchronic Variable Stress. *Neuroscience* 376, 108–116 (2018).
50. Hodes, G. E. *et al.* Sex Differences in Nucleus Accumbens Transcriptome Profiles Associated with Susceptibility versus Resilience to Subchronic Variable Stress. *J Neurosci* 35, 16362–16376 (2015).
51. Paden, W. *et al.* Sex differences in adult mood and in stress-induced transcriptional coherence across mesocorticolimbic circuitry. *Transl Psychiat* 10, 59 (2020).
52. Lei, Y. *et al.* SIRT1 in forebrain excitatory neurons produces sexually dimorphic effects on depression-related behaviors and modulates neuronal excitability and synaptic transmission in the medial prefrontal cortex. *Mol Psychiatr* 56, 1–18 (2019).
53. Song, L. *et al.* Hippocampal PPAR $\alpha$  is a novel therapeutic target for depression and mediates the antidepressant actions of fluoxetine in mice. *Brit J Pharmacol* 175, 2968–2987 (2018).
54. Noh, K. *et al.* Negr1 controls adult hippocampal neurogenesis and affective behaviors. *Mol Psychiatr* 24, 1189–1205 (2019).
55. Dion-Albert, L. *et al.* Sex-specific blood-brain barrier alterations and vascular biomarkers underlie chronic stress responses in mice and human depression. *Biorxiv* 2021.04.23.441142 (2021) doi:10.1101/2021.04.23.441142.
56. Dudek, K. A. *et al.* Molecular adaptations of the blood-brain barrier promote stress resilience vs. depression. *P Natl Acad Sci Usa* 117, 3326–3336 (2020).



57. Aguilar-Valles, A. *et al.* Translational control of depression-like behavior via phosphorylation of eukaryotic translation initiation factor 4E. *Nat Commun* 9, 2459 (2018).
58. Wainberg, M. *et al.* Clinical laboratory tests and five-year incidence of major depressive disorder: a prospective cohort study of 433,890 participants from the UK Biobank. *Transl Psychiat* 11, 380 (2021).
59. Cambiasso, M. J. *et al.* Interaction of sex chromosome complement, gonadal hormones and neuronal steroid synthesis on the sexual differentiation of mammalian neurons. *J Neurogenet* 1–7 (2017) doi:10.1080/01677063.2017.1390572.
60. Georgiou, P. *et al.* Estradiol mediates stress-susceptibility in the male brain. *Biorxiv* 2022.01.09.475485 (2022) doi:10.1101/2022.01.09.475485.
61. Jaric, I., Rocks, D., Grealley, J. M., Suzuki, M. & Kundakovic, M. Chromatin organization in the female mouse brain fluctuates across the oestrous cycle. *Nat Commun* 10, 2851 (2019).
62. Frye, H. E. *et al.* Sex Differences in the Role of CNIH3 on Spatial Memory and Synaptic Plasticity. *Biol Psychiat* 90, 766–780 (2021).
63. Jaric, I., Rocks, D., Cham, H., Herchek, A. & Kundakovic, M. Sex and Estrous Cycle Effects on Anxiety- and Depression-Related Phenotypes in a Two-Hit Developmental Stress Model. *Front Mol Neurosci* 12, 74 (2019).
64. Gegenhuber, B., Wu, M. V., Bronstein, R. & Tollkuhn, J. Regulation of neural gene expression by estrogen receptor alpha. *Biorxiv* 2020.10.21.349290 (2020) doi:10.1101/2020.10.21.349290.
65. Guillamón, A., Blas, M. R. de & Segovia, S. Effects of sex steroids on the development of the locus coeruleus in the rat. *Dev Brain Res* 40, 306–310 (1988).
66. Luque, J. M., Blas, M. R. de, Segovia, S. & Guillamón, A. Sexual dimorphism of the dopamine- $\beta$ -hydroxylase-immunoreactive neurons in the rat locus ceruleus. *Dev Brain Res* 67, 211–215 (1992).
67. Serova, L., Rivkin, M., Nakashima, A. & Sabban, E. L. Estradiol Stimulates Gene Expression of Norepinephrine Biosynthetic Enzymes in Rat Locus coeruleus. *Neuroendocrinology* 75, 193–200 (2002).
68. Borodovitsyna, O. & Chandler, D. J. Age- and sex- dependent changes in locus coeruleus physiology and anxiety-like behavior in response to acute stress. *Biorxiv* 2020.11.10.377275 (2020) doi:10.1101/2020.11.10.377275.
69. Valentino, R. J. & Bangasser, D. A. Sex-biased cellular signaling: molecular basis for sex differences in neuropsychiatric diseases. *Dialogues Clin Neurosci* 18, 385–393 (2016).
70. Bangasser, D. A., Wiersielis, K. R. & Khantsis, S. Sex differences in the locus coeruleus-norepinephrine system and its regulation by stress. *Brain Res* 1641, 177–88 (2016).

71. McGuffin, P., Katz, R., Watkins, S. & Rutherford, J. A Hospital-Based Twin Register of the Heritability of DSM-IV Unipolar Depression. *Arch Gen Psychiat* 53, 129 (1996).
72. Lyons, M. J. *et al.* A Registry-Based Twin Study of Depression in Men. *Arch Gen Psychiat* 55, 468 (1998).
73. Sullivan, P. F., Neale, M. C. & Kendler, K. S. Genetic Epidemiology of Major Depression: Review and Meta-Analysis. *Am J Psychiat* 157, 1552–1562 (2000).
74. Consortium, C.-D. G. of the P. G. *et al.* Genetic relationship between five psychiatric disorders estimated from genome-wide SNPs. *Nat Genet* 45, 984–994 (2013).
75. Clements, C. C. *et al.* Genome-wide association study of patients with a severe major depressive episode treated with electroconvulsive therapy. *Mol Psychiatr* 1–11 (2021) doi:10.1038/s41380-020-00984-0.
76. Levey, D. F. *et al.* Bi-ancestral depression GWAS in the Million Veteran Program and meta-analysis in >1.2 million individuals highlight new therapeutic directions. *Nat Neurosci* 1–10 (2021) doi:10.1038/s41593-021-00860-2.
77. Howard, D. M. *et al.* Genome-wide meta-analysis of depression identifies 102 independent variants and highlights the importance of the prefrontal brain regions. *Nat Neurosci* 22, 343–352 (2019).
78. Wray, N. R. *et al.* Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nat Genet* 50, 668–681 (2018).
79. Li, X. *et al.* Common variants on 6q16.2, 12q24.31 and 16p13.3 are associated with major depressive disorder. *Neuropsychopharmacol* 43, 2146–2153 (2018).
80. Cai, N. *et al.* Sparse whole-genome sequencing identifies two loci for major depressive disorder. *Nature* 523, 588–591 (2015).
81. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nat Genet* 47, 1236–41 (2015).
82. Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 511, 421–7 (2014).
83. Purcell, S. M. *et al.* Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature* 460, 748–52 (2009).
84. Gaugler, T. *et al.* Most genetic risk for autism resides with common variation. *Nat Genet* 46, 881–5 (2014).
85. Pettersson, E. *et al.* Genetic influences on eight psychiatric disorders based on family data of 4 408 646 full and half-siblings, and genetic data of 333 748 cases and controls. *Psychol Med* 49, 1–8 (2018).

86. Ryan, N. M. *et al.* DNA sequence-level analyses reveal potential phenotypic modifiers in a large family with psychiatric disorders. *Mol Psychiatr* 23, 2254–2265 (2018).
87. Howrigan, D. P. *et al.* Exome sequencing in schizophrenia-affected parent-offspring trios reveals risk conferred by protein-coding de novo mutations. *Nat Neurosci* 1–9 (2020) doi:10.1038/s41593-019-0564-3.
88. Gusev, A. *et al.* Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *Am J Hum Genet* 95, 535–52 (2014).
89. Gamazon, E. R. *et al.* Using an atlas of gene regulation across 44 human tissues to inform complex disease- and trait-associated variation. *Nat Genet* 50, 956–967 (2018).
90. Schaub, M. A., Boyle, A. P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease associations with regulatory information in the human genome. *Genome Res* 22, 1748–1759 (2012).
91. Doolittle, W. F. Is junk DNA bunk? A critique of ENCODE. *P Natl Acad Sci Usa* 110, 5294–300 (2013).
92. Eddy, S. R. The ENCODE project: missteps overshadowing a success. *Curr Biology Cb* 23, R259-61 (2013).
93. Benton, M. L., Talipineni, S. C., Kostka, D. & Capra, J. A. Genome-wide enhancer annotations differ significantly in genomic distribution, evolution, and function. *Bmc Genomics* 20, 511 (2019).
94. Kwasnieski, J. C., Fiore, C., Chaudhari, H. G. & Cohen, B. A. High-throughput functional testing of ENCODE segmentation predictions. *Genome Res* 24, 1595–1602 (2014).
95. Brown, A. A. *et al.* Predicting causal variants affecting expression by using whole-genome sequencing and RNA-seq from multiple human tissues. *Nat Genet* 49, 1747–1751 (2017).
96. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* 518, 317–30 (2015).
97. An, J.-Y. *et al.* Genome-wide de novo risk score implicates promoter variation in autism spectrum disorder. *Science* 362, eaat6576 (2018).
98. Mao, F. *et al.* Post-transcriptionally impaired de novo mutations contribute to the genetic etiology of four neuropsychiatric disorders. *Biorxiv* 175844 (2019) doi:10.1101/175844.
99. Turner, T. N. & Eichler, E. E. The Role of De Novo Noncoding Regulatory Mutations in Neurodevelopmental Disorders. *Trends Neurosci* 42, 115–127 (2018).
100. Satterstrom, F. K. *et al.* Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism. *Cell* (2020) doi:10.1016/j.cell.2019.12.036.

101. Mayr, C. Regulation by 3'-Untranslated Regions. *Annu Rev Genet* 51, 171–194 (2017).
102. Price, A. J. *et al.* Characterizing the nuclear and cytoplasmic transcriptomes in developing and mature human cortex uncovers new insight into psychiatric disease gene regulation. *Genome Res* 30, 1–11 (2019).
103. Kinney, J. B., Murugan, A., Callan, C. G. & Cox, E. C. Using deep sequencing to characterize the biophysical mechanism of a transcriptional regulatory sequence. *P Natl Acad Sci Usa* 107, 9158–63 (2010).
104. Patwardhan, R. P. *et al.* High-resolution analysis of DNA regulatory elements by synthetic saturation mutagenesis. *Nat Biotechnol* 27, 1173–1175 (2009).
105. Baeza-Centurion, P., Miñana, B., Schmiedel, J. M., Valcárcel, J. & Lehner, B. Combinatorial Genetics Reveals a Scaling Law for the Effects of Mutations on Splicing. *Cell* 176, 549-563.e23 (2019).
106. Rosenberg, A. B., Patwardhan, R. P., Shendure, J. & Seelig, G. Learning the sequence determinants of alternative splicing from millions of random sequences. *Cell* 163, 698–711 (2015).
107. Wong, M. S., Kinney, J. B. & Krainer, A. R. Quantitative Activity Profile and Context Dependence of All Human 5' Splice Sites. *Mol Cell* 71, 1012-1026.e3 (2018).
108. Safra, M., Nir, R., Farouq, D., Slutzkin, I. V. & Schwartz, S. TRUB1 is the predominant pseudouridine synthase acting on mammalian mRNA via a predictable and conserved code. *Genome Res* 27, 393–406 (2017).
109. Matreyek, K. A. *et al.* Multiplex assessment of protein variant abundance by massively parallel sequencing. *Nat Genet* 50, 874–882 (2018).
110. Grossman, S. R. *et al.* Systematic dissection of genomic features determining transcription factor binding and enhancer function. *Proc National Acad Sci* 114, E1291–E1300 (2017).
111. Gisselbrecht, S. S. *et al.* Highly parallel assays of tissue-specific enhancers in whole *Drosophila* embryos. *Nat Methods* 10, 774–780 (2013).
112. King, D. M. *et al.* Synthetic and genomic regulatory elements reveal aspects of cis-regulatory grammar in Mouse Embryonic Stem Cells. *Elife* 9, e41279 (2020).
113. Fiore, C. & Cohen, B. A. Interactions between pluripotency factors specify cis-regulation in embryonic stem cells. *Genome Res* 26, 778–86 (2016).
114. Levo, M. *et al.* Systematic Investigation of Transcription Factor Activity in the Context of Chromatin Using Massively Parallel Binding and Expression Assays. *Mol Cell* 65, 604-617.e6 (2017).

115. Maricque, B. B., Dougherty, J. D. & Cohen, B. A. A genome-integrated massively parallel reporter assay reveals DNA sequence determinants of cis -regulatory activity in neural cells. *Nucleic Acids Res* 45, gkw942 (2016).
116. White, M. A., Myers, C. A., Corbo, J. C. & Cohen, B. A. Massively parallel in vivo enhancer assay reveals that highly local features determine the cis-regulatory function of ChIP-seq peaks. *Proc National Acad Sci* 110, 11952–11957 (2013).
117. Cottrell, K. A., Chaudhari, H. G., Cohen, B. A. & Djuranovic, S. PTRE-seq reveals mechanism and interactions of RNA binding proteins and miRNAs. *Nat Commun* 9, 301 (2018).
118. Litterman, A. J. *et al.* A massively parallel 3' UTR reporter assay reveals relationships between nucleotide content, sequence conservation, and mRNA destabilization. *Genome Res* 29, 896–906 (2019).
119. Rabani, M., Pieper, L., Chew, G.-L. & Schier, A. F. A Massively Parallel Reporter Assay of 3' UTR Sequences Identifies In Vivo Rules for mRNA Degradation. *Mol Cell* 70, 565 (2018).
120. Rieger, M. A., King, D. M., Cohen, B. A. & Dougherty, J. D. CLIP-Seq and massively parallel functional analysis of the CELF6 RNA binding protein reveals a role in destabilizing synaptic gene mRNAs through interaction with 3'UTR elements in vivo. *Biorxiv* 401604 (2018) doi:10.1101/401604.
121. Sample, P. J. *et al.* Human 5' UTR design and variant effect prediction from a massively parallel translation assay. *Nat Biotechnol* 37, 803–809 (2019).
122. Castaldi, P. J. *et al.* Identification of Functional Variants in the FAM13A Chronic Obstructive Pulmonary Disease Genome-Wide Association Study Locus by Massively Parallel Reporter Assays. *Am J Resp Crit Care* 199, 52–61 (2019).
123. Myint, L. *et al.* A screen of 1,049 schizophrenia and 30 Alzheimer's-associated variants for regulatory potential. *Am J Medical Genetics Part B Neuropsychiatric Genetics* 183, 61–73 (2020).
124. Shen, S. Q. *et al.* Massively parallel cis -regulatory analysis in the mammalian central nervous system. *Genome Res* 26, 238–255 (2016).
125. Shen, S. Q. *et al.* A candidate causal variant underlying both higher intelligence and increased risk of bipolar disorder. *Biorxiv* 580258 (2019) doi:10.1101/580258.
126. Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants using a Multiplexed Reporter Assay. *Cell* 165, 1519–1529 (2016).
127. Ulirsch, J. C. *et al.* Systematic Functional Dissection of Common Genetic Variation Affecting Red Blood Cell Traits. *Cell* 165, 1530–1545 (2016).
128. Klein, J. C. *et al.* Functional testing of thousands of osteoarthritis-associated variants for regulatory activity. *Nat Commun* 10, 2434 (2019).

129. Arnold, C. D. *et al.* Genome-Wide Quantitative Enhancer Activity Maps Identified by STARR-seq. *Science* 339, 1074–1077 (2013).
130. Mogno, I., Kwasnieski, J. C. & Cohen, B. A. Massively parallel synthetic promoter assays reveal the in vivo effects of binding site variants. *Genome Res* 23, 1908–1915 (2013).
131. Kwasnieski, J. C., Mogno, I., Myers, C. A., Corbo, J. C. & Cohen, B. A. Complex effects of nucleotide variants in a mammalian cis-regulatory element. *P Natl Acad Sci Usa* 109, 19498–503 (2012).
132. Myint, L., Avramopoulos, D. G., Goff, L. A. & Hansen, K. D. Linear models enable powerful differential activity analysis in massively parallel reporter assays. *Bmc Genomics* 20, 209 (2019).
133. (DGT), F. C. and the R. P. and C. *et al.* A promoter-level mammalian expression atlas. *Nature* 507, 462–70 (2014).
134. Aguet, F. *et al.* Genetic effects on gene expression across human tissues. *Nature* 550, 204–213 (2017).
135. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature* 489, 57–74 (2012).
136. Song, M. *et al.* Mapping cis-regulatory chromatin contacts in neural cells links neuropsychiatric disorder risk variants to target genes. *Nat Genet* 51, 1252–1262 (2019).
137. Song, M. *et al.* Cell-type-specific 3D epigenomes in the developing human cortex. *Nature* 1–6 (2020) doi:10.1038/s41586-020-2825-4.
138. Nott, A. *et al.* Brain cell type-specific enhancer-promoter interactome maps and disease-risk association. *Sci New York N Y* 366, 1134–1139 (2019).
139. Blankvoort, S., Witter, M. P., Noonan, J., Cotney, J. & Kentros, C. Marked Diversity of Unique Cortical Enhancers Enables Neuron-Specific Tools by Enhancer-Driven Gene Expression. *Curr Biol* 28, 2103-2114.e5 (2018).
140. Nair, R. R., Blankvoort, S., Lagartos, M. J. & Kentros, C. Enhancer-Driven Gene Expression (EDGE) Enables the Generation of Viral Vectors Specific to Neuronal Subtypes. *Isience* 23, 100888 (2020).
141. Graybuck, L. T. *et al.* Enhancer viruses and a transgenic platform for combinatorial cell subclass-specific labeling. *Biorxiv* 525014 (2020) doi:10.1101/525014.
142. Hrvatin, S. *et al.* A scalable platform for the development of cell-type-specific viral drivers. *Elife* 8, e48089 (2019).
143. Hughes, A. E. O., Myers, C. A. & Corbo, J. C. A massively parallel reporter assay reveals context-dependent activity of homeodomain binding sites in vivo. *Genome Res* 28, 1520–1531 (2018).

144. Nguyen, T. A. *et al.* High-throughput functional comparison of promoter and enhancer activities. *Genome Res* 26, 1023–33 (2016).
145. Inoue, F., Kreimer, A., Ashuach, T., Ahituv, N. & Yosef, N. Identification and Massively Parallel Characterization of Regulatory Elements Driving Neural Induction. *Cell Stem Cell* (2019) doi:10.1016/j.stem.2019.09.010.
146. Osterwalder, M. *et al.* Enhancer redundancy provides phenotypic robustness in mammalian development. *Nature* 554, 239–243 (2018).
147. McClymont, S. A. *et al.* Parkinson-Associated SNCA Enhancer Variants Revealed by Open Chromatin in Mouse Dopamine Neurons. *Am J Hum Genet* 103, 874–892 (2018).
148. Vockley, C. M. *et al.* Direct GR Binding Sites Potentiate Clusters of TF Binding across the Human Genome. *Cell* 166, 1269–1281.e19 (2016).
149. Liang, D. *et al.* Cell-type-specific effects of genetic variation on chromatin accessibility during human neuronal differentiation. *Nat Neurosci* 24, 941–953 (2021).
150. Kopp, N., Climer, S. & Dougherty, J. D. Moving from capstones toward cornerstones: successes and challenges in applying systems biology to identify mechanisms of autism spectrum disorders. *Frontiers Genetics* 6, 301 (2015).
151. Jaffe, A. E. *et al.* Profiling gene expression in the human dentate gyrus granule cell layer reveals insights into schizophrenia and its genetic risk. *Nat Neurosci* 23, 510–519 (2020).
152. Shigaki, D. *et al.* Integration of Multiple Epigenomic Marks Improves Prediction of Variant Impact in Saturation Mutagenesis Reporter Assay. *Hum Mutat* 40, 1280–1291 (2019).
153. Choi, J. *et al.* Massively parallel reporter assays of melanoma risk variants identify MX2 as a gene promoting melanoma. *Nat Commun* 11, 2718 (2020).
154. Bourges, C. *et al.* Resolving mechanisms of immune-mediated disease in primary CD4 T cells. *Embo Mol Med* 12, e12112 (2020).
155. Kvon, E. Z. *et al.* Comprehensive In Vivo Interrogation Reveals Phenotypic Impact of Human Enhancer Variants. *Cell* 180, 1262–1271.e15 (2020).
156. Consortium, S. W. G. of the P. G. *et al.* Schizophrenia risk from complex variation of complement component 4. *Nature* 530, 177–183 (2016).
157. Liu, W. *et al.* Identification of a functional human-unique 351-bp Alu insertion polymorphism associated with major depressive disorder in the 1p31.1 GWAS risk loci. *Neuropsychopharmacol Official Publ Am Coll Neuropsychopharmacol* 1–11 (2020) doi:10.1038/s41386-020-0659-2.
158. Vockley, C. M. *et al.* Massively parallel quantification of the regulatory effects of noncoding genetic variation in a human cohort. *Genome Res* 25, 1206–1214 (2015).

159. Wang, X. *et al.* High-resolution genome-wide functional dissection of transcriptional regulatory regions and nucleotides in human. *Nat Commun* 9, 5380 (2018).
160. Arensbergen, J. van *et al.* High-throughput identification of human SNPs affecting regulatory element activity. *Nat Genet* 51, 1160–1169 (2019).
161. Klein, J. C. *et al.* A systematic evaluation of the design and context dependencies of massively parallel reporter assays. *Nat Methods* 1–9 (2020) doi:10.1038/s41592-020-0965-y.
162. Lee, D. *et al.* STARRPeaker: Uniform processing and accurate identification of STARR-seq active regions. *Biorxiv* 694869 (2019) doi:10.1101/694869.
163. Liu, Y. *et al.* Functional assessment of human enhancer activities using whole-genome STARR-sequencing. *Genome Biol* 18, 219 (2017).
164. Ernst, J. *et al.* Genome-scale high-resolution mapping of activating and repressive nucleotides in regulatory regions. *Nat Biotechnol* 34, 1180–1190 (2016).
165. Maricque, B. B., Chaudhari, H. G. & Cohen, B. A. A massively parallel reporter assay dissects the influence of chromatin structure on cis-regulatory activity. *Nat Biotechnol* 37, 90–95 (2018).
166. Madeira, C. *et al.* Nonviral Gene Delivery to Neural Stem Cells with Minicircles by Microporation. *Biomacromolecules* 14, 1379–1387 (2013).
167. Skene, N. G. *et al.* Genetic identification of brain cell types underlying schizophrenia. *Nat Genet* 50, 825–833 (2018).
168. Walker, R. L. *et al.* Genetic Control of Expression and Splicing in Developing Human Brain Informs Disease Mechanisms. *Cell* 179, 750-771.e22 (2019).
169. Uffelmann, E. & Posthuma, D. Emerging methods and resources for biological interrogation of neuropsychiatric polygenic-signal. *Biol Psychiat* (2020) doi:10.1016/j.biopsych.2020.05.022.
170. Uebbing, S. *et al.* Massively parallel discovery of human-specific substitutions that alter enhancer activity. *Proc National Acad Sci* 118, e2007049118 (2021).
171. Doan, R. N. *et al.* Mutations in Human Accelerated Regions Disrupt Cognition and Social Behavior. *Cell* 167, 341-354.e12 (2016).
172. Ryu, H. *et al.* Massively parallel dissection of human accelerated regions in human and chimpanzee neural progenitors. *Biorxiv* 256313 (2018) doi:10.1101/256313.
173. Lu, X. *et al.* Global discovery of lupus genetic risk variant allelic enhancer activity. *Nat Commun* 12, 1611 (2021).
174. Dimidschstein, J. *et al.* A viral strategy for targeting and manipulating interneurons across vertebrate species. *Nat Neurosci* 19, 1743–1749 (2016).



175. Dickel, D. E. *et al.* Genome-wide compendium and functional assessment of in vivo heart enhancers. *Nat Commun* 7, 12923 (2016).
176. Dixit, A. *et al.* Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens. *Cell* 167, 1853-1866.e17 (2016).
177. Alyagor, I. *et al.* Combining Developmental and Perturbation-Seq Uncovers Transcriptional Modules Orchestrating Neuronal Remodeling. *Dev Cell* 47, 38-52.e6 (2018).
178. Jin, X. *et al.* In vivo Perturb-Seq reveals neuronal and glial abnormalities associated with Autism risk genes. *Biorxiv* 791525 (2019) doi:10.1101/791525.
179. Geller, E. *et al.* Massively parallel disruption of enhancers active during human corticogenesis. *bioRxiv* 852673 (2019) doi:10.1101/852673.
180. Brandt, M., Gokden, A., Ziosi, M. & Lappalainen, T. A polyclonal allelic expression assay for detecting regulatory effects of transcript variants. *Biorxiv* 794081 (2019) doi:10.1101/794081.
181. Potting, C. *et al.* Genome-wide CRISPR screen for PARKIN regulators reveals transcriptional repression as a determinant of mitophagy. *P Natl Acad Sci Usa* 115, E180–E189 (2017).
182. Tian, R. *et al.* CRISPR Interference-Based Platform for Multimodal Genetic Screens in Human iPSC-Derived Neurons. *Neuron* 104, 239-255.e12 (2019).
183. Lalli, M. A., Avey, D., Dougherty, J. D., Milbrandt, J. & Mitra, R. D. Multiplexed single-cell autism modeling reveals convergent mechanisms altering neuronal differentiation. *Biorxiv* 862680 (2019) doi:10.1101/862680.
184. Fulco, C. P. *et al.* Activity-by-contact model of enhancer-promoter regulation from thousands of CRISPR perturbations. *Nat Genet* 51, 1664–1669 (2019).
185. Zheng, Y. *et al.* CRISPR interference-based specific and efficient gene inactivation in the brain. *Nat Neurosci* 21, 447–454 (2018).
186. Heman-Ackah, S. M., Bassett, A. R. & Wood, M. J. A. Precision Modulation of Neurodegenerative Disease-Related Gene Expression in Human iPSC-Derived Neurons. *Sci Rep-uk* 6, 28420 (2016).
187. Ho, S.-M. *et al.* Evaluating Synthetic Activation and Repression of Neuropsychiatric-Related Genes in hiPSC-Derived NPCs, Neurons, and Astrocytes. *Stem Cell Rep* 9, 615–628 (2017).
188. Colasante, G. *et al.* In vivo CRISPRa decreases seizures and rescues cognitive deficits in a rodent model of epilepsy. *Brain J Neurology* 143, 891–905 (2020).
189. Zhou, H. *et al.* In vivo simultaneous transcriptional activation of multiple genes in the brain using CRISPR–dCas9-activator transgenic mice. *Nat Neurosci* 21, 440–446 (2018).

190. Lau, C.-H., Ho, J. W.-T., Lo, P. K. & Tin, C. Targeted Transgene Activation in the Brain Tissue by Systemic Delivery of Engineered AAV1 Expressing CRISPRa. *Mol Ther Nucleic Acids* 16, 637–649 (2019).
191. Schrode, N. *et al.* Synergistic effects of common schizophrenia risk variants. *Nat Genet* 51, 1475–1485 (2019).
192. Won, H., Huang, J., Opland, C. K., Hartl, C. L. & Geschwind, D. H. Human evolved regulatory elements modulate genes involved in cortical expansion and neurodevelopmental disease susceptibility. *Nat Commun* 10, 2396 (2019).
193. Diao, Y. *et al.* A new class of temporarily phenotypic enhancers identified by CRISPR/Cas9-mediated genetic screening. *Genome Res* 26, 397–405 (2016).
194. Rajagopal, N. *et al.* High-throughput mapping of regulatory DNA. *Nat Biotechnol* 34, 167–174 (2016).
195. Yasuda, K. *et al.* In Vivo Imaging of Human MDR1 Transcription in the Brain and Spine of MDR1-Luciferase Reporter Mice. *Drug Metabolism Dispos Biological Fate Chem* 43, 1646–54 (2015).
196. Consortium, C.-D. G. of the P. G. *et al.* Genomic Relationships, Novel Loci, and Pleiotropic Mechanisms across Eight Psychiatric Disorders. *Cell* 179, 1469-1482.e11 (2019).
197. Consortium, T. B. *et al.* Developmental and genetic regulation of the human cortex transcriptome illuminate schizophrenia pathogenesis. *Nat Neurosci* 21, 1117–1125 (2018).
198. Deming, Y. *et al.* The MS4A gene cluster is a key modulator of soluble TREM2 and Alzheimer’s disease risk. *Sci Transl Med* 11, eaau2291 (2019).
199. Dobbyn, A. *et al.* Landscape of Conditional eQTL in Dorsolateral Prefrontal Cortex and Colocalization with Schizophrenia GWAS. *Am J Hum Genet* 102, 1169–1184 (2018).
200. Chatterjee, S. *et al.* Enhancer Variants Synergistically Drive Dysfunction of a Gene Regulatory Network In Hirschsprung Disease. *Cell* 167, 355-368.e10 (2016).
201. Catizone, A. N. *et al.* Locally acting transcription factors regulate p53-dependent cis-regulatory element activity. *Nucleic Acids Res* (2020) doi:10.1093/nar/gkaa147.
202. Amiri, A. *et al.* Transcriptome and epigenome landscape of human cortical development modeled in organoids. *Sci New York N Y* 362, eaat6720 (2018).
203. Ouwenga, R. L. & Dougherty, J. Fmrp targets or not: long, highly brain-expressed genes tend to be implicated in autism and brain disorders. *Mol Autism* 6, 16 (2015).
204. Pardiñas, A. F. *et al.* Common schizophrenia alleles are enriched in mutation-intolerant genes and in regions under strong background selection. *Nat Genet* 50, 381–389 (2018).

205. Leeuw, C. A. de, Neale, B. M., Heskes, T. & Posthuma, D. The statistical properties of gene-set analysis. *Nat Rev Genet* 17, 353–364 (2016).
206. Sey, N. Y. A. *et al.* A computational tool (H-MAGMA) for improved prediction of brain-disorder risk genes by incorporating brain chromatin interaction profiles. *Nat Neurosci* 1–11 (2020) doi:10.1038/s41593-020-0603-0.
207. Nostrand, E. L. V. *et al.* A large-scale binding and functional map of human RNA-binding proteins. *Nature* 583, 711–719 (2020).
208. Nostrand, E. L. V. *et al.* Principles of RNA processing from analysis of enhanced CLIP maps for 150 RNA binding proteins. *Biorxiv* 807008 (2019) doi:10.1101/807008.
209. Gandal, M. J. *et al.* Shared molecular neuropathology across major psychiatric disorders parallels polygenic overlap. *Sci New York N Y* 359, 693–697 (2018).
210. Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the human brain. *Science* 362, eaat8464 (2018).
211. Consortium, T. Gte. *et al.* The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science* 348, 648–660 (2015).
212. Fromer, M. *et al.* Gene expression elucidates functional impact of polygenic risk for schizophrenia. *Nat Neurosci* 19, 1442–1453 (2016).
213. Kim-Hellmuth, S. *et al.* Cell type specific genetic regulation of gene expression across human tissues. *Biorxiv* 806117 (2019) doi:10.1101/806117.
214. Lonsdale, J. *et al.* The Genotype-Tissue Expression (GTEx) project. *Nat Genet* 45, 580–5 (2013).
215. Gandal, M. J. *et al.* Transcriptome-wide isoform-level dysregulation in ASD, schizophrenia, and bipolar disorder. *Sci New York N Y* 362, eaat8127 (2018).
216. Hoffman, G. E. *et al.* CommonMind Consortium provides transcriptomic and epigenomic data for Schizophrenia and Bipolar Disorder. *Sci Data* 6, 180 (2019).
217. Moore, J. E., Pratt, H. E., Purcaro, M. J. & Weng, Z. A curated benchmark of enhancer-gene interactions for evaluating enhancer-target gene prediction methods. *Genome Biol* 21, 17 (2020).
218. Li, M. *et al.* Integrative functional genomic analysis of human brain development and neuropsychiatric risks. *Sci New York N Y* 362, eaat7615 (2018).
219. Zhou, J. *et al.* Whole-genome deep-learning analysis identifies contribution of noncoding mutations to autism risk. *Nat Genet* 51, 973–980 (2019).
220. Li, Y. *et al.* Genome-wide regulatory model from MPRA data predicts functional regions, eQTLs, and GWAS hits. *Biorxiv* 110171 (2017) doi:10.1101/110171.

221. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat Genet* 48, 245–52 (2016).
222. Consortium, Gte. *et al.* Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. *Nat Commun* 9, 1825 (2018).
223. Loo, L. *et al.* Single-cell transcriptomic analysis of mouse neocortical development. *Nat Commun* 10, 134 (2019).
224. Hodge, R. D. *et al.* Conserved cell types with divergent features in human versus mouse cortex. *Nature* 573, 61–68 (2019).
225. Machiela, M. J. & Chanock, S. J. LDlink: a web-based application for exploring population-specific haplotype structure and linking correlated alleles of possible functional variants. *Bioinformatics* 31, 3555–3557 (2015).
226. Li, D., Hsu, S., Purushotham, D., Sears, R. L. & Wang, T. WashU Epigenome Browser update 2019. *Nucleic Acids Res* 47, W158–W165 (2019).

## Chapter 2: Molecular and functional sex differences of noradrenergic neurons in the mouse locus coeruleus

This chapter was previously published as a journal article:

Mulvey, B. *et al.* Molecular and Functional Sex Differences of Noradrenergic Neurons in the Mouse Locus Coeruleus. *Cell Reports* 23, 2225–2235 (2018).

Preclinical work has long focused on male animals, though biological sex clearly influences risk for certain diseases, including many psychiatric disorders. Such disorders are often treated by drugs targeting the CNS norepinephrine system. Despite roles for noradrenergic neurons in behavior and neuropsychiatric disease models, their molecular characterization has lagged relative to other brain monoaminergic populations. We profiled mouse noradrenergic neurons *in vivo*, defining over 3,000 high-confidence transcripts expressed therein, including druggable receptors. We uncovered remarkable sex differences in gene expression, including female upregulation of the EP3 receptor—which we leverage to illustrate the behavioral and pharmacologic relevance of these findings—and of *Slc6a15* and *Lin28b*, both MDD-associated genes. Broadly, we present a means of transcriptionally profiling locus coeruleus under baseline and experimental conditions. Our findings underscore the need for preclinical work to include both sexes, and suggest that sex differences in noradrenergic neurons may underlie behavioral differences relevant to disease.

### 2.1 Introduction

Numerous neuropsychiatric and neurodevelopmental diseases demonstrate a skew in incidence between sexes, including a female predominance of major depressive and generalized anxiety disorders (MDD and GAD, respectively)<sup>1</sup>, and a male predominance of attention-deficit hyperactivity disorder and autism spectrum disorders (ADHD and ASDs, respectively)<sup>2,3</sup>. Sex

differences in reproductive behavior have been thoroughly attributed to sexually dimorphic brain regions; however, questions remain as to whether more modest behavioral differences—especially those relevant to common psychiatric disorders—are mediated by transcriptional sex differences in key neuronal populations.

The locus coeruleus (LC)—a small nucleus of neuromodulatory neurons whose projections release norepinephrine (NE) throughout the CNS—is implicated in a broad range of functions, including learning, novelty detection, arousal, anxiety, and fever<sup>4</sup>. Given these extensive neurobehavioral roles, it is perhaps unsurprising that the LC-NE system has been broadly implicated in psychiatric disorders and animal models of them. For example, depression is often modeled in rodents using pro-inflammatory compounds including interleukins and lipopolysaccharide (LPS). LPS, besides inducing fever, is known to induce substantial activation of the LC and other noradrenergic cell types (established by<sup>5</sup>). Interleukin-6, a common upstream pro-inflammatory messenger, was recently found to directly trigger tonic firing of the LC, eliciting depressive behaviors via activation of alpha-adrenergic receptors<sup>6</sup>. Tonic firing of the LC in response to corticotropin-releasing factor (CRF) or optogenetic stimulation can also induce acute aversive and anxiety-like behaviors<sup>7,8</sup>. Sex differences in behavioral responses to stress have been attributed to molecular-level differences in CRF signaling via the LC<sup>9,10</sup>. In the clinical setting, evidence-based practices offer a robust demonstration of the involvement of LC dysregulation—or its ability to normalize dysregulation occurring elsewhere—in psychiatric disease, given the use of NE-modulating drugs in depression, anxiety, attention-deficit hyperactivity disorder, and addiction. Altogether, a comprehensive transcriptomic profile of the LC is of substantial interest for understanding, and potentially targeting, the molecular functions of these cells.

To date, sex differences in the rodent LC have been observed at both single-gene and structural level. Sex differences in stress response have been attributed to differential CRF sensitivity and CRF receptor trafficking in mouse LC<sup>11,12</sup>, including sex-differential effects of CRF1 agonism on LC excitability<sup>13</sup>.  $\mu$  opioid receptors are also highly expressed in the LC;  $\mu$  agonism completely suppresses firing of the LC in male, but not female, mice<sup>14</sup>. Reports of structural dimorphism in the rodent LC have been ambiguous, with reports of the LC being larger in either sex depending on the strain of rat (compare<sup>9,15,16</sup>). These repeated demonstrations of sex differences in particular aspects of LC structure and function compelled us to study both sexes in our pursuit of characterizing LC gene expression.

In order to transcriptionally profile the LC, we generated a translating ribosome affinity purification (TRAP) line. We identified dozens of potential LC-specific drug targets, and validated a subset of these with independent methods. We identified differentially-expressed genes (DEGs) in the LC following LPS stimulation, demonstrating the utility of this line and method to detect pharmacologically-mediated changes in gene expression in the LC. To our surprise, we discovered a comparable number of DEGs between sexes as well. In order to demonstrate that these transcript-level sex differences in the LC correspond to consequent physiologic differences, we modulated a receptor upregulated in female LC, EP3 (encoded by *Ptger3*). Using electrophysiology and behavior experiments in cannulated mice, we demonstrate that the EP3 agonist sulprostone acts more strongly in female mice to suppress tonic firing of LC neurons *in vitro* and to specifically inhibit an LC-mediated stress response in females *in vivo*. Thus, we demonstrate previously unidentified sex differences in gene expression in NE neurons at a magnitude capable of

influencing neurophysiology and pharmacologic responses. These molecular sex differences at the level of LC neurons may guide future investigations into models, mechanisms, or treatments for sex-skewed psychiatric diseases.

## 2.2 Results

### 2.2.1 Generation and validation of reagents for transcriptional profiling of noradrenergic neurons

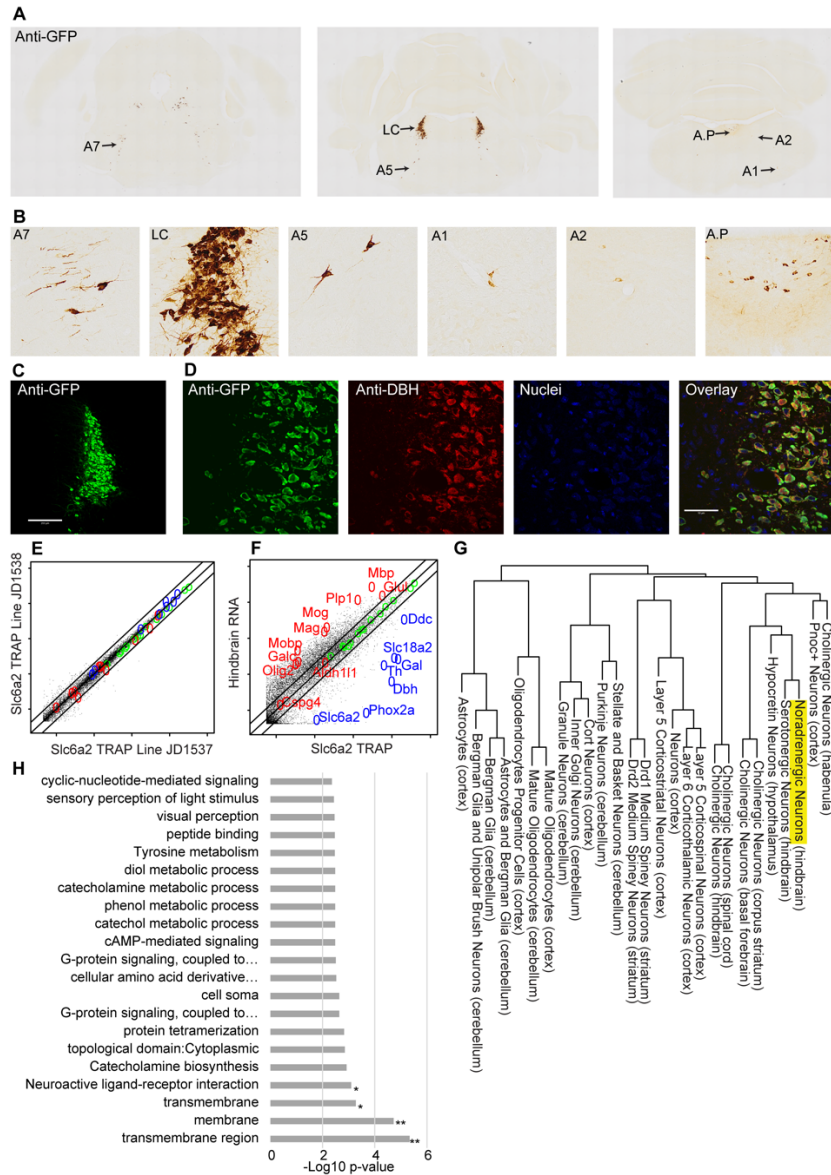
We generated a mouse line for transcriptional profiling of noradrenergic neurons by expressing EGFP/RPL10A from a NE reuptake transporter (*Slc6a2*) bacterial artificial chromosome (BAC). Neuroanatomical characterization revealed robust transgene expression in the A4 and A6 subdivisions of the LC (**Figure 2.1A-C**), where EGFP/RPL10A perfectly co-localized with the LC-specific NE-synthesizing protein, DBH (**Figure 2.1D**), along with reasonably robust co-labeling in the A5 and A7 groups. EGFP/RPL10A labeling was weak and sparse in more caudal DBH+ populations (e.g., A1, A2), consistent with prior immunofluorescence studies of SLC6A2 expression<sup>17</sup>. Little EGFP expression was seen in ependymal cells, except rarely in cells caudal to the 4<sup>th</sup> ventricle (*not shown*), in contrast to previously reported ependymal expression of SLC6A2. In total, this anatomical characterization asserts that mRNA collected by TRAP will be from the most robustly labeled and populous cells: the A4-A7 groups, predominantly the LC.

We then performed TRAP on two *Slc6a2* founder lines to evaluate consistency, to confirm enrichment of known LC-specific transcripts by TRAP, and to identify novel transcripts enriched in LC compared to the hindbrain. Reproducibility was strong between the lines (Pearson correlation >.99, **Figure 2.1E**). Relative to whole-hindbrain RNA from the same mice, TRAP enriched for genes with known specificity and functionality in the LC (**Figure 2.1F**), including **1**)



enzymes related to NE turnover (*Th*, *Ddc*, *Maoa*, and *Dbh*), **2**) Vesicular monoamine (*Slc18a2*) and NE (*Slc6a2*) transporters, **3**) Galanin (*Gal*) and its receptor (*GalR1*), and 4) a transcriptional regulator of LC development (*Phox2a*). After conservative filtering for expression and background, at least 3139 transcripts were detected with high confidence in NE neurons; 526 were enriched >2-fold compared to hindbrain.

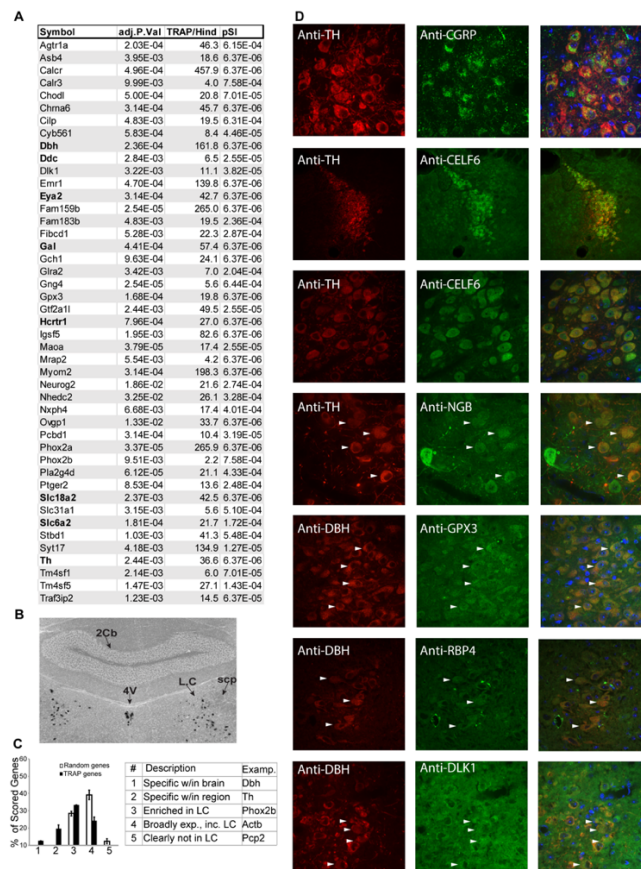
Gene Ontologies (GO) analysis was applied to broadly characterize NE neuron-enriched transcripts, revealing enrichment of transmembrane receptors and ligands (**Figure 2.1H**), consistent with prior observations that CNS cell-type specific genes often include receptors<sup>18</sup>. These LC-enriched receptors are notable given the importance of the LC—and extrinsic modulation thereof—in behavior. Next, transcriptional profiles of LC and other CNS cell subtypes previously characterized by TRAP were compared to predict genes specific to the LC compared to other cells in the brain and the ‘transcriptional ontology’ of the LC using our previously described pSI algorithm (see *Experimental Procedures*). Compared to all other available cell types, 162 genes scored as LC-enriched (pSI<0.05), 78 highly so (pSI<.005) (**Figure 2.2A**). Hierarchical clustering placed noradrenergic cells with other neuromodulatory populations (**Figure 2.1G**), including serotonergic, hypocretinergic (Hcrt), and forebrain cholinergic neurons. The shared transcriptional relationship with Hcrt neurons is striking; despite their anatomic separation and the use of distinct neurotransmitters, they form reciprocal connections and share functional roles in control of sleep and arousal (reviewed in<sup>19</sup>), supporting the notion that gene expression in neuronal cell subtypes is a strong predictor of a subtype’s functional roles.



**Figure 2.1. Characterization of noradrenergic bacTRAP lines.** **A)** Anti-GFP immunohistochemistry demonstrates EGFP/Rpl10a labeling in the hindbrain. **B)** Anti-GFP staining is most robust in anterior groups (A4-A7), especially LC. **C)** Immunofluorescence for GFP (green) labels entire LC (scale bar: 200uM) **D)** GFP colabels completely with DBH (red) (scale bar, 50uM). **E)** Comparison of two *Slc6a2* TRAP lines demonstrates reproducibility. **F)** *Slc6a2* TRAP mRNA vs. total hindbrain mRNA enriches NE neuron markers (blue) and depletes unrelated cell-type (glial) markers (red). Lines at 0.5, 1, 2 fold. Log10 scale. **G)** Hierarchical clustering of *Slc6a2* neurons. **H)** *Slc6a2* TRAP enriches for transmembrane proteins and receptors. (Hypergeometric test, Benjamini-Hochberg corrected \* $p < .06$ , \*\* $p < .01$ ).

Transcripts identified as LC-enriched (**Figure 2.2A**) were then validated using standard RNA and protein detection methods in wild-type mice. *In-situ* hybridization (ISH) for *Calcr* selectively and

robustly stained the LC (**Figure 2.2B**), consistent with microarray results, suggesting very high enrichment of *Calcr* in LC compared to hindbrain (over 300-fold). As the characteristic ‘quarter-moon’ anatomy of the LC could be readily discerned by ISH, additional transcripts were systematically evaluated for enrichment using The Allen Brain Atlas (**Figure 2.2C** and **Supplementary Figure 2.1**). 70% of TRAP transcripts showed enriched *in-situ* staining in LC; over 19% scored as having ‘marker-like’ expression. Finally, we confirmed protein translation in LC of several identified genes using immunofluorescence (**Figure 2.2D**).



**Figure 2.2. Transcript and protein expression in LC neurons. A)** Top 45 named genes enriched by TRAP over hindbrain (adj.p.value), fold change (TRAP/Hind), and Specificity index p-value (*p*SI) comparing *Slc6a2* to all cell populations from **Figure 2.1G**. **B)** ISH confirms LC enrichment of calcitonin receptor. **C)** Blind comparison of TRAP-identified and random genes confirms TRAP transcript presence in LC ( $p < 2.7E-38$ ,  $\chi^2$  test, normalized to number of scorable ISH,  $n = 53,94$ ). (See also **Supplementary Figure 2.1**). **D)** Immunofluorescence confirms translation of TRAP identified transcripts (green) in LC cells (arrowheads), labeled by either TH or DBH (red). 4V = 4th ventricle; scp = superior cerebellar peduncle; 2cb = cerebellar vermis.

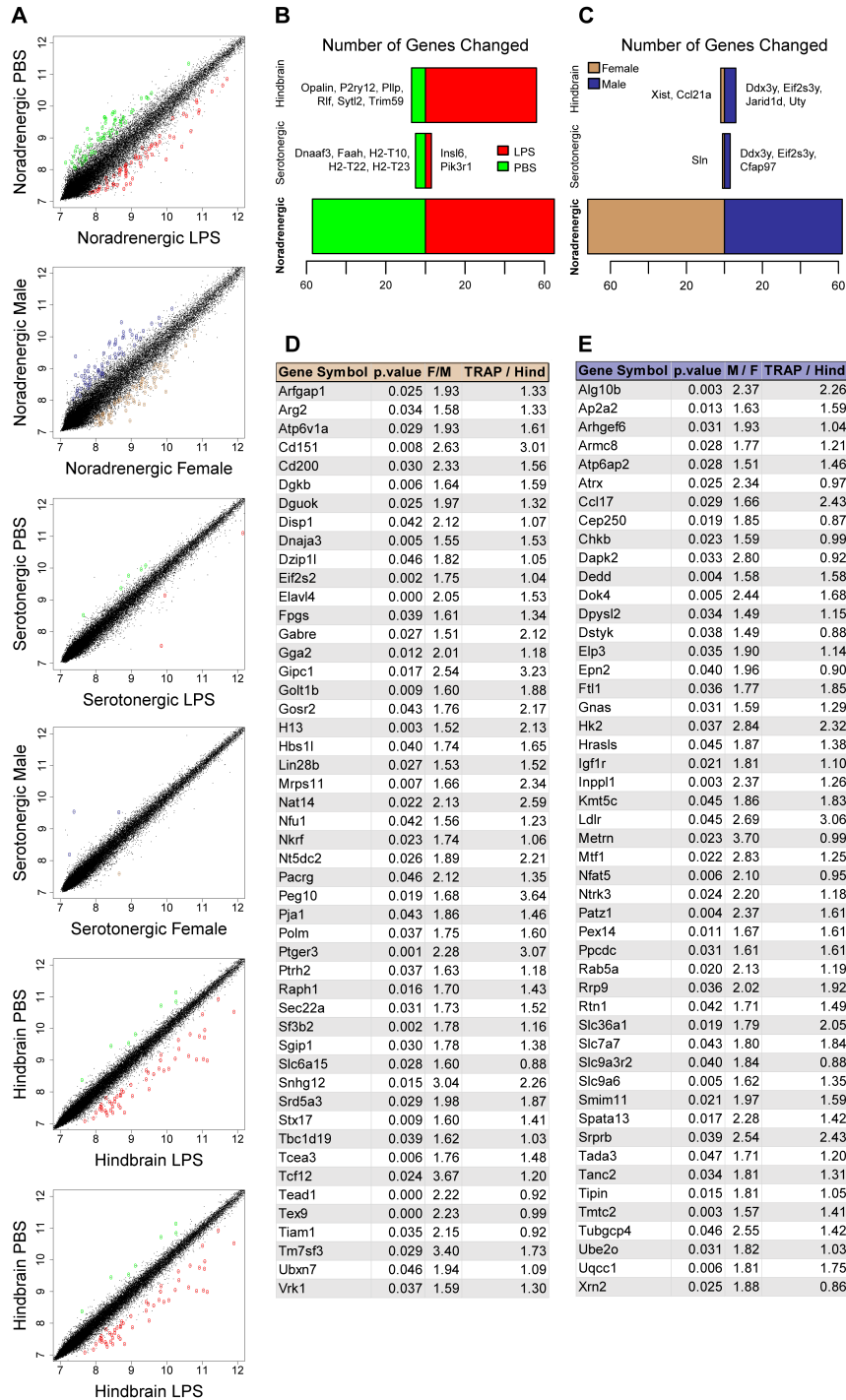
### 2.2.2 Transcriptional responses of NE neurons to LPS can be identified with *Slc6a2* TRAP

Having verified that TRAP characterizes baseline transcriptional features of the LC, we next sought to demonstrate the utility of this mouse line for profiling changes consequent to stimulation of the LC. LPS is a well-characterized example of an LC-activating stimulus: it strongly increases *FOS* expression in the LC, and the LPS-induced febrile response is lost with LC ablation<sup>20</sup>. We therefore injected individual TRAP mice with LPS or vehicle in a sex-balanced design. Whole-hindbrain RNA from the same lysates served as controls, as did a parallel TRAP experiment with *Slc6a4* TRAP mice targeting hindbrain serotonergic neurons (**Figure 2.3**).

We first examined whole-hindbrain changes in response to LPS (**Figure 2.3B**); 56 genes showed a response, predominantly upregulation. GO analysis identified a significant increase of interferon-induced transmembrane proteins (**Supplementary Figure 2.2**), consistent with a broad pro-inflammatory transcriptional response. However, examination of *Slc6a2* TRAP revealed an even greater response (**Figure 2.3C**), largely distinct from that of the hindbrain (only two genes were shared between hindbrain and LC). In contrast, response of serotonin neurons was limited (**Figure 2.3B**). Given the multitude of psychopharmacotherapeutics and environmental stimuli—including inflammation, pain, and acute stress—known to modulate the LC, these findings establish a useful system for identifying LC-specific molecular responses to whole-animal manipulations.

To our surprise, sex-stratified analyses of these same experimental data revealed substantial sex differences in noradrenergic neurons. First, analysis of transcriptional profiles of the whole hindbrain varied remarkably little between sexes, with the exception of a few sex chromosomal

genes (**Figure 2.3C**). Likewise, serotonin neurons showed no appreciable differences outside of sex chromosomal transcripts (e.g. *Ddx3y*, *Eif2s3y*). In contrast, noradrenergic neurons showed substantial molecular sex differences: a total of 152 LC-enriched transcripts were also sex DEGs, mostly autosomal (**Figure 2.3D-E**), and these did not overlap with the serotonin neuron-enriched sex-differential differential transcripts. These transcripts also do not significantly overlap with those stimulated by LPS, nor are they characterized by similar functional categories (**Supplementary Figure 2.2**). This indicates the observed molecular divergence is not likely due to sex differences in baseline activity level of the LC, but rather reflects a more complex molecular distinction between the sexes. Motivated by the clear role of sex as a risk factor for psychiatric disorders, we focused on these sex differences for additional study.



**Figure 2.3. Differentially-expressed genes by sex.** A) Scatterplots contrasting transcriptional data between LPS and vehicle (PBS) or sex; DEGs indicated by color. B-C) Number of DEGs after LPS (B: red=up, green=down) or between vehicle-treated sexes (C) in each sample type. D-E) Top 50 named sex-DEGs in LC of females (D) or males (E). M/F or F/M: Fold change between sexes. TRAP/Hind: Fold change, TRAP vs hindbrain. (See also **Supplementary Figure 2.2**).

### 2.2.3 Sex-differential LC genes and putative *cis*-elements underlying differential regulation

As a preliminary investigation into possible gene-regulatory mechanisms underlying sex differences in LC gene expression, I characterized DNA sequence motifs in *cis* with these 152 genes. Performing *de novo* motif discovery for DEGs from each sex, I identified twelve motifs (**Figure 2.4**). Compared to their frequency near 1,000 randomly selected genes, six of these were significantly enriched near the DEGs. Thus, at least a portion of the sex differences in gene expression could be explained by conserved *cis*-regulatory elements in the surrounding genome. Known transcription factor binding sites predicted in these motifs included *OTX2*, *NR2F6*, and *MTF1* (**Figure 2.4**).

Motif	Total # motif occurrences found	# of unique sequences matching	E-value	Chi-sq p-value (vs. 1,000 random)	Select TFBSes Predictions (TomTom)	Sequence
<b>Female LC-Enriched TSS regions</b>						
Fz4	26	26 / 77	1.7E-58	4.59E-03	NR4A2, NR2F6, OTX2	
Fz9**	19	19 / 77	2.7E-24	3.60E-11	(None)	
Fz10	13	13 / 77	3.4E-21	0.5518	OTX2	
Fz4	37	22 / 77	8.0E-97	0.1426	OTX2	
Fz10**	22	17 / 77	1.1E-41	8.28E-13	FOS, FOS::JUN, TBX1, MAZ	
Fz12**	31	19 / 77	6.3E-63	5.43E-07	(None)	
<b>Male LC-Enriched TSS regions</b>						
Mz1	15	15 / 70	1.5E-46	0.5645	(None)	
Mz2**	20	20 / 70	1.5E-19	4.50E-04	MTF1, OTX2	
Ma3*	31	24 / 70	2.1E-104	5.17E-03	NR2F6, OTX2	
Ma4*	23	18 / 70	3.1E-61	0.0198	OTX2	
Mv1	13	13 / 70	1.0E-44	1.00	Let-7, NR2F6, MAZ	
Mv3	11	11 / 70	2.1E-07	0.1531	(None)	

**Figure 2.4: Motifs discovered in conserved, noncoding regions near sex-DEGs.** *Motif*: Refers to the sex, algorithm (allowing  $\leq 1$  or  $\geq 1$  motif occurrence per query sequence) and the result # in the MEME output. *Total # motif occurrences found*: total number of significant matches for the motif. *# Unique sequences matching*: the number of queried sequences with  $\geq 1$  motif match out of the total number of sequences queried. *E-value*: a MEME measure of probability compared to

a shuffled version of the input sequences. *Chi-sq p value*: from comparative abundance of the identified motif in 1,000 randomly selected, protein-coding gene regions subjected to the same masking paradigm. *Select TFBS predictions*: TFBSes identified using TomTom with previously described expression in neural cell types and of functional and/or cell type relevance to the LC. *Sequence*: the position-weight matrix (PWM) given in the MEME results.

#### **2.2.4 Molecular differences predict functional differences between sexes**

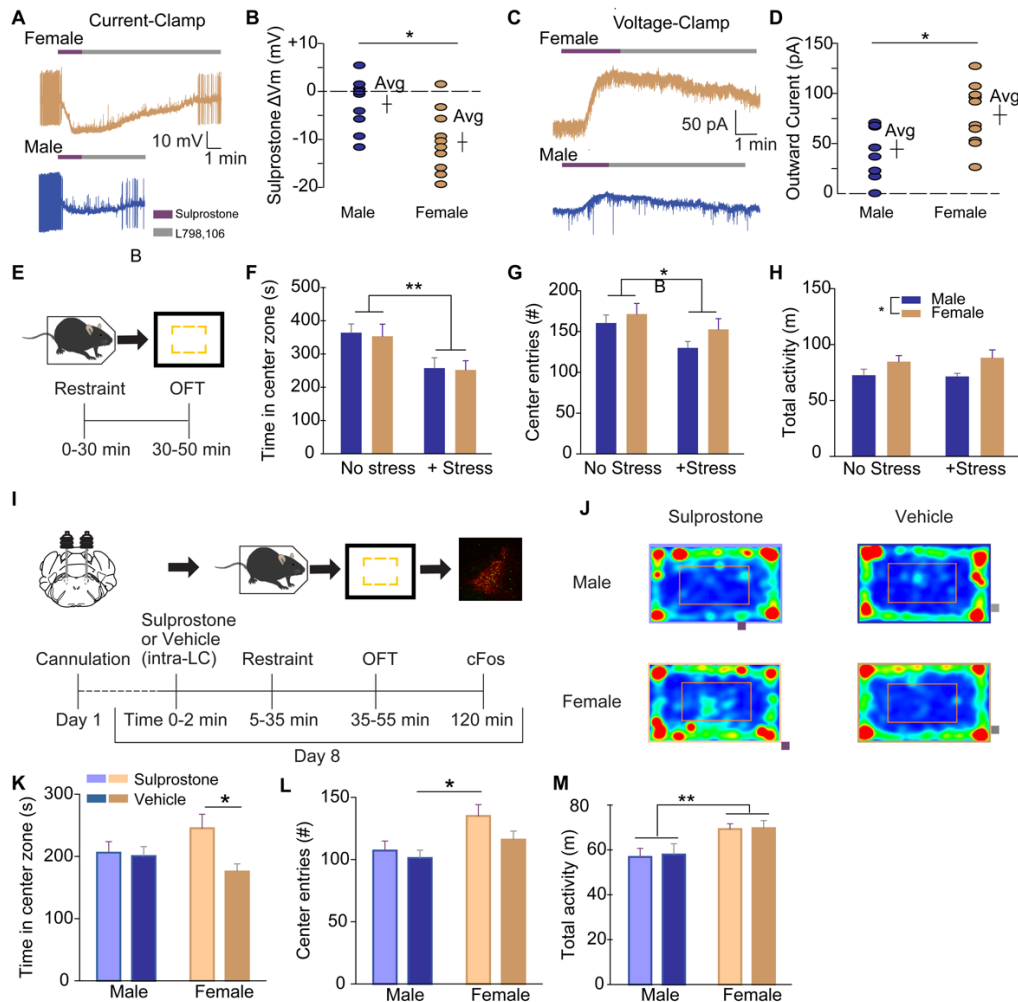
Finally, we noted LC specificity and female LC enrichment (>2-fold) of the *Ptger3* gene, encoding prostaglandin E2 (PGE2) receptor EP3. Given the specificity of this gene's expression to the LC within the hindbrain and the existence of a known, selective agonist, sulprostone, we selected EP3 to pharmacologically test whether the magnitudes of detected sex differences in receptor gene expression were adequate to alter LC electrophysiology and/or behavior.

We first assessed whether pharmacologic manipulation of the EP3 receptor resulted in sex-differential electrophysiologic responses by performing whole-cell recordings from LC neurons in *ex vivo* slices. EP3 presence in the LC was confirmed—and subsequently manipulated—by bath application of sulprostone (agonist), followed by L798,106 (antagonist) to displace sulprostone, halting its effects. Sulprostone suppressed baseline tonic firing of LC neurons in both sexes, but with a greater magnitude and duration of hyperpolarization in female LC compared to male (**Figure 2.5A-B**). Voltage-clamp recordings from a second cohort of mice revealed larger outward current from female LC neurons (**Figure 2.5C-D**), verifying that the magnitude of LC inhibition by EP3 corresponds to sex differences in *Ptger3* expression.

We previously showed in male mice that LC silencing with Gi-coupled DREADDs prevents anxiety-like behavior in the open-field task (OFT) after restraint stress<sup>7</sup>. This behavioral paradigm provided a robust model system in which we could activate the LC *in vivo*, and subsequently



attempt to suppress LC pharmacologically, with behavior as the outcome measure. We first validated that restraint stress robustly induces anxiety-like behavior (avoidance of center) in the OFT in mice of *both* sexes (**Figure 2.5E-H**). Restraint stress did not impact total activity or the sex difference therein, but clearly induced avoidance of center in both sexes. Thus, we hypothesized that EP3 agonists could be used to reduce stress-induced anxiety specifically in wild-type female mice. Indeed, administration of sulprostone via cannula to male and female LC immediately before restraint stress and OFT resulted in selectively reduced anxiety-like behavior in females (**Figure 2.5I-L**). Post-hoc analysis of phosphorylated-FOS (p-FOS) expression in TH+ neurons of the LC in the same animals revealed similar staining intensity of p-FOS and numbers of double-positive cells between sexes (*data not shown*). The robust immediate-early response in both conditions are consistent with intact stimulation of LC neurons by stress, suggesting sulprostone inhibits noradrenergic output. Sulprostone did not affect the baseline sex difference in total ambulatory activity (**Figure 2.5M**).



**Figure 2.5. Sex differences in *Ptger3* expression can be reflected in LC-mediated behavior.** **A, C)** Representative current-clamp (**A**) and voltage-clamp (**C**) traces from LC slices exposed to sulprostone (200nM) followed by L798,106 (300 nM). **B)** Maximum change in membrane voltage (mV) after sulprostone, by sex ( $-10.5 \pm 2.1$  mV,  $n=10$  cells from 5 females,  $-2.6 \pm 1.8$  mV,  $n=9$  cells from 4 males;  $p < 0.05$ , Mann-Whitney). **D)** Maximum change in outward current (pA) after sulprostone, by sex ( $78.5 \pm 9.7$  pA,  $n=10$  cells from 3 females,  $44.4 \pm 8.9$  pA,  $n=9$  cells from 3 males,  $p < 0.05$ , Mann-Whitney). **E)** Schematic of validated stress-anxiety paradigm. **F-H)** Effects of restraint stress on OFT task performance ( $n = 6-7$ ). **I)** Timeline of LC pharmacology-behavior experiments. **J)** Representative OFT traces from each sex and treatment condition. **K-L)** Sulprostone administration prevents stress-induced anxiety in female ( $n=17,17$ ) but not male ( $n=16, 14$ ) mice in OFT (center zone entries and time, respectively). (See also **Supplementary Figure 2.4** for estrous-stage specific behaviors). **M)** Total activity was unaffected by sulprostone for both sexes (\*:  $p < 0.05$ , \*\*:  $p < 0.01$ ).

## 2.3 Discussion

Characterizing gene expression in noradrenergic neurons—regardless of whether they are at the etiologic root of neuropsychiatric diseases and disease models—is key to narrowing down possible mechanisms by which NE signaling may be dysregulated in disease states. We have presented a mouse line enabling transcriptional profiling of LC neurons at baseline and after physiologic manipulations (LPS). In characterizing this line, we discover and herein report a breadth of previously unidentified sex differences in molecular features of the mouse LC. These findings highlight the LC as an area of focus for future studies in neuropsychiatry—especially in domains where sex differences are observed in modeled behaviors or diseases. In contrast, we find that serotonergic neurons show few sex differences in gene expression, despite their hypothesized role in behavior and psychiatric disease. Sex-differential expression of one such receptor, *PTGER3*, was adequate in magnitude to sex-differentially affect electrophysiologic and behavioral pharmacologic responses. This independent verification of our transcriptomic findings suggests sex-differentially expressed genes in LC 1) may underlie sex differences in behavior and behavioral pathology and 2) can be targeted to sex-specifically modulate LC-mediated behaviors. Thus, we conclude that the LC is an interesting candidate for mediating sex differences in monoamine-associated psychiatric phenotypes. We further envision that this mouse line could provide an invaluable tool in studies aimed at identifying mechanisms of existing NE-targeting drugs at the transcriptomic level, and enable prioritization of new, precise drug targets aimed at the same transcriptional endpoints.

Our profiling extends previous work illustrating discrete molecular sex differences in the LC. Perhaps best characterized is the trafficking of receptor CRF1 and response to its ligand, CRF<sup>11-</sup>

<sup>13,21</sup>. Likewise, the expression of the  $\mu$  opioid receptor and response to opioid agonism in the LC shows a sex difference<sup>14</sup>. Estrogen regulates genes required for NE synthesis in a sex-specific fashion (see below). Structural dimorphism in the rat LC has also been observed, though the direction of effect depends on the strain of rat<sup>9,15,16</sup>. Finally, postnatal citalopram exposure in rats causes ectopic projection of LC fibers into the neocortex and increased LC excitability in males, but not females<sup>22</sup>. We expand upon this body of research by identifying thousands of genes expressed in LC and sex differences therein, as well as recurrent, conserved motifs in *cis* with these genes.

Among the transcripts with sex differences identified in the LC, we identified a number of genes and putative *cis*-regulators notable for their previous implications in behavior and brain development. Putative regulators in *cis* with the DEGs included three striking candidates: *OTX2*, *NR2F6*, and *MTF1*. *OTX2* was recently shown to regulate depression-related consequences of early life stress in male mice (females were untested) through actions in dopaminergic neurons<sup>23</sup>. *MTF1* is notable for its role in binding and responding to heavy metals, perturbations of which have been implicated in ASDs<sup>24,25</sup>; furthermore, this transcription factor was itself enriched in male LC, providing hints of a potential regulator of some of our observed sex differences. Finally, *NR2F6* is a nuclear receptor known to be required for LC differentiation, consistent with LC enrichment of the genes used for motif analysis. As the motif analysis only utilized conserved regions of mammalian sequence near these genes, these regulatory mechanisms, and thus sex differences, may be conserved in humans.

Intriguingly, we also noted a previously unidentified female enrichment of the prostaglandin E2 (PGE2) receptor *Ptger3* (EP3) in the LC. PGE2 and PTGER family receptors are known to mediate sexually dimorphic neurodevelopment in the preoptic area of the hypothalamus<sup>26,27</sup>; sex-differential expression of these receptors in a separate, adult brain region was thus intriguing. The enrichment of PGE2 receptors is interesting in the context of LPS, which we used here to stimulate LC, but also stimulates fever; PGE2 and the LC are major effectors of LPS-induced fever via the EP3 and EP4 receptors<sup>20,28</sup>. Follow-up studies are merited to explore whether the EP3 receptor plays a role in fever effects on behavior via the LC, and whether its differential expression is cause or consequence of the broader transcriptomic sex differences presented here.

This expression difference was sufficient to modulate behavior in a sex-specific way, which we validated by LC-targeted pharmacologic manipulation. Using sulprostone to agonize EP3, we identified strong inhibitory effects on LC firing in female LC neurons, consistent with the pattern of its increased female expression. We then utilized restraint stress as a validated means of activating the LC and triggering LC-mediated behavioral changes<sup>7,8,29,30</sup>. In turn, we aimed to suppress restraint-driven LC activation by administering sulprostone beforehand, ameliorating behavioral signs of stress-induced LC activity in female, but not male, mice. We thus demonstrate that the sex differences in receptor expression measured in the LC by TRAP are of an adequate magnitude to manipulate an LC-regulated behavior in a sex-specific manner.

Our stress paradigm was only used to robustly activate the LC, rather than to investigate stress *per se*. Whether EP3 plays a role in—or undergoes transcriptional or translational regulation in response to—physiologic sex differences in the stress response of LC remains unclarified. We also

note that the mice in which DEGs were identified were singly-housed (potentially a stressor, i.e. social isolation), and housed at an unstressful, thermonormal 30°C for fever experiments. We note, however, that our electrophysiologic and behavioral findings regarding *Ptger3* were consistent with the observed expression changes, despite the mice for the later experiments being group-housed at a normal room temperature. Using these TRAP mice to deliberately study sex-specific transcriptional/translational responses to stress will be an interesting application of this mouse line for future investigations.

It is interesting to speculate that sex differences in LC gene expression may specifically influence increased female risk of disorders like GAD and MDD, where NE-modulating drugs have seen use for decades. If higher baseline expression of some genes in the female LC promotes risk for MDD, then common variants that elevate expression of those same genes may likewise confer depression risk. Indeed, when we examined the 15 documented MDD-associated loci<sup>31</sup> for the presence of sex-differential LC genes, we find two genes are associated with MDD and enriched in female LC of mouse: *Slc6a15* and *Lin28b*. **This striking coincidence may imply that certain sex-and variant-mediated MDD risk factors converge in the LC. Future research is warranted to explore whether sex and disease-associated regulatory variants concordantly affect gene expression and psychiatric disease risk via LC and other cell populations.**

Overall, the marked molecular sex differences present interesting areas for future inquiry. Most notably, the mechanism of establishing these sex differences (hormonal-developmental, sex chromosomal, or post-pubertal hormonal) remains unclarified. Previous work has shown that estrogen regulates expression of *Th* and *Dbh*, and thus NE synthesis, in a sex-differential manner

in adult rodents<sup>32,33</sup>. In the present study, estrous cycling was not examined for effects on transcriptional sex differences. Regarding behavior, we note that estrous cycling does not appear to play a role in most behaviors<sup>34</sup>, including center time in the open field task for C57BL6 mice<sup>35</sup>; binning behavior data by estrus stage revealed no substantial differences (**Supplementary Figure 2.4**). Another possible mechanism, the perinatal masculinizing hormonal surge in male rodents that organizes other dimorphic regions, might be equally important in the LC—indeed, structural differences in female rat LC can be attenuated by perinatal testosterone administration<sup>36</sup>. Overall, it is possible that multiple mechanisms contribute to the molecular sex differences we detected. Thus, future studies are warranted, both focusing on identifying these mechanisms and examining the potential conservation of these differences.

## 2.4 Experimental Procedures

*Animal Research Statement.* All procedures involving animals were approved by the Institutional Animal Care and Use Committees of Rockefeller University, Case Western Reserve University, and Washington University in St. Louis.

*Immunofluorescence microscopy.* PFA-perfused mouse brains were dissected, cryoprotected, and cut by cryostat into floating sections for immunostaining with primary antibodies and Alexa fluorophore-coupled secondary antibodies. Antibodies used are described in **Supplementary Table 2.1**.

*TRAP for initial description of LC.* Replicate pools of five mixed sex adult mice from each of the *Slc6a2* lines were sacrificed. Brains were removed for collection of hindbrain posterior to the

pontine/hypothalamic junction (discarding the cerebellum). All array data were analyzed in R using Bioconductor packages. *GCRMA* was used to normalize within replicates, and to biotinylated spike in probes (green dots, **Figure 2.1F**) between conditions. Fold change, Specificity Index (SI) and pSI were calculated for genes expressed above non-specific background, defined as the mean + 2 SD of the TRAP:hindbrain fold change of negative control transcripts (**Figure 2.1F**, red dots). LC-enriched transcripts were identified using the empirical Bayesian statistic with FDR correction in *limma*. The pSI algorithm was used with default settings to compare LC TRAP to other cells profiled by TRAP (**Figure 2.1G**). Hierarchical clustering across cell types was conducted in R utilizing expression values from genes with pSI <0.01 in any cell type.

*Scoring of the Allen Brain Atlas for LC gene specificity.* Any transcript that scored between 1 and 3 (examples in **Supplementary Figure 2.1**) was considered to be marker-like. A  $\chi^2$  test was performed comparing observed counts of each score with expected counts (based on the random gene set) for each score.

*Single animal TRAP for sex-differential and LPS-responsive gene expression.* Samples were prepared and hybridized in two batches counterbalanced for mouse strain, sex, and LPS. After processing, one *Slc6a4* TRAP sample and two hindbrain samples were excluded due to poor hybridization. Remaining samples were normalized using the *lumi* package. Appropriate clustering of replicates was confirmed with multidimensional scaling (MDS) plots (**Supplementary Figure 2.5**). Differential expression was defined as  $p < .05$  on a paired T-test (paired on batch and covariate) and Log2 change was +/- .585 across 3 of 4 paired comparisons. For balance, the single *Slc6a4*



replicate was used twice to replace the low quality sample. To validate this analytic approach, we performed standard differential expression analysis in *limma*. The complete gene-wise analysis result tables are available with the published version of this work as supplemental tables.

*Motif Analysis of peri-TSS Sequences of Sex-DEGs in LC.* For each sex-differentially expressed gene identified in the LC, mm10 genomic sequence 10kb 5' and 10kb 3' to the TSS were acquired. Exonic bases non-conserved regions of sequence (based on PhyloP scores) were masked out. Masked flanking sequences were submitted to MEME<sup>37</sup> for *de novo* discovery of motifs 8-20bp long. The associated tool, TomTom, was used to compare motifs to known TF binding sites<sup>38-40</sup>. Motif frequency, consensus sequence, and predicted TFs discussed are provided (**Figure 2.4**). To assess motif enrichment near LC transcripts, 1000 random protein-coding TSSes were selected and processed identically. These were searched for motif matches using FIMO at the same *p* cutoff for a “match” used by MEME during discovery. The number of unique loci containing  $\geq 1$  motif match among all loci was then compared using chi-square analysis, followed by Benjamini-Hochberg correction.

*Electrophysiology of LC neurons exposed to PTGER3 agonist/antagonist.* Whole-cell recordings were made using an Axopatch 200B amplifier (*Molecular Devices*). LC neurons were identified by location, capacitance > 40 pF, an input resistance < 100 M $\Omega$ , and a tonic firing rate of 0.5 – 4 Hz.

*Stereotaxic cannulation of LC.* Mice were allowed to recover from surgery 7-9 days prior to behavioral testing. Animals were also habituated to handling and connection to tubing for 3 consecutive days prior to behavioral testing.

*Stress-induced anxiety behavioral paradigm.* Anymaze was used for video recording of animal movements for center and periphery analysis. The center zone was defined as a concentric rectangle comprising 50% of the OFT area. Cannula placement was confirmed by cryostat sectioning of perfused brains (**Supplementary Figure 2.3**) to determine mice for inclusion in the final behavioral and c-FOS analyses.

*c-FOS quantification in LC following sulprostone/vehicle, restraint, and OFT.* Gain, light intensity, and exposure time were identical for all prepared microscope slides. Using ImageJ: background was subtracted, ROIs made around the LC based on TH staining, and average pixel intensity for c-Fos fluorescence was measured. Fos-TH double positive cells were counted manually by a blinded experimenter.

## **2.5 Acknowledgments**

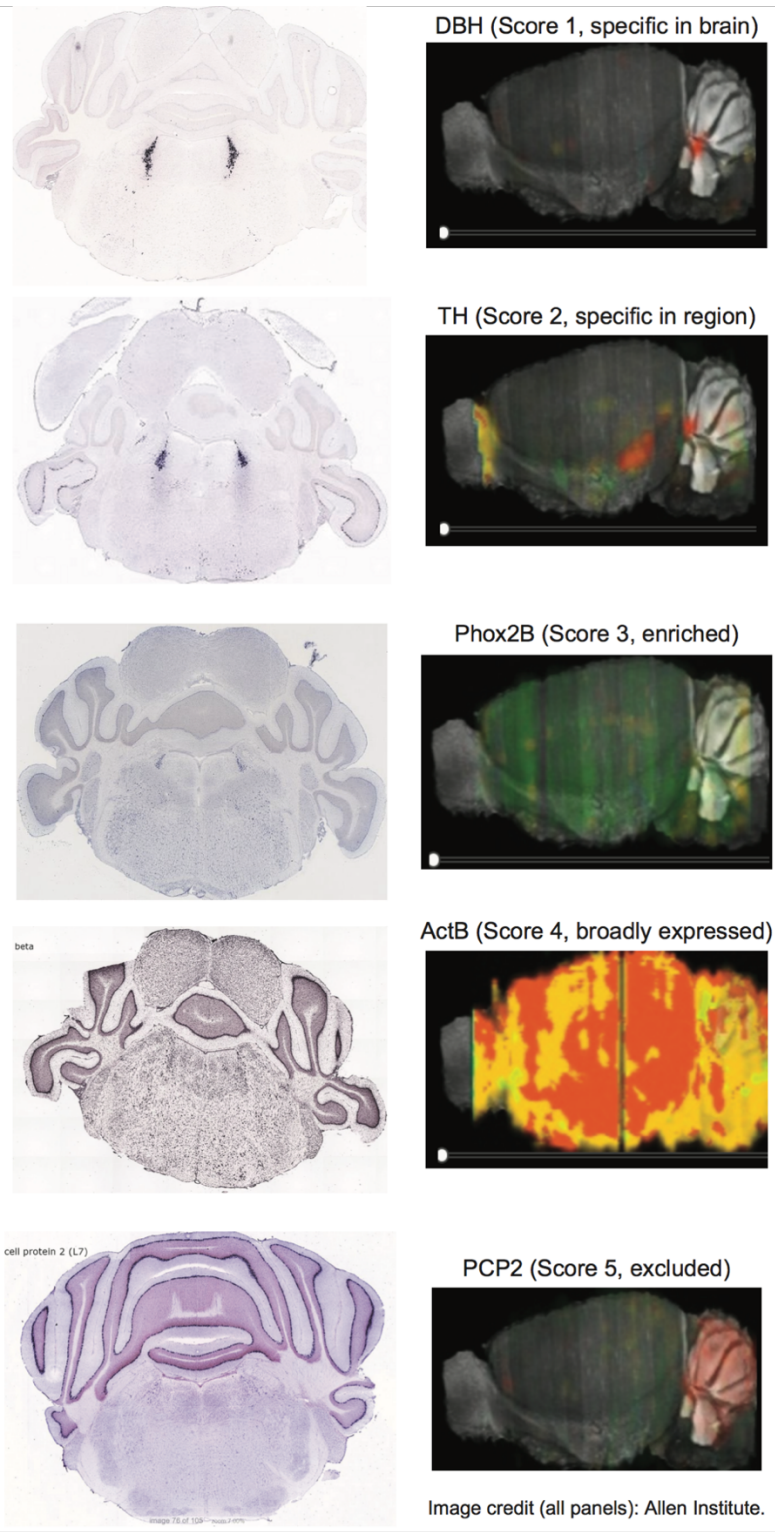
Thanks to R. Jaswany, E. Park, C. Jakes, K. McCullough, and L. Broestl for assistance in performing lab experiments; C. Weichelbaum for editorial assistance; the Rockefeller University Genomics Resource Center and Bioimaging Cores. This chapter is dedicated to T.C. Mazer, for discussion of potential hormonal mechanisms of these findings; rest in peace. This work was supported by the NIH (5R01HG008687, 4R00NS067239, 5R21DA038458, R01DA035821, R01NS095809), the Simons Foundation, the Brain and Behavior Research Foundation, Ludwig

Cancer Research, BBSRC (BB/M001873/1), and a CDI Microgrant from the Washington University Center for Cellular Imaging (WUCCI). J.D.D. has previously received royalties for patents related to TRAP technology. The remaining authors declared no competing interests. All transcriptomic data are available at GEO: <https://www.ncbi.nlm.nih.gov/geo/>, accession GSE100005.

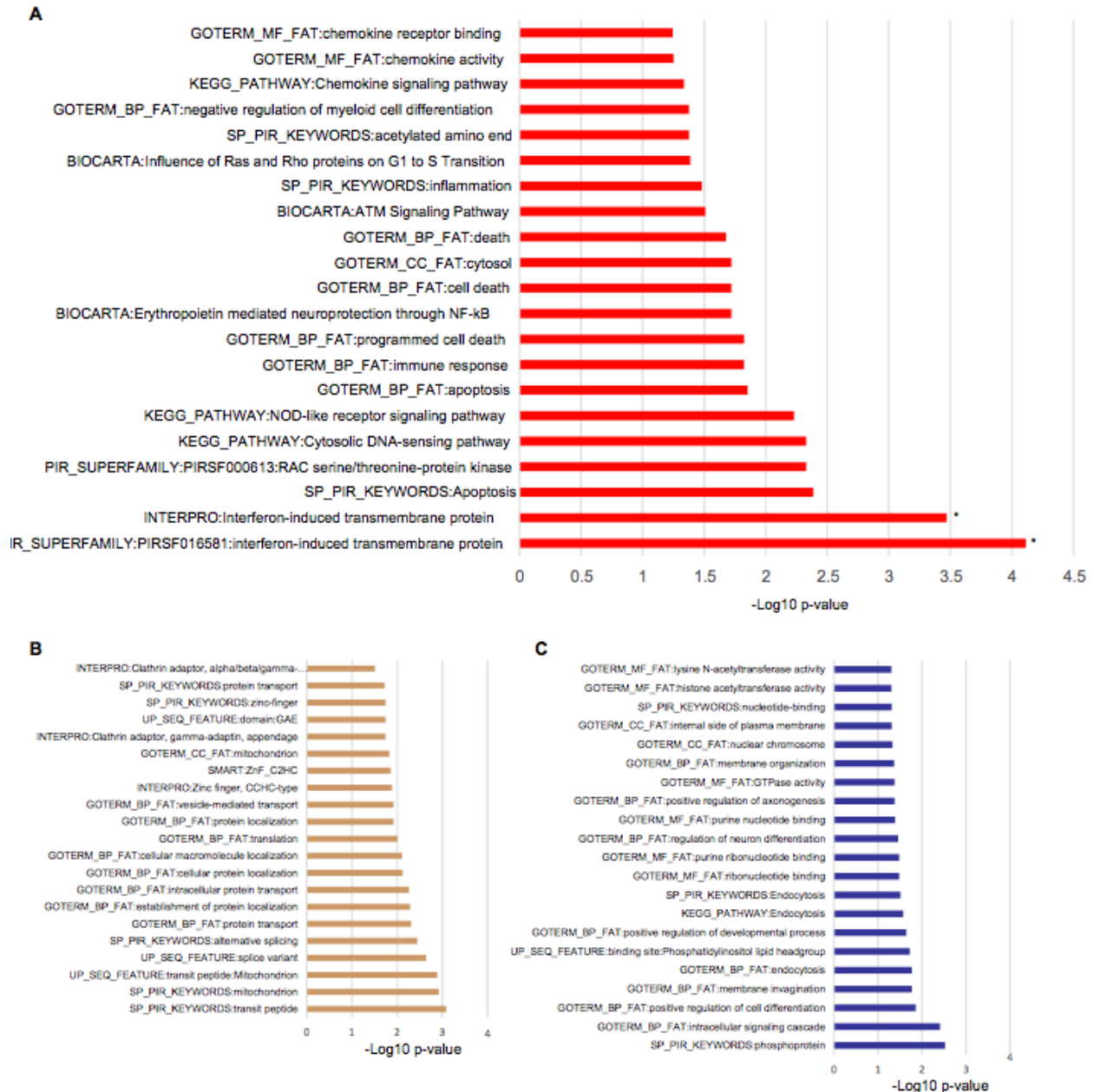
## 2.6 Supplementary Material

**Supplementary Table 2.1. Antibodies used for immunofluorescence and ISH.** *Target:* Protein or molecule targeted by the antibody. *Concentration:* Ratio of volume of antibody (at its provided concentration) to buffer for staining. *Species:* Species in which the antibody was generated.

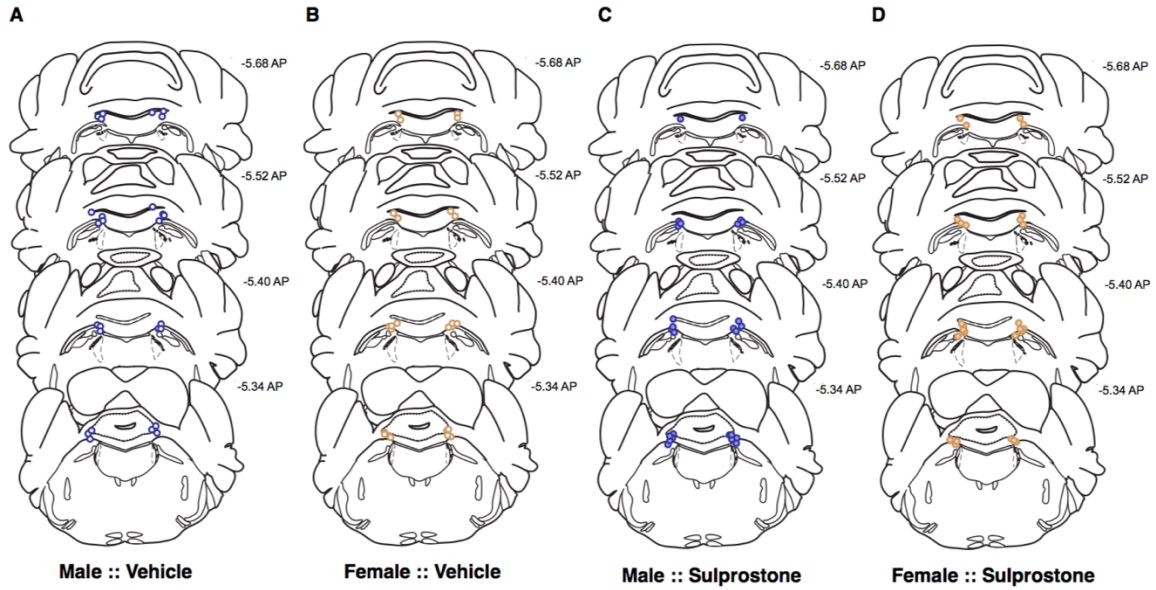
Target (Mouse)	Concentration	Species	Figure	Vendor	Product No.	Notes
TH	1 : 500	Mouse	2	MilliPore Sigma	MAB318	
DBH	1 : 500	Rabbit	2	Immunostar	22806	
CGRP	1 : 2,000	Rabbit	2	MilliPore Sigma	AB1971	
CELF6	1 : 500	Rabbit	2	N/A	N/A	Custom produced. Original citation: Dougherty, JD, et. al. "The Disruption of Celf6, a Gene Identified by Translational Profiling of Serotonergic Neurons, Results in Autism-Related Behaviors." <i>J Neurosci</i> , 2016. doi.org/10.1523/JNEUROSCI.4762-12.2013
NG2	1 : 100	Rabbit	2	Santa Cruz Biotechnology	SC-30144	Discontinued.
GPX3	1 : 50	Mouse	2	Santa Cruz Biotechnology	SC-58361	
RBP4	1 : 50	Rabbit	2	Genetex	EP3657	Now supplied by Abcam.
DLK1	1 : 50	Mouse	2	Santa Cruz Biotechnology	SC-80024	Discontinued.
Digoxigenin	1 : 7,500	Sheep	2	Roche	11093274910	Coupled to alkaline phosphatase (for <i>in situ</i> hybridization staining with NBT+BCIP).
Phospho-Ser32 c-Fos	1 : 500	Rabbit	N/A	Cell Signaling Technologies	5348S	Discussed in results section.
Rabbit IgG	1 : 1,000	Donkey	2	Invitrogen	A21206	Coupled to Alexa Fluor 488.
Rabbit IgG	1 : 1,000	Donkey	2	Invitrogen	A10040	Coupled to Alexa Fluor 546.
Mouse IgG	1 : 1,000	Donkey	2	Invitrogen	A21202	Coupled to Alexa Fluor 488.
Mouse IgG	1 : 1,000	Donkey	2	Invitrogen	A10036	Coupled to Alexa Fluor 546.



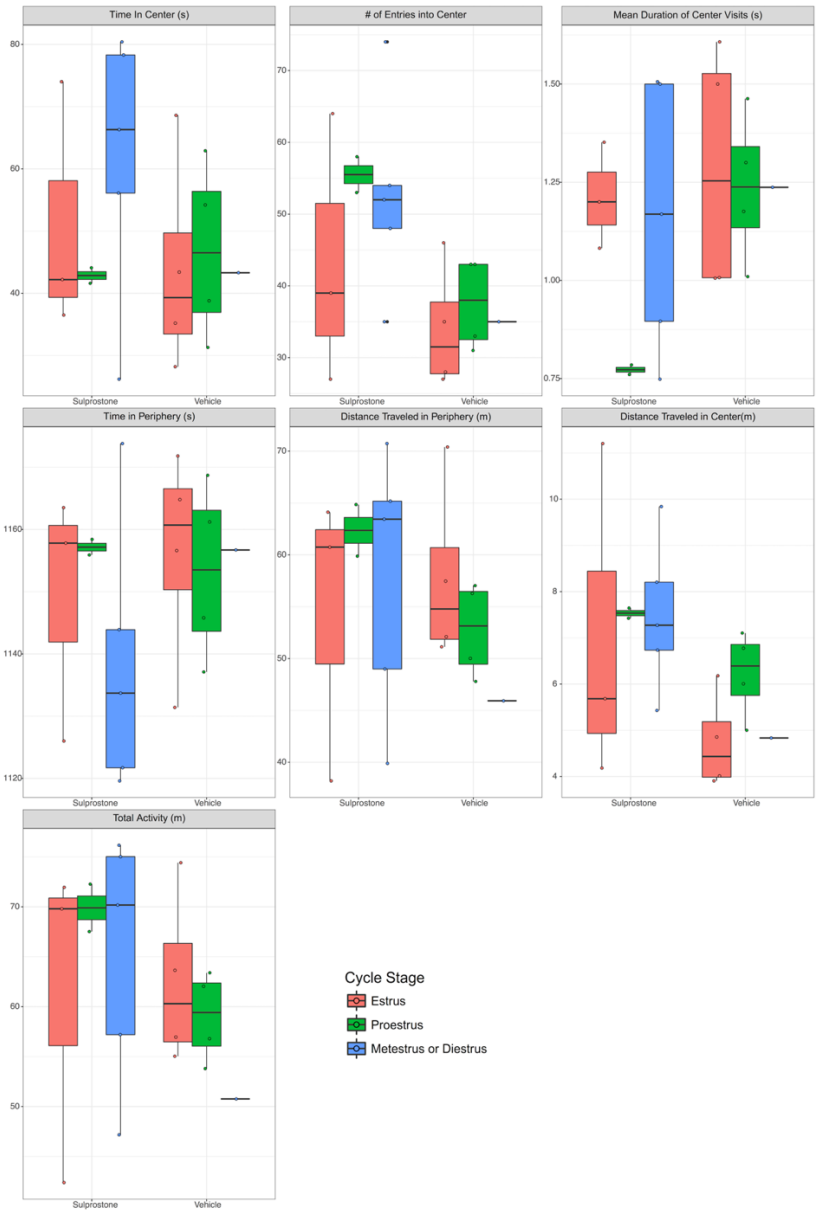
**Supplementary Figure 2.1. Examples of Allen Brain Atlas ISH scoring.** Images from Allen Brain Atlas for examples of *in situ* hybridizations that would score a 1-5, respectively.



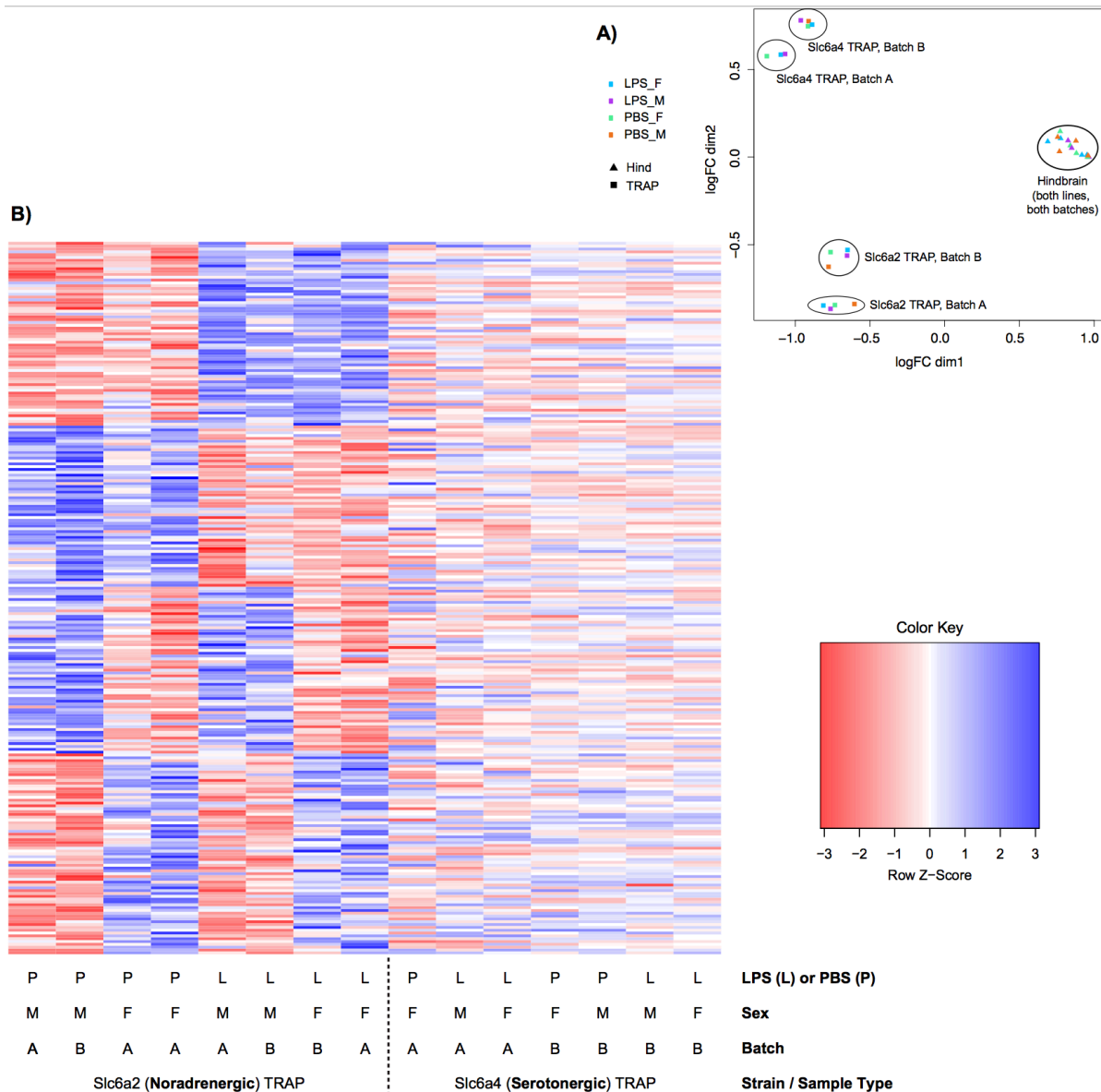
**Supplementary Figure 2.2. Pathway analysis of transcripts altered in hindbrain by LPS and of sex-differentially expressed transcripts in noradrenergic neurons. A)** A pathway analysis using DAVID reveals the hindbrain showed a significant increase of interferon related gene expression, and trends in a variety of chemokine and inflammatory pathways. **B-C)** An exploratory pathway analysis using DAVID illustrates the trends in male and female noradrenergic gene lists (-Log<sub>10</sub> p-values. No  $p < 0.05$  after Benjamini-Hochberg correction for multiple testing). Though no individual pathway survived correction for multiple testing, the trend towards more mitochondrial genes and vesicular transport transcripts in female neurons may suggest they have slightly higher metabolic demands. Male neurons had slightly more nuclear factors, driven by transcription factors such as *Atrx1*, *Mtf1* and *Pbx1*.



**Supplementary Figure 2.3. Map of LC cannula placements from post-mortem brain tissue.** Tissue from all mice cannulated and given sulprostone or vehicle was sliced, and the deepest point of the two cannula tracks was noted for each mouse. These points are collectively mapped here for mice across all conditions. If cannulae were not in the target area, the mice/tissue were excluded from behavioral analysis, c-FOS quantification, and this diagram. Anatomy is shown in the coronal plane with mm along anterior-posterior (AP) axis; dotted lines outline location of the LC. Circle outlines: blue = male, orange = female; fill color: purple = sulprostone-treated, empty/white = vehicle-treated.



**Supplementary Figure 2.4. OFT data by estrous stage from one cohort of sulprostone-behavior mice.** Boxplots for seven OFT measures in counts, meters (m), or seconds (s) from females in a cohort of sulprostone-behavior mice with estrus staging data. Metestrus and diestrus were combined per Silva, *et. al*<sup>41</sup>.



**Supplementary Figure 2.5. Multidimensional scaling (MDS) and heatmaps illustrate the clustering of samples and gene sets.** **A)** Hierarchical clustering confirms the expected segregation of samples by cell/sample type (*Slc6a2* vs. *Slc6a4* vs. hindbrain, or TRAP vs. hindbrain), and show a degree of batch effects, guiding the analytical strategies employed. **B)** Clustering genes that were DE between at least one comparison in the *Slc6a2* TRAP samples (i.e., comparing expression between sexes or between LPS/PBS conditions) illustrates the presence of distinct LPS- and sex-dependent DEG sets in noradrenergic neurons. The adjacent *Slc6a4* samples illustrate that sex and LPS do not affect these genes in serotonergic neurons in an analogous manner.



## 2.7 References

1. Kessler, R. C. *et al.* Lifetime Prevalence and Age-of-Onset Distributions of DSM-IV Disorders in the National Comorbidity Survey Replication. *Arch Gen Psychiat* 62, 593 (2005).
2. Christensen, D. L. *et al.* Prevalence and Characteristics of Autism Spectrum Disorder Among Children Aged 8 Years — Autism and Developmental Disabilities Monitoring Network, 11 Sites, United States, 2012. *Mmwr Surveillance Summ* 65, 1–23 (2016).
3. Fombonne, E. Epidemiology of pervasive developmental disorders. *Pediatr Res* 65, 591–8 (2009).
4. Sara, S. J. The locus coeruleus and noradrenergic modulation of cognition. *Nat Rev Neurosci* 10, 211–23 (2009).
5. Hare, A. S., Clarke, G. & Tolchard, S. Bacterial Lipopolysaccharide-Induced Changes in FOS Protein Expression in the Rat Brain: Correlation with Thermoregulatory Changes and Plasma Corticosterone. *J Neuroendocrinol* 7, 791–799 (1995).
6. Kurosawa, N., Shimizu, K. & Seki, K. The development of depression-like behavior is consolidated by IL-6-induced activation of locus coeruleus neurons and IL-1 $\beta$ -induced elevated leptin levels in mice. *Psychopharmacology* 233, 1725–37 (2016).
7. McCall, J. G. *et al.* CRH Engagement of the Locus Coeruleus Noradrenergic System Mediates Stress-Induced Anxiety. *Neuron* 87, 605–20 (2015).
8. Seo, D. & Bruchas, M. R. Polymorphic computation in locus coeruleus networks. *Nat Neurosci* 20, 1517–1519 (2017).
9. Bangasser, D. A., Zhang, X., Garachh, V., Hanhauser, E. & Valentino, R. J. Sexual dimorphism in locus coeruleus dendritic morphology: A structural basis for sex differences in emotional arousal. *Physiol Behav* 103, 342–351 (2011).
10. Valentino, R. J. & Bangasser, D. A. Sex-biased cellular signaling: molecular basis for sex differences in neuropsychiatric diseases. *Dialogues Clin Neurosci* 18, 385–393 (2016).
11. Curtis, A. L., Bethea, T. & Valentino, R. J. Sexually Dimorphic Responses of the Brain Norepinephrine System to Stress and Corticotropin-Releasing Factor. *Neuropsychopharmacol* 31, 544–554 (2006).
12. Bangasser, D. A. & Valentino, R. J. Sex Differences in Molecular and Cellular Substrates of Stress. *Cell Mol Neurobiol* 32, 709–723 (2012).
13. Bangasser, D. A., Wiersielis, K. R. & Khantsis, S. Sex differences in the locus coeruleus-norepinephrine system and its regulation by stress. *Brain Res* 1641, 177–88 (2016).
14. Guajardo, H. M., Snyder, K., Ho, A. & Valentino, R. J. Sex Differences in  $\mu$ -Opioid Receptor Regulation of the Rat Locus Coeruleus and Their Cognitive Consequences. *Neuropsychopharmacol* 42, 1295–1304 (2016).

15. Babstock, D., Malsbury, C. W. & Harley, C. W. The dorsal locus coeruleus is larger in male than in female Sprague–Dawley rats. *Neurosci Lett* 224, 157–160 (1997).
16. Luque, J. M., Blas, M. R. de, Segovia, S. & Guillamón, A. Sexual dimorphism of the dopamine- $\beta$ -hydroxylase-immunoreactive neurons in the rat locus ceruleus. *Dev Brain Res* 67, 211–215 (1992).
17. Schroeter, S. *et al.* Immunolocalization of the cocaine- and antidepressant-sensitive 1-norepinephrine transporter. *J Comp Neurology* 420, 211–232 (2000).
18. Doyle, J. P. *et al.* Application of a translational profiling approach for the comparative analysis of CNS cell types. *Cell* 135, 749–62 (2008).
19. Carter, M. E., Lecea, L. de & Adamantidis, A. Functional wiring of hypocretin and LC-NE neurons: implications for arousal. *Front Behav Neurosci* 7, 43 (2013).
20. Almeida, M. C., Steiner, A. A., Coimbra, N. C. & Branco, L. G. S. Thermoeffector neuronal pathways in fever: a study in rats showing a new role of the locus coeruleus. *J Physiology* 558, 283–294 (2004).
21. Valentino, R. J., Reyes, B., Bockstaele, E. V. & Bangasser, D. Molecular and cellular sex differences at the intersection of stress and arousal. *Neuropharmacology* 62, 13–20 (2012).
22. Darling, R. D. *et al.* Perinatal citalopram exposure selectively increases locus ceruleus circuit function in male rats. *J Neurosci Official J Soc Neurosci* 31, 16709–15 (2011).
23. Peña, C. J. *et al.* Early life stress confers lifelong stress susceptibility in mice via ventral tegmental area OTX2. *Sci New York N Y* 356, 1185–1188 (2017).
24. Arora, M. *et al.* Fetal and postnatal metal dysregulation in autism. *Nat Commun* 8, 15493 (2017).
25. Hagemeyer, S., Mangus, K., Boeckers, T. M. & Grabrucker, A. M. Effects of Trace Metal Profiles Characteristic for Autism on Synapses in Cultured Neurons. *Neural Plast* 2015, 1–16 (2015).
26. Amateau, S. K. & McCarthy, M. M. Induction of PGE2 by estradiol mediates developmental masculinization of sex behavior. *Nat Neurosci* 7, 643–650 (2004).
27. Wright, C. L., Burks, S. R. & McCarthy, M. M. Identification of prostaglandin E2 receptors mediating perinatal masculinization of adult sex behavior and neuroanatomical correlates. *Dev Neurobiol* 68, 1406–19 (2008).
28. Oka, T. *et al.* Relationship of EP1–4 prostaglandin receptors with rat hypothalamic cell groups involved in lipopolysaccharide fever responses. *J Comp Neurology* 428, 20–32 (2000).
29. McCall, J. G. *et al.* Locus coeruleus to basolateral amygdala noradrenergic projections promote anxiety-like behavior. *Elife* 6, e18247 (2017).
30. Uematsu, A. *et al.* Modular organization of the brainstem noradrenaline system coordinates opposing learning states. *Nat Neurosci* 20, 1602–1611 (2017).

31. Hyde, C. L. *et al.* Identification of 15 genetic loci associated with risk of major depression in individuals of European descent. *Nat Genet* 48, 1031–1036 (2016).
32. Serova, L., Rivkin, M., Nakashima, A. & Sabban, E. L. Estradiol Stimulates Gene Expression of Norepinephrine Biosynthetic Enzymes in Rat Locus coeruleus. *Neuroendocrinology* 75, 193–200 (2002).
33. Thanky, N. R., Son, J. H. & Herbison, A. E. Sex differences in the regulation of tyrosine hydroxylase gene transcription by estrogen in the locus coeruleus of TH9-LacZ transgenic mice. *Mol Brain Res* 104, 220–226 (2002).
34. Prendergast, B. J., Onishi, K. G. & Zucker, I. Female mice liberated for inclusion in neuroscience and biomedical research. *Neurosci Biobehav Rev* 40, 1–5 (2014).
35. Meziane, H., Ouagazzal, A.-M., Aubert, L., Wietrych, M. & Krezel, W. Estrous cycle effects on behavior of C57BL/6J and BALB/cByJ female mice: implications for phenotyping strategies. *Genes Brain Behav* 6, 192–200 (2007).
36. Guillamón, A., Blas, M. R. de & Segovia, S. Effects of sex steroids on the development of the locus coeruleus in the rat. *Dev Brain Res* 40, 306–310 (1988).
37. Bailey, T. L., Johnson, J., Grant, C. E. & Noble, W. S. The MEME Suite. *Nucleic Acids Res* 43, W39–W49 (2015).
38. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of transcription factor binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids Res* 46, D252–D259 (2018).
39. Mathelier, A. *et al.* JASPAR 2016: a major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* 44, D110-5 (2016).
40. Weirauch, M. T. *et al.* Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* 158, 1431–1443 (2014).
41. Silva, A. F. *et al.* Sex and estrous cycle influence diazepam effects on anxiety and memory: Possible role of progesterone. *Prog Neuro-psychopharmacology Biological Psychiatry* 70, 68–76 (2016).

## Chapter 3: From sex differences in gene expression to sex differences in expression regulation

Two findings in Chapter 2 form critical bridges from sex differences in the LC to the remaining work described in this dissertation. The first is the identification of conserved, recurrent sequences near sex-differentially expressed genes of the LC (**Figure 2.4**). Each mouse genomic sequence matching one of the twelve motifs were designed into an MPRA oligonucleotide library intended for *in vivo* testing to verify sex-differential regulatory function of those particular sequences.

The second finding, described in **Chapter 2.3** (discussion) highlights that two early genome-wide candidate genes for MDD were also among the <100 genes with female-biased, LC-enriched expression in the TRAP dataset. This was an exciting confluence of events that spurred me to additionally design an allelic MPRA oligonucleotide library for each human genomic sequence, in three sequence contexts relative to the SNP (with a tile such that the SNP was at the 5' end, center, or 3' end) in modest linkage disequilibrium with the GWAS tag SNPs corresponding to these two genes' (*SLC6A15*, *LIN28B*) loci. Given the sex differences in these genes' expression *and* candidacy as mediators of MDD genetic risk, these two loci seemed like strong candidates for identifying sex-interacting risk variants for the disorder.

Sadly, the oligonucleotide library purchased for these two experiments was irreparably ill-synthesized, a fact which took nearly a year to identify. In the meantime, regulatory genomics data from postmortem human brain and *in vitro* human cell types had grown abundant, especially thanks to the PsychENCODE consortium's first data tranche released in December of 2018. Simultaneously, the gold-standard psychiatric GWAS group, the Psychiatric Genomics

Consortium (PGC), had published new, better-powered, and well-phenotyped GWASes for several psychiatric disorders, including MDD, expanding the number of associated loci. Simultaneously, work leveraging UK Biobank data with PGC data identified new GWAS loci for depression and highly correlated traits, like neuroticism. Human genomics research now had a wealth of psychiatric GWAS statistics to work with, and far more brain-based -omics data to annotate these findings than had been available just a couple years prior.

In light of the newly available human neural multiomic data and rapidly increasing number of genetic associations for MDD and closely related traits, a second library spanning many additional GWAS loci was constructed, which forms the basis of the remainder of this work.

## **Chapter 4: Transcriptional-regulatory convergence across functional MDD risk variants identified by massively parallel reporter assays**

This chapter was previously published:

Mulvey, B. & Dougherty, J. D. Transcriptional-regulatory convergence across functional MDD risk variants identified by massively parallel reporter assays. *Transl Psychiat* 11, 403 (2021).

Family and population studies indicate clear heritability of major depressive disorder (MDD), though its underlying biology remains unclear. The majority of single-nucleotide polymorphism (SNP) linkage blocks associated with MDD by genome-wide association studies (GWASes) are believed to alter transcriptional regulators (e.g., enhancers, promoters) based on enrichment of marks correlated with these functions. A key to understanding MDD pathophysiology will be elucidation of which SNPs are functional and how such functional variants biologically converge to elicit the disease. Furthermore, retinoids can elicit MDD in patients and promote depressive-like behaviors in rodent models, acting via a regulatory system of retinoid receptor transcription factors (TFs). I therefore sought to simultaneously identify functional genetic variants and assess retinoid pathway regulation of MDD risk loci. Using Massively Parallel Reporter Assays (MPRAs), I functionally screened over 1 000 SNPs prioritized from 39 neuropsychiatric trait/disease GWAS loci, selecting SNPs based on overlap with predicted regulatory features—including expression quantitative trait loci (eQTL) and histone marks—from human brains and cell cultures. I identified >100 SNPs with allelic effects on expression in a retinoid-responsive model system. Functional SNPs were enriched for binding sequences of retinoic acid-receptive transcription factors (TFs), with additional allelic differences unmasked by treatment with all-*trans* retinoic acid (ATRA). Finally, motifs overrepresented across functional SNPs corresponded to TFs

highly specific to serotonergic neurons, suggesting an *in vivo* site of action. Our application of MPRA to screen MDD-associated SNPs suggests a shared transcriptional regulatory program across loci, a component of which is unmasked by retinoids.

## 4.1 Introduction

Major depressive disorder (MDD) is a common but debilitating psychiatric disorder affecting hundreds of millions worldwide<sup>1</sup>, exacting substantial tolls on both individuals<sup>2</sup> and societies<sup>3</sup>. Despite the global burden of MDD, nearly half of patients do not clinically respond to treatment<sup>4</sup>, in part due to limited understanding of its biological underpinnings. Family studies have demonstrated that MDD risk is at least 30% heritable<sup>5,6</sup>. More recently, genome-wide association studies (GWASes) have demonstrated similar estimates for severe MDD<sup>7</sup>, and have helped narrow in on associated single nucleotide polymorphisms (SNPs)<sup>8–11</sup>, a tantalizing entry point for understanding the biology of MDD. However, GWAS studies do not identify causal variants, but rather implicate wider co-inherited regions consisting of many SNPs in linkage disequilibrium (LD). Most disease-associated SNPs are found outside of protein-coding sequences, suggesting probable roles in transcriptional regulation (TR)<sup>12–15</sup>. Which linked SNPs have functional allelic impacts on TR—and how they act together across loci to result in disease—remains unresolved.

It is thought that undetected, small-effect SNPs acting across the genome—including conditional SNPs within GWAS-significant loci<sup>16</sup>—contribute to the substantial heritability *not* caught by GWAS-significant SNPs alone<sup>17</sup>. Early support for multiple linked variants underlying GWAS signals came from examination of cell line histone marks in loci from six autoimmune disorder GWASes ; all six showed enrichment of TR-suggestive marks at LD SNPs only in a pertinent cell

type (B lymphocytes). Strikingly, 65% of the loci with  $\geq 1$  SNP overlapping lymphocyte histone marks contained multiple SNP-mark pairs, and over half of these loci contained at least three such SNPs<sup>18</sup>. Altogether, these findings implied that GWAS regions likely affect several TR features via several linked variants, especially in relevant cell types. More recently, GWASes have identified what are now called “conditional SNPs” associated with MDD<sup>19</sup>. However, despite predictions of multiple TR SNPs within GWAS loci, functional demonstration of this phenomenon has been sparse to date. The largest functional TR assay of MDD-associated variants examined 34 SNPs using luciferase assays<sup>20</sup>, representing successful but low-throughput identification of functional MDD SNPs. However, in terms of broad linkage, these loci constitute well over 10 000 SNPs, which will ultimately require higher-throughput approaches.

Furthermore, how functional SNPs—even once identified—biologically result in disease remains unclear, given their individually small effects on risk. The polygenic<sup>21</sup> and omnigenic<sup>17</sup> models were conceived of to address these aspects of complex disease genetics, establishing a guiding principle for GWAS interpretation. In brief, these theories posit that consistent emergence of a specific phenotype via widespread genomic variation necessarily requires common biological endpoints of those variants’ effects. At the molecular level, these points of convergence could be either upstream (shared regulation across loci)<sup>22</sup> or downstream (common biological pathways across loci). For downstream analyses, myriad approaches have been developed to nominate gene targets of putative TR SNPs using proximity<sup>23</sup>, chromatin structure<sup>24–26</sup>, or expression quantitative trait loci (eQTLs)<sup>27–29</sup>, yielding gene sets tested for enrichment in biological pathways<sup>27,28,30</sup> and cell types<sup>31</sup>. However, no analogous approaches exist for identifying convergent upstream (*i.e.*, TR) molecular effects of genetic risk, in part because a prerequisite is defining the functional SNPs.



One possible point of upstream TR convergence of MDD risk variants is retinoic acid and related compounds (retinoids). Retinoids drive transcriptional responses via several retinoid-binding nuclear receptor transcription factors (TFs) and heterodimerizing partners<sup>32,33</sup>. Besides their critical role in neurodevelopment, including of depression-implicated limbic structures<sup>34</sup>, retinoids have been associated with MDD onset and suicidality by epidemiological studies of the retinoid agonist isotretinoin<sup>35</sup>. Moreover, thyroid hormone is often used as an adjunctive treatment in MDD, and thyroid receptor TR effects are frequently carried out cooperatively with RXR family retinoid receptors<sup>36</sup>. Additional evidence for retinoid pathway activity in the adult brain—and its overactivity as a risk factor for depression—comes from rodent pharmacology and genetic models. For example, knockdown of *Cyp26b1*—which metabolizes retinoids—in adult mouse anterior insula suppresses interest in social novelty by reducing spontaneous activity of excitatory neurons<sup>37</sup>. Likewise, depressive symptoms have been observed in rats after intracerebroventricular all-*trans* retinoic acid (ATRA) administration<sup>38</sup>. In addition, RARA is more abundant in CRH neurons of affective disorder hypothalami<sup>39</sup>, where it both upregulates corticotropin releasing hormone (*CRH*) expression and blocks glucocorticoid negative feedback on *CRH*<sup>40</sup>, suggesting a link between retinoid TFs and elevated hypothalamic-pituitary-adrenal axis activity in MDD. Finally, given the substantial shared genetic risk across psychiatric disorders<sup>41</sup>, it is notable that schizophrenia GWAS loci show enrichment for retinoid TR<sup>42</sup>, and that circulating retinoids are dysregulated in schizophrenia patients<sup>43</sup>. Similarly, retinoid pathway genes, including *CYP26B1*, are dysregulated in postmortem brain from autism spectrum disorders, bipolar disorder, and schizophrenia patients<sup>44</sup>. Interestingly, retinoid deficiencies have been associated with these diseases, including recent observations of reduced serum levels of retinoic acid and its precursor,

retinol, in schizophrenia<sup>43</sup>; similarly, reductions in serum retinol and expression of all three *RAR* genes were shown in autism spectrum disorders<sup>45</sup>. These findings led us to speculate that a component of MDD-associated genetic risk may likewise demonstrate an upstream convergence via recurrent retinoid-mediated TR disruptions across loci.

Massively parallel reporter assays (MPRAs) provide a solution to both experimentally identify functional variants and, consequently, their shared TR features. MPRAs assess thousands of DNA elements for transcriptional-regulatory functions and allelic differences simultaneously by pairing each short genomic sequence element of interest to several unique barcodes, with a constant promoter and reporter gene placed in between<sup>46-49</sup>. Delivery of a library of DNA elements to cells, followed by RNA collection and sequencing, enables quantitative estimation of the expression driven by each element as a ratio of expressed RNA barcode to delivered DNA barcode. These assays have recently been adapted to systematically identify SNPs with functional allelic TR differences from GWAS loci for several diseases<sup>50-58</sup>. Two key features make MPRAs advantageous for identifying both functional SNPs and their TR interactions. First, the assay is carried out via transfection and targeted RNA-sequencing, meaning it can be executed in unmodified cell lines appropriate to the application. Second, MPRAs can be conducted to define TR effects of experimental manipulations in these systems, such as drug administration<sup>59,60</sup>.

Therefore, I sought to experimentally identify functional TR SNPs from 39 GWAS loci associated with MDD, neuroticism, and broader psychiatric disease risk, with the hypothesis that functional SNPs converge at the level of retinoid-mediated TR. From broad linkage regions, I selected over 1 000 SNPs based on overlapping human brain and neural epigenomic signals suggestive of TR

activity. Critically, selection of neither the loci nor the SNPs was predicated on retinoid involvement, allowing for unbiased functional screening of a cross-section of MDD GWAS loci. To ensure I could detect SNPs subject to retinoid-mediated TR, I used neuroblastoma (N2a) cells, as they are strongly and rapidly retinoid-responsive<sup>61,62</sup>. My initial assay identified over 75 functional SNPs from 29 GWAS regions, confirming that GWAS loci contain several functional SNPs. I then examined whether these functional SNPs possessed shared upstream TR features—namely, transcription factor (TF) binding motifs. Remarkably, there was indeed enrichment of retinoic acid binding TFs among the MPRA-functional vs. -nonfunctional SNPs, supporting my hypothesis. To further characterize retinoid effects on TR at MDD-associated SNPs, I performed a second assay using all-*trans* retinoic acid (ATRA), known to stall division of N2a and other neuroblastoma cells by inducing neuronal-like differentiation<sup>61</sup>. First, I found that functional SNPs containing retinoid receptor motifs had increased magnitudes of effect in the presence of ATRA, consistent with bonafide retinoid receptor TR activity. More broadly, ATRA led to striking rearrangements of the baseline regulatory landscape, including altered magnitude and reversed direction of allelic effects. Additionally, it revealed new SNPs with allelic TR differences unmasked by ATRA treatment. Significant ATRA-allele interaction SNPs largely overlapped RXRA binding sites from chromatin immunoprecipitation (ChIP)-seq, as well as motifs of several known retinoid-induced TFs, indicating broad roles of both retinoid TFs and their downstream TR systems at functional MDD-associated SNPs.

Finally, I explored the cell type-specificity of TFs predicted to regulate my functionally identified SNPs. Strikingly, I found TFs highly specific to serotonin neurons were strongly enriched among those I predicted to be recurrently involved in retinoid-dependent SNP function. These findings

suggest that the broad transcriptional-regulatory systems engaged by retinoids—and as I illustrate, the genetic component of MDD risk *they* engage—may converge on serotonergic neurons. In summary, I identify MDD-associated functional SNPs with both baseline and ATRA-mediated allelic differences in TR, and these disproportionately show upstream convergent regulation by retinoid receptors and TFs they induce. This highlights a striking potential point of convergence between genetic risk loci and an environmental risk factor for MDD.

## 4.2 Methods

### 4.2.1 Identifying candidate psychiatric GWAS regulatory variants.

To prioritize putative regulatory variants from neuropsychiatric disease GWAS regions (predominantly MDD; **Figure 4.1A**), SNPs in linkage disequilibrium (LD) with GWAS tag variants at  $R^2 > 0.1$  were collected and intersected with histone modification, eQTL, Hi-C, and enhancer segmentation datasets from human postmortem tissue and neural lineage cell lines (see *Supplementary Methods*, **Figure 4.1B**). SNPs were manually selected based on diversity and density of annotation overlap within each locus (*Supplementary Methods*). As a negative control, I identified candidates from one additional locus associated with several anthropomorphic traits<sup>63</sup>, in a trait-blinded manner. Altogether, 1453 SNPs were selected. Final LD of selected SNPs was distributed similarly to starting SNPs (**Figure 4.1D**). To confirm that I could detect CNS-relevant regulatory SNPs, a positive control TR SNP functionally demonstrated in mouse retina and brain<sup>56</sup> was also included.

Human genomic sequence (hg19) tiles up to 126bp were taken centered on the 1454 candidate enhancer SNPs, each paired to ten unique 10bp barcode sequences per allele and ordered as an oligonucleotide (oligo) pool from Twist Bioscience (San Francisco, CA). Also in the pool were

110 “basal” barcodes (no human genomic sequence cloned upstream of the minimal promoter), such that the only variable sequence between reporter clones was the barcode itself. The oligos were PCR amplified, then cloned into plasmid (**Figure 4.1E**); subsequently, a reporter cassette containing a minimal promoter (*hsp68*) driving the dsRed reporter gene<sup>64</sup> and the untranslated “woodchuck” element (for RNA stabilization, to improve signal)<sup>65</sup> was cloned in.

#### **4.2.2 Massively parallel reporter assays**

N2A cells were grown in uncoated 6-well plates in medium consisting of 0.1  $\mu$ M vacuum-filtered DMEM with 10% Fetal Bovine Serum (2% fetal bovine serum for the ATRA assay, based on media conditions from the literature<sup>66,67</sup>). For transfection, cells were reverse transfected by plating in antibiotic-free medium onto pre-plated 400  $\mu$ L mixtures of 2.5  $\mu$ g plasmid with Lipofectamine 2000. In the first assay, n=6 replicate wells were transfected and co-prepared for sequencing. A power analysis of these results using the 25<sup>th</sup>, 50<sup>th</sup>, and 75<sup>th</sup> percentile standard deviation of sequence expression measurements indicated I were  $\geq 80\%$  powered to detect Bonferroni-corrected  $p < 0.1$  variant effects as low as  $\text{abs}(\log_2\text{FC})$  1.1 (*Supplementary Methods*). In the drug MPRA experiment, n=12 wells were transfected, harvested, and prepared for sequencing together, with n=6 ATRA-treated and n=6 vehicle-treated.

After transfection, cells incubated for 7 hours at 37°C and 5% CO<sub>2</sub>. Medium was replaced with the respective medium containing antibiotics, and in the second assay, a final concentration of 20  $\mu$ M ATRA dissolved in DMSO, or equivalent volume of vehicle (DMSO). Medium was not replaced before RNA collection in the first assay; in the second assay, it was refreshed every 24 hours. 72 hours after transfection, cells were collected and RNA extracted using the Zymo (Irvine, CA) Clean-and-Concentrator 5 kit per manufacturer instructions. Eluted RNA was treated with

Turbo DNA-free kit to remove any residual plasmid to prevent contaminating DNA reads during sequencing, and cleaned a second time using the Zymo kit as above.

### **4.2.3 Targeted sequencing of RNA and input plasmid**

Briefly, equal amounts of RNA (1 $\mu$ g) from each sample were prepared for sequencing by targeted cDNA synthesis using a primer against the distal 3'UTR of the reporter. These, along with input plasmid, were subjected to PCR, enzymatic digestion, ligation of Illumina sequencing adapters, and a final PCR to add sample indexes for sequencing. Enzymes, and size-selection cleanup steps used in this process are fully detailed in *Supplementary Methods*. No-reverse-transcriptase controls utilizing sample RNA were co-prepared for both experiments and did not generate detectable product, indicating sequencing amplicons generated from RNA samples were exclusively representative of RNA content. Samples were sequenced to an average depth of ~8 million reads (first assay) or ~20 million reads (second assay).

### **4.2.4 Analysis**

Allelic SNP effects on expression in the first assay and in single-condition analyses of the second assay were assessed by *t*-testing the element's expression of each allele across replicates. In the first assay, over 90% of SNPs had normally distributed expression values (Shapiro-Wilk test,  $p > 0.05$ ). Uncorrected *t*-test *p*-values and Mann-Whitney U test *p*-values were well-correlated for the 89 non-normally distributed SNPs (Pearson's  $r = 0.825$ ). Nonetheless, for *t*-test significant SNPs ( $p_{\text{emp}} < 0.05$ ) not passing the Shapiro-Wilk test, I verified the result by checking for a nominally significant Mann-Whitney U test at  $p < 0.05$ . No SNPs were excluded from analysis on this basis. dbSNP-assigned reference ("ref") and alternative ("alt") alleles for each SNP were used to define comparison direction (the difference of activity under the alt allele vs. the ref allele). For the first

MPRA and single-condition analysis of vehicle samples from the second assay,  $p$  values were adjusted using empirical p-value correction via simulated allelic comparisons between random subsets of “basal” barcodes (see *Supplementary Methods*) following an analogous procedure from a multiplex CRISPR study<sup>68</sup>, with significance defined as  $p_{\text{emp}} < 0.05$  unless specified otherwise. This ensures that a representative cross-section of expression variability driven by barcode sequences is accounted for when assessing TR differences. Single-condition analysis of ATRA samples utilized standard Benjamini-Hochberg FDR correction, as primary effects of interest in these samples were ascertained by linear modeling. For analysis of ATRA effects, I verified that variances were similar between the drug and vehicle conditions; indeed, the median barcode expression level standard deviation was 0.1216 in ATRA-treated and 0.1226 in vehicle-treated samples (with respective 25<sup>th</sup> and 75<sup>th</sup> %ile standard deviations also matched within 0.005 expression units). I calculated samplewise barcode-level expression values passing the “single-condition” filtering steps used for t-testing (*Supplementary Methods*) were fitted using a linear mixed model (LMM) requiring a minimum of 40% (96) of the 240 possible expression measurements per SNP. The LMM included a random term for replicate (to account for well-specific effects), expressed as: barcode expression  $\sim$  allele + drug + allele:drug + (1|replicate). Empirical q-value correction for LMM  $F$  statistics was performed in an analogous manner to the prior experiment, generating a vector of  $F$  statistics for each coefficient from 20 000 randomized basal-only comparisons. All SNPs with an interaction  $p_{\text{emp}} < 0.05$  also had a likelihood ratio test (LRT)  $p < 0.051$  comparing a maximum-likelihood (ML) interaction model to an ML LMM with additive terms only, indicating that the interactive model was more predictive but not overfit compared to an additive model. For SNPs with significant allele and interaction coefficients, a meaningful allele main effect was considered present if the single-condition vehicle and ATRA

analyses showed the same allelic direction of effect, with a vehicle  $p_{\text{emp}} < 0.1$  and ATRA FDR < 0.1 (*i.e.*, near-significant within each condition of  $n=6$ , thus reasonably capable of achieving significance in the LMM analysis of the two conditions combined).

#### **4.2.5 MotifbreakR analysis and functional SNP enrichment for perturbed motifs**

The motifbreakR<sup>69</sup> package was used to identify TF binding motifs significantly different between alleles of each SNP. Briefly, the number of MPRA-identified functional SNPs matching a given TF's motif(s) for at least one allele was compared to the number of non-functional SNPs matching across 10 000 random draws of  $n$  (number of significant) SNPs. A second version of this analysis focused on the concordance rate—that is, whether the frequency of functional variants experiencing concurrent strengthening of motif and expression or vice versa—was significant compared to 10 000 draws of  $n$  random SNPs from the analyzed set. Analysis of the first assay's SNPs defined functionality based on a  $p_{\text{emp}}$  value threshold of 0.05. I performed two motif analyses of the second assay results, one comparing allele main effect SNPs ( $p_{\text{emp}} < 0.1$ ) to those with  $p_{\text{emp}} > 0.1$  for allele, drug, and interaction effects, representing the breadth of functional variant-susceptible *cis*-regulators. The second analysis compared interaction SNPs ( $p_{\text{emp}} < 0.05$ ) to SNPs with an allele main effect (allele  $p_{\text{emp}} < 0.1$ ) but no interaction (interaction  $p_{\text{emp}} > 0.1$ ).

#### **4.2.6 Analyses of functional SNP-enriched TF expression in human brain and Chromatin Immunoprecipitation (ChIP)-seq**

I utilized outside ChIP-seq datasets to validate motif-based predictions of retinoid receptor binding and refine prediction of involved TFs. I intersected my functional SNPs to 25 tracks of ChIP-seq for retinoid receptors (19 human<sup>70,71</sup>, 6 mouse (3 ATRA treated, 3 vehicle treated) converted to



hg19 coordinates using UCSC's LiftOver<sup>72</sup>); 11 tracks of RXR heterodimerization partners (10 human THRA/THRB<sup>70,71</sup> and one aggregate analysis of human VDR<sup>73</sup>); and human genome-wide predictions of DR5<sup>74</sup>, a canonical RAR•RXR heterodimer binding sequence. For functional SNPs implicated at an RAR, RXR, VDR, or THRA/B site by either motifbreakR or CHIP, I identified potential target genes using chromatin-conformation<sup>75</sup> and eQTL<sup>76–79</sup> data. I performed broad-scope gene enrichment analysis of this gene set using Enrichr<sup>80</sup>. To examine shared biology of TFs implicated by motifbreakR enrichment at functional variants, I utilized PantherDb<sup>81</sup>. I finally examined TFs for enrichment among highly-expressed genes in adult and developing human brain using the ABAEnrichment package's Wilcoxon approach<sup>82</sup>, effectively weighting TFs by the number of functional SNPs implicated by motifbreakR (see *Supplementary Methods*).

#### **4.2.7 Code availability**

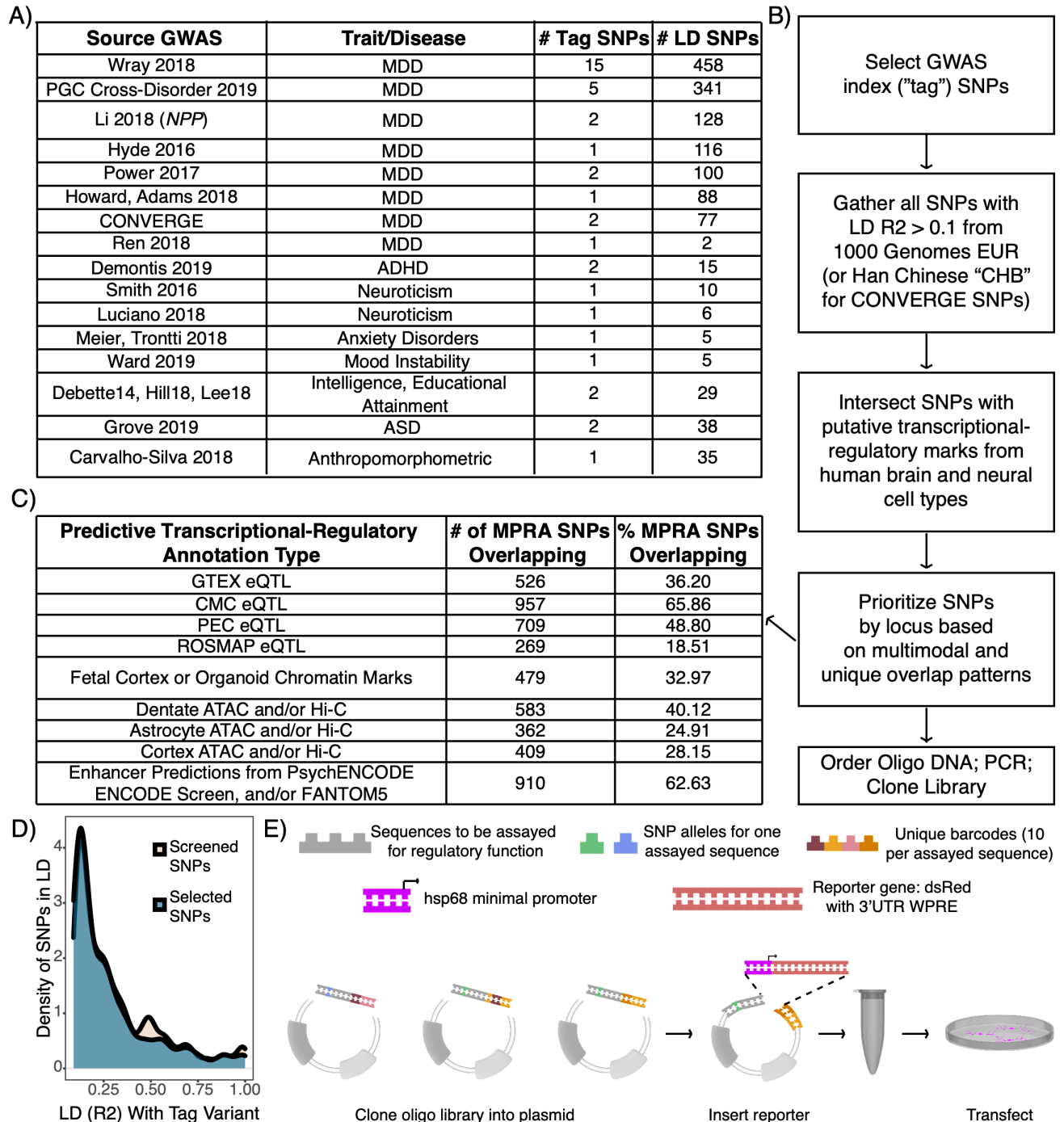
A summary spreadsheet of all significant SNPs identified in one or both assays, along with full analysis results, including barcode-wise expression in each sample, single-condition allelic effect tests, linear modeling results, and significantly enriched TFs in each of the comparisons executed, along with the code utilized to execute these analyses, is available at [https://bitbucket.org/jdlabteam/n2a\\_atra\\_mdd\\_mpra\\_paper/src](https://bitbucket.org/jdlabteam/n2a_atra_mdd_mpra_paper/src).

### **4.3 Results**

#### **4.3.1 Many MDD loci contain more than one functional SNP**

I identified >1 000 SNPs from MDD-associated GWAS loci, prioritizing SNPs overlapping with epigenetic data from neural samples, and cloned them into an MPRA library (**Figure 4.1**). I included one positive control SNP, shown to alter neural tissue gene expression, and one control locus near *CDKALI* not *a priori* associated with psychiatric disease. To identify functional variants

from these SNPs, the library was transfected into N2a cells (n=6 replicates, **Figure 4.1E**). Variant activity was assessed by RNA sequencing and barcode counts compared to input plasmid barcode counts. After filtering for read depth and barcode representation, 1013 SNPs spanning all 40 LD regions remained for analysis. Results were highly replicable across samples (Pearson  $r$  0.63-0.85 for barcode expression; 0.90-0.96 for elements, **Supplementary Figure 4.1**). I use “element” to signify the set of barcodes corresponding to one unique sequence of interest (1 SNP = 2 elements).



**Figure 4.1. Design of an MPRA library to identify candidate functional SNPs in MDD loci.**

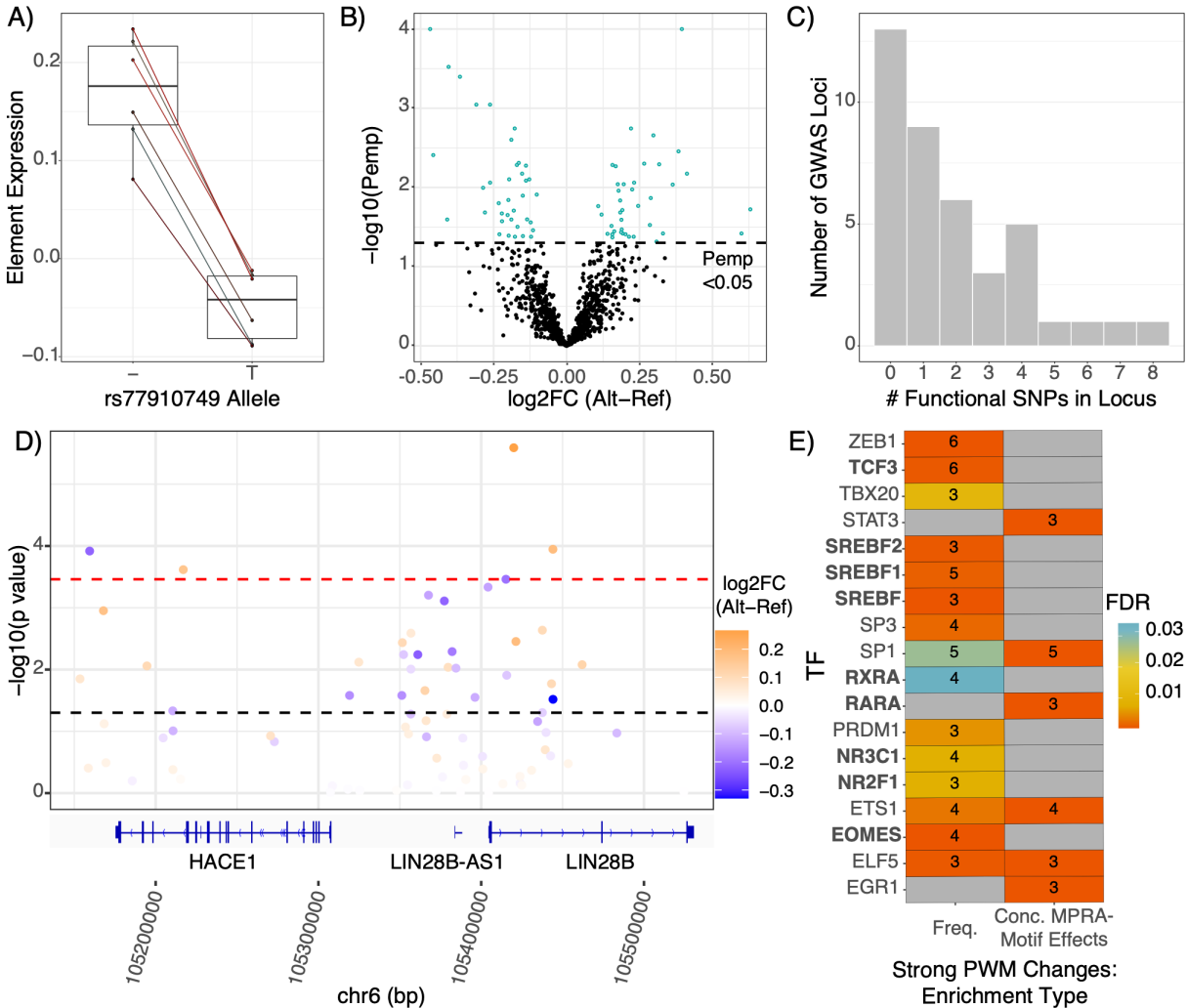
**A)** Table of GWAS studies and number of loci covered in the MPRA library. **B)** Flow chart of design and prioritization process. **C)** Brain and neural transcriptional-regulatory predictive annotation overlap with SNPs included in MPRA library. Fraction and number of SNPs in designed MPRA library intersecting each transcriptional-regulatory predictive annotation type. **D)** The manual prioritization process was not LD biased. The subset of prioritized SNPs are spread over LD space similarly to the full set of screened SNPs. **E)** Schematic of library construction and delivery. Panel adapted from<sup>49</sup>.

Of 1013 SNPs analyzed, I identified significant allelic TR ( $p_{\text{emp}} < 0.05$ ) at 76 SNPs (65 from MDD loci; 1 from the control *CDKALI* locus) across 27 of the 40 analyzed GWAS regions, with effects ranging 0.1 to 0.63 (median 0.2) log<sub>2</sub> fold-change (**Figure 4.2B**). Interestingly, the functional variant from the control locus is suggestively associated (GWAS  $p < 5 \cdot 10^{-6}$ ) with “*Poisoning by analgesics, antipyretics, and antirheumatics*” in UK Biobank<sup>83</sup>. As this likely includes attempted suicides, the SNP was retained for analyses. The positive control SNP, which I utilized to confirm my ability to detect small effect sizes expected of regulatory SNPs, showed the expected lower expression of the T allele at a  $p_{\text{emp}}$  of  $< 0.051$  (**Figure 4.2A**)<sup>56</sup>.

While my assay was designed to broadly examine wide LD regions around GWAS index variants, I did identify one functional variant, rs11209952, in a fine-mapped credible set of variants for seeking general practitioner care for depression in UK Biobank<sup>84</sup>. Moreover, consistent with prior studies of “conditional” or “secondary” SNP associations—wherein additional LD SNPs have associations independent of their linked, larger-effect variant<sup>19,85,86</sup>—I identified several loci with multiple functional SNPs (**Figure 4.2C**) (range 1-8, mean 2.8, median 2). Notably, I identified as functional rs1806153, a recently defined “conditional SNP” for MDD<sup>19</sup>. My findings support models predicting multiple functional SNPs in GWAS loci, and directly validate one such finding from association analysis.

One notable TR SNP I identified, rs314267, comes from a “*LIN28B*” (nearest gene) GWAS locus repeatedly linked to MDD<sup>8,87</sup> as well as cross-psychiatric disorder risk<sup>41</sup>. MPRA significance and effect size are illustrated for the region, showing that this locus contains several functional SNPs (**Figure 4.2D**). All significant MPRA SNPs in the locus had effect directions consistent with brain

eQTLs. rs314267 is the most significant *LIN28B* eQTL SNP (eSNP) in the region in PsychENCODE<sup>77</sup>, and is a CommonMind Consortium (CMC) eSNP for both *LIN28B* and *HACE1*<sup>76</sup>. *HACE1* is also downregulated in postmortem MDD hippocampal CA1<sup>88</sup>. Hi-C data from human neural cell cultures suggest rs314267 is within a neuron-specific *LIN28B* regulator, with promoter chromatin contacts found in dentate and cortical neurons, but not astrocytes<sup>75</sup>. *LIN28B* plays broad roles in neurodevelopment<sup>89</sup> and has potentially sex-differentiated functions<sup>90-93</sup>; considering sex differences in MDD prevalence and severity<sup>94,95</sup>, *LIN28B* constitutes an especially interesting gene target from this locus. Finally, I examined potential upstream TR mechanisms for SNP activity using VARAdb<sup>96</sup>. Query of rs314267 revealed a two order of magnitude allelic difference in the motif match p-value for *TCF4*—a gene itself implicated in cross-psychiatric-disorder risk<sup>41,97</sup>. Overall, the identification of functional SNPs implicated in regulation of *HACE1* and *LIN28B* exemplifies the ability of MPRA to identify functional variants involving sensible TR mechanisms and target genes.



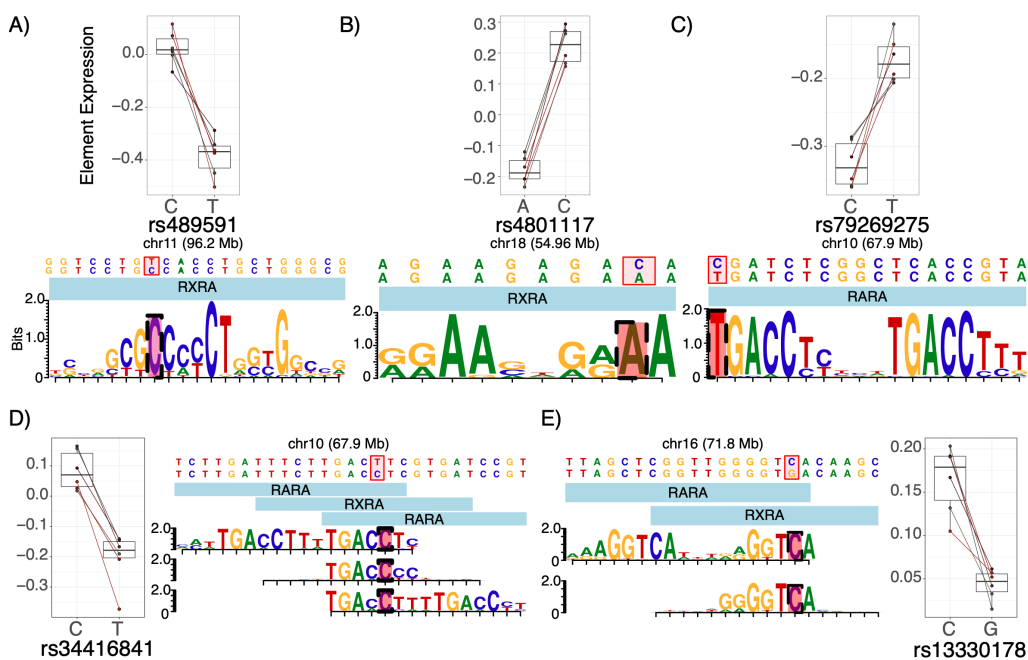
**Figure 4.2. MPRA defines SNPs with a functional effect on gene expression.** **A)** MPRA results of positive control SNP. Shen, et. al. found that the T allele drove decreased expression relative to the deletion (“-”) allele, which was robustly reproduced in the present assay. **B)** Volcano plot of allelic differences in reporter expression. Points represent one SNP’s composite  $\log_2$  allelic fold change (alt vs. ref), determined as the mean of samplewise alternative allele barcode expression minus the matched mean of reference allele barcodes. The dotted line indicates the statistically corrected significance threshold. **C)** Number of functional SNPs (MPRA significant SNPs) per GWAS locus in the assay. Number of loci (y-axis) containing a given number of MPRA-significant ( $p_{emp} < 0.05$ ) SNPs (x-axis). **D)** The LIN28B locus harbors several functional SNPs. SNPs are plotted according to their chromosomal position (hg19) and colored based on their composite  $\log_2$  allelic fold change. Refseq genes are visualized by the Integrative Genomics Viewer<sup>98</sup>. **E)** TF binding motifs involved in retinoid signaling, steroid synthesis and response, and neural activity are enriched among functional SNPs. Boxes are colored by FDR-corrected significance of enrichment for motifbreakR-defined “strong” allelic perturbations to binding motifs among functional SNPs; the number of functional SNPs perturbing (left column) and/or with concordant motif and MPRA effects (right column) are shown. Concordant effects were defined by greater MPRA expression driven by the allele better-matched to the corresponding TF motif and vice versa—the expected behavior of strictly activating TFs.

### 4.3.2 Shared regulatory architecture across distinct loci

I next sought to test my hypothesis that functional MDD risk variants shared retinoic acid-related TR architecture. If so, functional SNPs should disproportionately disrupt binding sites of retinoid-binding TFs compared to SNPs without an allelic effect on TR. Such data would indicate that MDD risk is mediated in part through perturbations of specific upstream transcriptional circuits and may highlight how risk conferred through retinoids converges with risk conferred through genetics to perturb downstream gene expression.

To take an unbiased approach to my retinoid hypothesis, I broadly analyzed all motifs showing enrichment at TR SNPs. Motifs for several dozen TFs were perturbed by the functional SNPs more frequently than expected, often with ‘strong’ perturbations to motifs and/or overrepresentation of concordant expression effects (**Figure 4.2E, Supplementary Figure 4.2**). This included several TFs aligned with biological processes relevant to psychiatric disease. For example, several TFs are involved in steroid pathways, from regulating biogenesis (*SREBF* family, 6 SNPs) to conveying downstream TR effects—most notably, via glucocorticoid receptor (*NR3C1*, 5 SNPs; **Supplementary Figure 4.3**), a central component of the stress response. Functional SNP overrepresentation of *SREBF* motifs is consistent with high expression of these TFs in N2as and related neuroblastomas<sup>99,100</sup>. A second group of transcription factors included three TFs involved in neural lineage commitment/development: *TCF3*<sup>101,102</sup>, *EOMES*, and *NR2F1*<sup>103</sup> (6,4, and 3 SNPs, respectively). Altogether, functional SNP enrichment for these TFs’ motifs bolster my confidence in this approach, as a) detected variation involves TFs known to be expressed in N2As (*SREBF*); b) functional variation involves TFs with roles in developing CNS, where disease variants likely act; and c) that the single-best characterized trigger of MDD (stress) is reflected in enrichment of alterations to *NR3C1* motifs.

Finally, consistent with my hypothesis of convergence on retinoid-mediated TR, functional variants were enriched for “strong” perturbations of retinoid receptor TF motifs (**Figure 4.3**), including RARA, RARB, and RXRA (5 SNPs from 4 MDD loci, **Figure 4.3**). Especially notable is the motif configuration at SNP rs34416841, which falls within three partially overlapping motifs for retinoid TFs. In addition, the elements overlapping rs489591 and rs13330178 appear to be functional human retinoid TF binding sites *in vivo* based on DNase hypersensitivity footprinting<sup>104</sup>.



**Figure 4.3. MPRA signal at SNPs disrupting putative retinoid TF motifs.** The five functional SNPs driving enrichment signals for retinoid TF motif perturbations are shown, along with the MPRA results for each variant. Each motif diagram shows only distinct position-weight matrices (PWMs). (motifbreakR uses a large meta-collection of motifs, which were often identical or nearly identical across retinoid TFs; such redundant motifs are not shown). Among functional-SNP enriched retinoid TFs, **A)** rs489591 and **B)** rs4801117 exclusively perturb RXRA motifs. **C)** rs79269275 perturbs an RARA (or RARB, identical but not shown) motif. **D)** rs34416841 alters several similar retinoid motifs across multiple positions and TFs. Not shown: near-identical motifs for RARB along the same sequence as the 5' RARA motif; near-identical RXRB and RARB motifs along the same sequence as the center RXRA motif; and a near-identical RARB motif along the same sequence as the 3' RARA motif. **E)** rs13330178 disrupts an RARA or RXRA binding site. Given that RXRA and RARs are known to heterodimerize, it is possible that this SNP disrupts the RXRA component of such a heteromer's binding sequence.



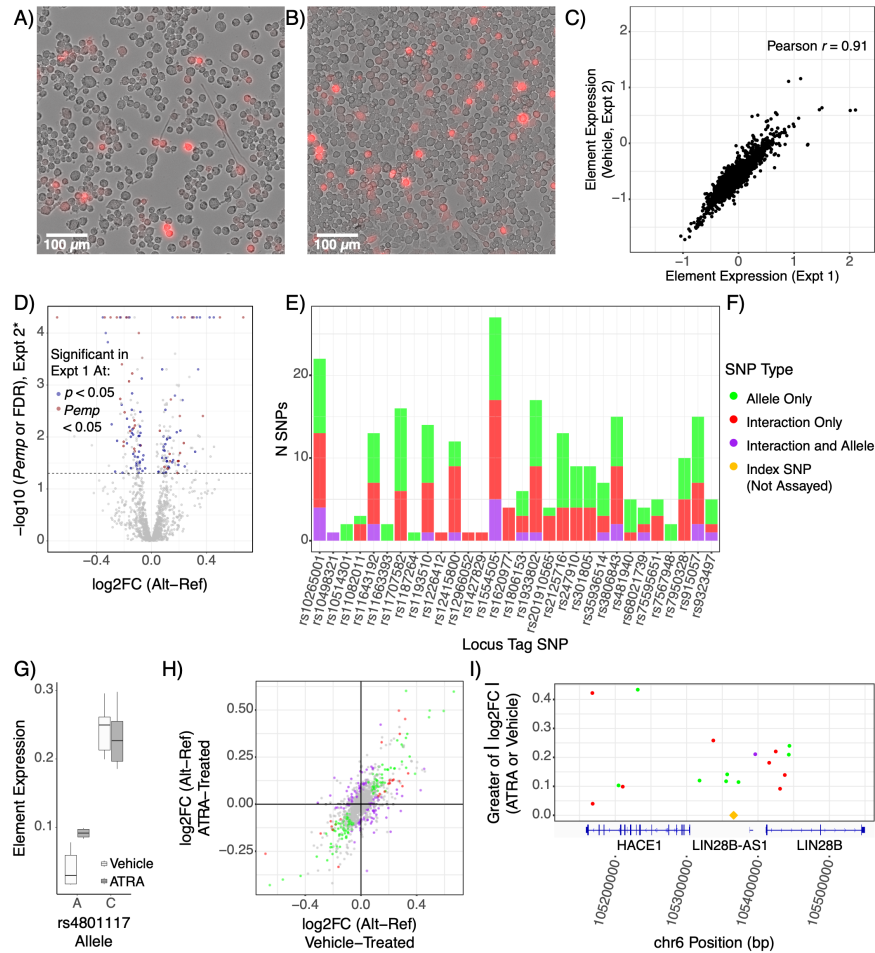
### 4.3.3 Retinoids unmask additional functional SNPs in MDD loci.

My findings supported the hypothesis that MDD-associated variants across multiple loci converge on TR, including that modulated by retinoids. I thus designed a pharmacological follow-up with two goals in mind. The first goal was to functionally verify that retinoids were involved in TR at SNPs where their motifs were found (in *cis*), and potentially unmask additional retinoid-targeted alleles. My second goal was to further assess retinoid signaling *trans* (i.e., indirect) effects on variants from these same GWAS regions, *e.g.*, via non-retinoid TF induction, co-regulation, or repression<sup>32</sup>. Therefore, I performed a second MPRA with an all-*trans* retinoic acid (ATRA) condition.

After 48 hours, cultures were imaged to verify drug activity (as ATRA is light-sensitive) based on known morphologic responses of N2as to ATRA, which include neurite outgrowth and mitotic arrest<sup>61,105,106</sup>. Indeed, drug-treated cells had a qualitatively lower cell density and produced neurite-like processes (**Figure 4.4A**) in comparison to vehicle-treated cells (**Figure 4.4B**). After RNA sequencing, I first analyzed vehicle-treated replicates alone to ensure replicability of the assay. Element expression levels in the vehicle condition strongly correlated to the first experiment (Pearson  $r = 0.91$ ; **Figure 4.4C**), and replicated the functional variants (**Figure 4.4D**); all 31 shared significant SNPs showed consistent directions of effect.

I next applied a linear mixed model (LMM) to identify SNPs responding to ATRA (that is, allele-drug interactions). A total of 1079 SNPs were analyzed after filtering for read and barcode depth. In part due to the effective doubling in power to detect allelic effects with 12 replicates and the LMM approach, I now identified 137 variants with a main effect of allele (129 from MDD loci). Four of the five retinoid receptor motif-perturbing variants from the first assay passed filtering; all

four of these variants again showed allelic main effects (all  $p_{\text{emp}} < 0.01$ ), as did many other functional variants identified in the previous experiment (**Figure 4.4D**). To my surprise, more variants showed a significant drug-allele interaction effect: a total of 128 SNPs (122 from MDD loci) (**Figure 4.4E-F**). Among the drug-allele interaction SNPs were one of the four retinoid-related SNPs identified from the first assay (rs4801117; interaction  $p_{\text{emp}} < 0.025$ , **Figure 4.4G**), while another trended towards interaction (rs489591; interaction  $p_{\text{emp}} = 0.117$ ). This strongly supports a role of retinoid TF activity at rs4801117 as predicted by the motif analysis. More broadly, comparison of changes between the two conditions reveals the striking extent to which the regulatory landscape of the N2As was altered by ATRA (**Figure 4.4F and 4.4H**). Notably, several additional functional variants were identified in the previously highlighted *LIN28B* locus, further illustrating multi-variant and context-dependent aspects of GWAS loci (**Figure 4.4E and 4.4I**). In all, this experiment highlights the ability of MPRA to detect contextual influences such as cell states and signaling on functional noncoding variation, and to unmask distinct, context-dependent functional SNPs.



**Figure 4.4. Retinoid treatment alters transcriptional regulation and unmasks additional functional variants.** **A)** ATRA-treated cells show growth arrest and neurite growth, demonstrating effective ATRA treatment, while **B)** Vehicle-treated cells continued to proliferate in a de-differentiated state. **C)** Results of the vehicle treatment replicate the initial MPRA findings. Element (single-allele) expression values for each sequence assessed in both assays is plotted. **D)** Significant and marginally significant functional SNPs from the first assay showed effects in the second assay. The larger allelic difference value from the ATRA and vehicle single-condition analyses is plotted for each SNP on the x-axis; the y-axis value is the corresponding corrected p-value (FDR correction for the ATRA-only analysis or empirical p-value correction for vehicle-only analysis). **E)** Retinoids unmask functional SNPs with additional or exclusive retinoid-mediated effects. **F)** SNP effect(s) color key for panels E, H, and I. SNPs with both effects were those with LMM interaction  $<0.05$ , LMM allele  $p_{emp} < 0.05$ , and both single-condition analyses showing the same allelic effect directionality at ATRA  $FDR < 0.1$  and vehicle  $p_{emp} < 0.1$ . **G)** rs4801117-A shows greater activity with ATRA treatment while the C allele is unaffected. The ATRA having an expression effect only on the A allele is consistent with the A allele matching the RXRA motif as shown in Figure 3B. **H)** Transcriptional-regulatory SNPs show a wide range of altered and unaltered effects with ATRA treatment. Single-condition  $\log_2FC$  values are shown. **I)** Several additional SNPs with retinoid-dependent function (i.e., allele-ATRA interaction) in the *LIN28B* locus. Only significant SNPs are illustrated. Notably, there are several functional SNPs clustered around the GWAS index SNP, suggesting association signal at this locus may be driven by multiple functional/conditionally functional variants.

#### **4.3.4 Retinoids reveal additional axes of convergent regulation at functional MDD-associated SNPs at levels of TF and cell type**

As the ATRA-based assay provided improved power to identify allelic variant effects on expression, I again employed my motifbreakR-based analyses to assess convergent transcriptional mechanisms underlying identified regulatory variants.. When examining SNPs with allelic effects in comparison to SNPs with no allelic, drug, or interaction effects, several retinoid receptor motifs were again overrepresented, including those of RXRA, RXRB, RARA, and RARG (**Figure 4.5A, Supplementary Figure 4.4**), totaling 11 of the 92 allele-main effect SNPs analyzed, spanning 10 MDD GWAS loci. These findings further support retinoid receptor binding sites as an upstream regulatory system recurrently involved in MDD risk genetics.

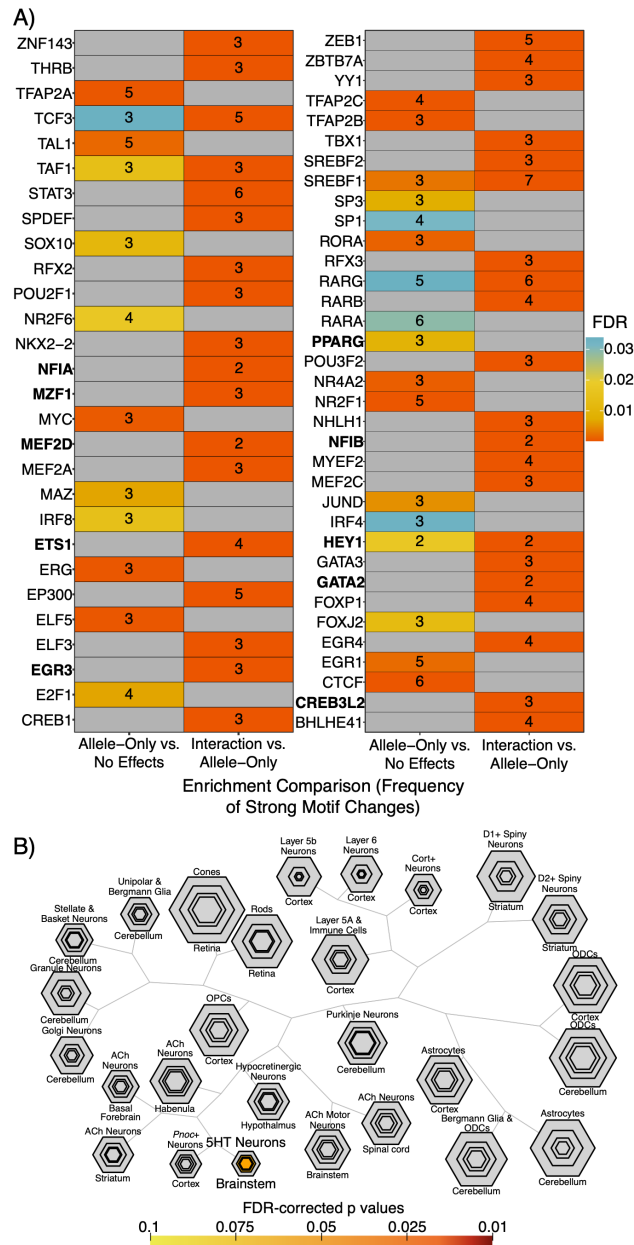
As retinoids resulted in stark changes across the transcriptional-regulatory landscape, I further sought to predict TFs potentially underlying allelic effects following retinoid exposure. Therefore, I also analyzed the interaction SNPs in comparison to allelic SNPs that were *not* subject to interactions. This revealed a novel set of TFs not observed in the preceding analyses, including TFs with roles in neural differentiation and maturation (**Figure 4.5A, Supplementary Figure 4.4**), as well thyroid hormone receptor THRB, an RXR binding partner. I compared the overrepresented motifs to TFs recently demonstrated to be upregulated in human neuroblastoma lines (KCNR, LAN5) by ATRA. Of the 26 TFs identified as ATRA-induced in these two lines, motifs were available for 18 in my analysis. Of these, 6 of the TFs were enriched among allele main-effect-only SNPs, while 12 of these TFs were enriched among the retinoid-allele interaction variants<sup>105</sup> (**Figure 4.5A, Supplementary Figure 4.4**), supporting my predictions of TFs playing ATRA-dependent roles at functional SNPs.

#### **4.3.5 Integrative analysis of TF sets at functional variants: TF binding, spatiotemporal brain enrichment, and putative target genes**

As retinoid receptors have highly redundant binding motifs, I sought to both validate motif-based implication of retinoid receptors and more finely identify the particular TFs binding at functional SNPs. I aggregated ChIP-seq data for RAR, RXR, and RXR-heterodimerization partners (VDR, THRA, THRB) and identified functional SNPs overlapping peaks for each TF. Altogether, 35 of my 277 functional SNPs from across the two assays were in at least one such binding site (**Supplementary Table 4.2**). 15/17 of the allele•ATRA interaction SNPs overlapped a ChIP peak for RXRA, suggesting RXRA may be the common mediator of the observed retinoid-dependent SNP effects.

I also performed Gene Ontology analysis of functional-variant enriched TFs against a background of all TFs in the motifbreakR tool using PANTHERdb but found no detailed Biological Processes of note. I next sought to examine whether TFs enriched at functional variants in my motif analyses corresponded to particular spatiotemporal expression patterns in the brain. To favor the most broadly-implicated TFs, I utilized the ABAEnrichment package's Wilcoxon analysis approach on the TF sets from the ATRA experiment using the number of motifbreakR SNPs as the TF gene "scores". In this analysis, several brain regions across developmental stages were nominally enriched (family-wide error rates < 0.05) in ATRA-dependent and -independent TF expression, with especially broad enrichment at high expression thresholds ( $\geq$  90th percentile) in adolescent brain. This does not appear to be an artifact of the cell model, considering that neuroblastomas are arrested in a neural crest progenitor (i.e., pre-/peri- natal cell type) stage. If replicated in future studies with larger adolescent sample numbers, this may suggest that retinoid-mediated aspects of

MDD genetic risk are especially active in the adolescent brain, perhaps contributing to frequent emergence of the disorder around this time.



**Figure 4.5. Distinct TFs underlying retinoid-dependent functional SNPs and implication of serotonergic neurons. A)** Motifs overrepresented among ATRA-independent (allele effect without interaction) SNPs (left columns) or among ATRA-dependent (interaction) SNPs (right column). The heatmap is shown in halves for visibility. TFs identified as ATRA-upregulated in human neuroblastoma lines<sup>105</sup> are in bold font. **B)** TFs implicated by ATRA-interacting SNPs significantly overlap TFs enriched in serotonergic neurons. Plot generated using the cell-specific expression analysis (CSEA) tool (<http://genetics.wustl.edu/jdlab/csea-tool-2/>)<sup>31</sup>. 5HT: Serotonin; ACh: Acetylcholine; ODC: Oligodendrocyte; OPC: Oligodendrocyte Progenitor.

I additionally utilized these sets of TFs as gene sets to investigate whether retinoid-dependent or -independent regulatory variants might be particularly active in certain cell types of the brain. I screened for enrichment of these TFs among genes with strong cell type-specific expression in brain as previously defined for over 20 cell type translomes<sup>31</sup>. Three TFs (spanning 8 ATRA-interacting SNPs) were discovered to be highly specific to serotonin neurons (**Figure 4.5B**): *GATA2*, *GATA3*, and *FEV*, while no cell type enrichments under FDR<0.1 were noted for TFs linked to ATRA-independent variants. Supporting these findings, an enrichment analysis of putative target genes (implicated by brain eQTL or neural Hi-C) of SNPs in retinoid TF motifs or ChIP peaks (**Supplementary Table 4.2**) revealed 5 genes nominally enriched for high regional expression in rhombomere 9, which gives rise to medullary populations of serotonin neurons<sup>107</sup>. (The full results can be explored at <https://maayanlab.cloud/Enrichr/enrich?dataset=27d6db2a8510a90ed0d78e6b60c59287>).

Using the R2 database (<http://r2.amc.nl>), I examined expression of *FEV*, *GATA2*, and *GATA3* TFs in 24 human neuroblastoma lines (GEO accession GSE28019), retinoic acid-treated human SH-SY5Y neuroblastoma cells<sup>108</sup>, alongside human neural progenitors<sup>109</sup> and melanoma lines as comparators<sup>110</sup>, confirming neuroblastomas strongly express all three of these TFs (**Supplementary Figure 4.5**). Single-cell RNA sequencing data from mouse brain confirms the specificity of these TFs, revealing that these TFs are only expressed in serotonergic, noradrenergic, peripheral autonomic, and midbrain inhibitory neurons—with all three expressed in serotonin neurons<sup>111</sup>. Furthermore, exogenous retinoids have been shown to lower circulating serotonin in humans<sup>45</sup> and to alter morphology of rat raphe neurons in slice culture<sup>112</sup>, suggesting these neurons are retinoid responsive. I do not believe my finding is an artifact of the N2a system, as I could

identify no evidence in the literature suggesting a serotonin-like identity of N2a cells with or without ATRA treatment. Altogether, these findings suggest serotonin neurons and closely related cell types<sup>93</sup> may be cellular points of convergence for several retinoid-mediated functional SNP effects on MDD risk.

## **4.4 Discussion**

To date, most functional investigations of SNPs in the context of psychiatric disorders have taken place in a low-throughput manner, such as single-variant classical reporter assays<sup>20</sup> or using CRISPR-Cas9 technology to edit limited positions for deep phenotyping<sup>113</sup>. Here, I leveraged MPRA to screen over 1,000 SNPs from loci associated with MDD, related phenotypes, and broader psychiatric disease, demonstrating the utility of this technique for dissecting the functional regulatory architecture of psychiatric GWAS loci, and defining shared upstream regulatory features across loci.

In doing so, I identify over 100 SNPs with allelic effects on expression, with most coming from loci containing  $\geq 2$  functional SNPs. These data provide experimental support for the prediction that multiple SNPs with allelic effects exist within GWAS loci as put forth in polygenic/omnigenic theory literature. I further examined the omnigenic hypothesis' more central prediction of regulatory convergence across loci. By examining the shared regulatory features (TF binding motifs) based on enrichment at functional SNPs, I were able to predict several TFs with TR activity recurrently altered across MDD-associated SNPs, highlighting retinoid receptors in particular.



Retinoids are encountered both exogenously (e.g., as ATRA in oncology, and as isotretinoin, carrying a black-box warning for suicidality) and endogenously, including during brain development. To investigate how SNP functions may be altered by retinoids, I repeated the assay with an ATRA condition. ATRA drastically rearranged the TR landscape of N2a cells, resulting in altered and novel allelic effects at over 100 SNPs and revealing ATRA-dependent mechanisms of function across 122 SNPs from 22 of 26 MDD GWAS loci assessed. Of 17 ATRA-interacting functional SNPs overlapping ChIP peaks for retinoid receptors, 15 overlapped ChIP sites of RXRA (**Supplementary Table 4.2**), suggesting it may be central in functional SNP activity at retinoid receptor binding sites in this system. Interestingly, single-cell epigenomics of human cortical cell types recently found RXRA motifs to be uniquely enriched in open chromatin of SST interneurons<sup>114</sup>, a strong candidate cell type for MDD<sup>115</sup>. These findings suggest that retinoid receptors—RXRA in particular—merit mechanistic follow-up regarding TR differences at MDD-associated SNPs. Future work may be able to leverage biobank-level datasets to ascertain whether retinoid-interacting SNPs are overrepresented in retinoid-treated patients experiencing adverse psychiatric side effects. While data on endogenous retinoids, *e.g.* plasma values, are not currently available in large genotype-phenotype-health record cohorts like UK Biobank, future datasets may enable investigation of circulating retinoids and their interaction with genotype in cognitive and psychiatric phenotypes.

The methodologic requirements of high-throughput assays such as MPRA bring inherent limitations to their results. The primary precaution in interpreting these results concerns cell type relevance. MPRA are subject to the TR landscape of the cell type used. Neuroblastoma cells, including N2As, are derived from peripheral neural crest progenitors—though they can be

differentiated into dopaminergic neurons<sup>106</sup> and commit to neuronal differentiation with ATRA<sup>61,67</sup>—and were selected for these assays based on intact retinoid signaling rather than representing a disease cell type per se. On the other hand, the neural crest-derived autonomic nervous system has received little consideration (relative to brain) in psychiatric genetics of MDD despite the well-appreciated role of stress in depression. These data may form an interesting foundation for future study of autonomic effects of MDD genetic risk.

Still, I can broadly speculate on brain cell types implicated by my findings. A notable prior pharmacology MPRA cleverly tested gDNA fragments for regulatory activity over a time course of dexamethasone treatment, while collecting epigenomic data in the same cell type over the same time course to compare MPRA signal and endogenous genomic marks. They found that endogenous genomic regulatory elements with repressive marks or depleted of glucocorticoid receptor binding were oftentimes active and/or dexamethasone-differentially active when assayed on the MPRA plasmids. This suggests that the transcriptional-regulatory capacity of an MPRA is not constrained by the epigenome of the model cell, but rather by its expressed TFs<sup>60</sup>. As such, retinoid receptor-mediated SNP functions observed are not limited to sequences that would be active in the N2a genome; as such, it is entirely plausible that the observed effects also occur in retinoid-receptor expressing brain populations. Mouse nervous system single-cell RNA-seq suggests retinoid receptor expression is absent in brain glia, but robust in many neuron types<sup>111</sup>. Thus, I suspect the directly-mediated retinoid receptor SNP effects I observe may be neuron-specific. Future studies may be able to address the interesting question of differences in neuronal subtypes exhibiting functional SNP effects.

I find that principles of the omnigenic model appear to hold true for MDD risk genetics, including the presence of far more functional variants (a total of 277 SNPs with allelic and/or interaction effects of 1178 assessed across the two assays) than there were GWAS loci (*i.e.*, tag SNPs). I find, interestingly, that functional SNPs form convergent subsets of upstream (transcription-regulatory) sequences and systems, which in turn have shared retinoid dependence and are collectively enriched in serotonin neurons via 8 ATRA-interacting functional SNPs in binding motifs of GATA2, GATA3, and FEV. It has previously been demonstrated that systemic administration of ATRA depletes serotonin by over 40% in the rat brain<sup>116</sup>, supporting the serotonin system as a convergent target of retinoid-regulated pathways. As GWAS of MDD begins to explore severe, treatment-refractory cases<sup>7</sup>, it will be interesting to see whether associated variation still shows such convergence, as treatment-resistant depression (generally, non-response to two or more classes of antidepressant) effectively signifies non-response to multiple serotonergic agents.

In all, I assessed the architecture of *cis*-regulatory variation in psychiatric disease risk loci, identifying at least one functional SNP in the majority of the 40 GWAS loci examined, largely corresponding to MDD-associated SNPs. Strikingly, retinoid receptor binding sites and TR systems subject to regulation by ATRA have a substantial impact on whether and how MDD-associated SNPs are functional. These findings constitute a robust experimental demonstration of the influence of physiological and environmental states on the molecular activities of disease-associated SNPs, and constitute a high-confidence set of MDD SNPs meriting deeper functional characterization of both their TR mechanisms and their environmental interactions.

## 4.5 Acknowledgments

I thank Stephen Plassmeyer and Tomás Lagunas, Jr. for technical assistance; Barak Cohen, PhD, Michael A. White, PhD, Dana King, PhD, Brett Maricque, PhD, and Ryan Friedman for library design, cloning, and analysis advice; and Genome Technology Access Center@McDonnell Genome Institute (GTAC@MGI) and Jessica Hoisington-Lopez and MariaLynn Crosby from the DNA Sequencing Innovation Lab at The Edison Family Center for Genome Sciences and Systems Biology for sequencing support. This work was supported by the NIH (1F30MH1116654 to BM, and 1R01MH116999 to JDD) and The Simons Foundation (571009 to JDD). GTAC@MGI is supported by UL1 TR002345. SNP annotation data resources are detailed in the supplementary text. I would also like to thank Idoya Lahortiga, Ph.D. and Luk Cox, Ph.D. curators of Somersault1824 (<https://gumroad.com/somersault1824>), for their open-access, Creative Commons BY-NC-SA 4.0 licensed libraries of high-quality biomedicine graphics (especially those from Graphite Life Explorer, ePMV, and Eyewire), adapted for Figure 1E. Color palettes for plots are from the R package wesanderson (<https://github.com/karthik/wesanderson>). A summary spreadsheet of all significant SNPs identified in one or both assays, along with full analysis results, including barcode-wise expression in each sample, single-condition allelic effect tests, linear modeling results, and significantly enriched TFs in each of the comparisons executed is available at [https://bitbucket.org/jdlabteam/n2a\\_atra\\_mdd\\_mpra\\_paper/src](https://bitbucket.org/jdlabteam/n2a_atra_mdd_mpra_paper/src). Raw sequencing files are deposited in GEO under accession GSE167519. Additional supplementary information is available at *Translational Psychiatry's* website. There were no conflicts of interest in the performing these experiments, analyses, or in the course of publishing these findings.

## 4.6 Supplementary Materials

### 4.6.1 Supplementary methods

**Library Design and SNP selection.** The library was designed by selecting tag SNPs of interest from neuropsychiatric trait and disease GWAS studies, predominantly for MDD or multi-diagnosis groups including MDD cohorts<sup>8–10,41,87,117–120</sup>, and from GWAS of traits with high SNP co-heritability with MDD—namely, neuroticism<sup>121,122</sup> and mood instability<sup>123</sup>. Additional variants discovered in other psychiatric disorder GWASes at or near MDD tag variants were included from studies of ASD<sup>120</sup>, anxiety disorders<sup>124</sup>, and attention-deficit hyperactivity disorder<sup>125</sup>. Two tag SNPs associated with intelligence<sup>126</sup> and educational attainment<sup>127</sup> near the gene *PTGER3*, which I previously illustrated to be sex-differentially expressed and functional in the mouse locus coeruleus<sup>93</sup>, were included. One negative control tag, rs1883640, close to the transcription start of gene *CDKALI* and associated with several anthropomorphic traits<sup>128</sup>, was also included as a negative control locus. In all, 38 tag SNPs were selected for LD expansion and epigenomic overlap screening in the LD neighborhood. All LD partners at  $R^2 > 0.65$  were included without consideration of epigenomic annotation intersections for two additional MDD-associated tag SNPs near sex-differentially expressed genes of mouse LC—*Slc6a15* and *Lin28b*<sup>93</sup>; SNPs from these two loci with  $R^2 < 0.65$  were also included solely on the basis of a RegulomeDB score of  $\geq 4$ , signifying a SNP overlapping both a TF binding site and a DNase hypersensitive site across epigenomic data curated by the tool<sup>129</sup>.

LDlink (using dbsnp 151) was queried for each tag SNP to acquire all biallelic EUR (or in the case of tags from Han Chinese CONVERGE<sup>9</sup>, Han Chinese (“CHB”)) SNPs in LD with the tag<sup>130</sup>. These SNPs were subsetted to those with a minor allele frequency  $\geq 1\%$  and with an LD  $R^2 > 0.1$  with the tag, as GWAS studies generally define ‘independent loci’ as SNPs in with LD  $R^2 < 0.1$ .

The hg19 coordinates of the retrieved, and subsetted LD SNPs were then intersected to those of myriad CNS epigenomic datasets (also in hg19 coordinates), including postmortem brain tissue eQTLs from GTEX v7<sup>78</sup>, PsychENCODE<sup>77</sup>, the CommonMind Consortium<sup>76</sup>, the Lieber Institute/Brainseq Consortium<sup>131,132</sup>, and ROSMap<sup>79</sup>; enhancers predicted based on human postmortem adult and fetal brain tissue histone marks or enhancer RNAs<sup>70,77,133–136</sup>; and chromatin contacts for human neural cell types identified *in vitro*<sup>75</sup>.

LD SNPs were selected by manual inspection of intersecting epigenomic annotations within an LD block, including the negative control *CDKALI* block. Negative control locus SNPs were selected while blinded to control status of the locus (by replacing the parent locus' name as one coming from a sub-region of interest during the selection process). 11 SNPs were removed from consideration due to overlap with a nonsynonymous coding SNP. For inclusion, a SNP was required to be an eQTL in at least one of the brain datasets mentioned and intersect with  $\geq 1$  additional annotation track; two exceptions to this were that no overlaps were required for the SNPs included by RegulomeDB score (see above), and a single overlap (eQTL or otherwise) was considered adequate for inclusion in a small minority of loci where most SNPs did not intersect any annotation. In annotation-rich regions, SNPs with the greatest diversity and abundance of intersecting annotations were selected. While computationally selecting SNPs based on the greatest number of intersections would be a more time-efficient approach, recent work suggests that regulatory variants may be better predicted by training deep learning models on local genomic sequences rather than sequences from across the full genome<sup>137</sup>.

To design the MPRA oligonucleotide sequences ordered, genomic tiles of 126bp, centered on the variant of interest were extracted from human reference genome hg19. For variants where the alleles were not of the same length (e.g., single-base deletions or multi-base alleles), the longer allele's sequence spanned 126bp, with the shorter allele spanning 126-(difference in length) bp. As oligonucleotide synthesis requires uniform sizing, all oligonucleotides were brought to a final length of 200bp by adding bases between cut sites used for directional insertion of the reporter gene, such that the bases are absent from the final plasmid library.

Sequences containing cut sites that would interfere with cloning were taken from the smallest up- or down-stream window to keep the SNP as near the center of the sequence as possible, by shifting enough to trim away 1-2 bases of the interfering cut site. When this required shifting sequences more than  $\pm 30$  bp, the selected SNP was removed from the design process (resulting in the removal of 8 SNPs from the selected set). The final, cut-site-free genomic tiles representing each allele of the 1453 selected SNPs were programmatically added to random ten, 10bp barcodes, which were pre-subsetted to be:  $\geq 2$  Hamming distances from one another, 25-75% GC, without  $>4$  of any one nucleotide in a row, and without restriction sites or partial restriction sites that would recreate a full site during the cloning process. Control sequences were likewise paired to 10 barcodes each. Finally, 'basal' oligonucleotides with 126 bases filling the space between gene-insertion cut sites (to ultimately place the minimal promoter directly adjacent to the reporter gene) were generated and paired to 110 barcodes to get an accurate measure of transcriptional output driven by the hsp68 minimal promoter in the absence of upstream mammalian genomic sequence.

***Cell culture, transfection, and RNA collection.*** Mouse neuroblastoma N2A cells were grown in uncoated 6-well plates in medium consisting of DMEM and 10% Fetal Bovine Serum (2% fetal bovine serum for the ATRA assay, based on media conditions from the literature)<sup>66,67</sup>. Cells were fed 2mL per well at each passage, and medium was refreshed as needed indicated by yellowing of the phenol red indicator contained in the DMEM solution. Cells were incubated at 37°C, 5% CO<sub>2</sub>.

For transfection, a tube containing 10µL per replicate of Lipofectamine 2000 in 200µL per replicate OptiMEM was prepared and incubated at room temperature for 5 minutes. Meanwhile, a second tube containing 2.5µg of MPRA plasmid library per replicate and 200µL of OptiMEM was prepared. The contents of the latter were added to the tube containing Lipofectamine 2000 and incubated at room temperature for 30 minutes. The mixture was then aliquoted out in volumes of 400µL to each well (replicate) and the plate re-lidded. The plates stayed 45-90 minutes in the biosafety cabinet until cells, grown and split from the same parent culture in both assays and collected from all wells of a well plate (first assay) or T150 culture flask (second assay), were aliquoted on top of the 400µL of plasmid•lipofectamine•optiMEM mixture.

Cells to be transfected were lifted from their existing wells by adding 1mL of 0.25% Trypsin•EDTA and incubating at 37°C for 5 minutes. 2mL of DMEM were added to each well, and the full 3mL of DMEM, trypsin, and cells were collected from each well into a tube containing  $n$  mL of 10% FBS in DMEM, where  $n$  was the number of wells being passaged. Cells were spun down at 700 rcf for 5 minutes at room temperature, then resuspended 10mL in antibiotic-free N2A medium (for initial MPRA, 0.1µM vacuum filter sterilized 10% FBS in DMEM; for retinoic acid MPRA, 0.1µM vacuum filter sterilized 2% FBS in DMEM). Duplicate counts of cells were taken



using the Countess automated cell counter. The average number of reported viable cells per mL was used to determine plating volumes to achieve target cell densities in 2mL volumes, and the required volume of the same medium used for resuspension (less 10mL) was added to the resuspended cells before aliquoting. For the initial MPRA, these target densities were three 9.6 cm<sup>2</sup> wells with  $5.7 \cdot 10^5$  cells/well and three 9.6 cm<sup>2</sup> wells with  $8.3 \cdot 10^5$  cells/well. Linear modeling of results from the first assay revealed singular fits when attempting to fit a variable for initial plating density, indicating this played no role in the detected expression. In the retinoic acid MPRA, twelve 9.6 cm<sup>2</sup> wells at  $7.8 \cdot 10^5$  cells/well. 2mL of resuspended cells were aliquoted to each well, then incubated at 37°C for (6.5 hours in initial MPRA, 7.5 hours in retinoid MPRA). At that time, medium was removed and replaced with the respective experiment's N2A medium, now containing antibiotics (and in the case of the retinoid experiment, either 20µM all-trans retinoic acid in DMSO, or an equivalent amount of DMSO alone for vehicle). For the initial MPRA, medium was not changed before cells were collected; for the retinoid MPRA, medium was changed every 24 hours with a freshly prepared aliquot of N2A medium supplemented with respective drug or vehicle. Retinoic acid and vehicle were prepared in a dark room in 500 and 200 µL aliquots, respectively, stored in black microfuge tubes at -20°C and handled opened only in biosafety cabinets without room or hood lights on.

In the initial MPRA, cells were collected 72 hours after transfection by rinsing wells twice with 1mL of DPBS, then adding 1mL of DPBS, thoroughly lifting cells with a scraper, and collecting the full 1mL of cells from the well into its own microfuge tube. Cells were spun down at 3000 rcf at 4°C for 10 minutes. 750µL of supernatant was removed, and replaced with 750µL Trizol. Samples were shaken vigorously by hand for 10 seconds to lyse cells and incubated on the

benchtop at room temperature for 10 minutes. In the retinoic acid MPRA, the same rinses were performed; followed by adding 250 $\mu$ L dPBS, then 750 $\mu$ L of Trizol LS to each well. The plate was swirled until cells were clearly lysed, at which point the full 1mL volume of each well was collected into a centrifuge tube, allowed to stand at room temperature for 10 minutes.

The handling of samples from both experiments was the same from that point forward: 200 $\mu$ L chloroform was added to each tube, which was then shaken vigorously by hand for 30 seconds. Tubes sat at room temperature for 7 minutes, and then were spun down at 12,000 rcf at 4°C for 30-40 minutes. 325 $\mu$ L of the aqueous phase was collected, and added to tubes pre-filled with two volumes (650 $\mu$ L) of Zymo RNA Binding Buffer and their combined volume (975 $\mu$ L) of 100% ethanol. The RNA solution was then purified using the Zymo Clean and Concentrator-5 kit per the manufacturer directions. RNA was eluted into 40 $\mu$ L of nuclease-free water.

To remove any residual plasmid from the samples, the Turbo DNA-Free kit was used to perform vigorous DNase treatment (2 $\mu$ L of DNase per RNA sample) according to the manufacturer instructions. RNA with DNase in solution was incubated at 37°C for 1 hour, followed by addition of the kit's DNase inactivating matrix and a 10 minute spin at 8000 rcf at room temperature. 44 $\mu$ L of supernatant was collected, and subjected to a second round of purification using the Zymo Clean and Concentrator-5 kit as described above. (A second purification was required in order to remove an unknown component of the Turbo kit which otherwise interfered with RNA concentration/integrity measurements).

**Power analysis of first MPRA to inform n for ATRA MPRA.** Using the code underlying the excellent MPRA tools package from Ghazi, et. al<sup>138</sup>, I performed power analysis of the first experiments' element (single allelic sequence)-level measurements as calculated after the filtering steps discussed in methods. From the first experiment, I find an SD of 0.0396 (25th %ile), 0.0556 (50th %ile), and 0.0801 (75th %ile). I assumed 6 barcodes per sequence, the median number of barcodes analyzed per sequence in the first experiment after filtering steps, with 6 replicates, to determine how well-powered I were for detecting effects at a Bonferroni-corrected p-value of 0.1 (though I used an empirical test statistic correction instead of Bonferroni correction, as described in the Methods section). The results of this analysis are shown in **Supplementary Table 4.1**. Based on these findings, I determined that I would be similarly well-powered to differentiate allelic effects and allele-drug interactions using an n of 6 vehicle and 6 ATRA samples in the second experiment.

**Sequencing library preparation.** 1µg of RNA was subjected to target-specific cDNA amplification using a primer reverse complementary to the proximal polyA signal sequence in the reporter transcript. 20µL reactions containing 1µg RNA, 4µL 5x SSRT first strand buffer, 0.6µL rRNAsin (Promega), 0.9µL SSRT III enzyme, 1µL 12µM antisense primer, 1µL 100 mM DTT, and 1µL 10mM dNTPs, brought to volume with water. Reactions were incubated at 50°C for 1 hour. Remaining enzyme, primers and RNA were then removed by adding 3.5µL Exosap-IT and incubating for 15 minutes at 37°C, followed by addition of 1µL 0.5M EDTA, addition of 3µL 1M NaOH with heating at 70°C for 12 minutes, and neutralization with 3µL 1M HCl.

The single stranded cDNA was then purified by washing 10 $\mu$ L MyOne silane beads per sample with Buffer RLT (Qiagen), resuspending beads in 3•sample volume of RLT, and adding to each sample, mixing end-over-end for 30 minutes to bind cDNA. Beads were magnetized and washed twice with 80% ethanol, air dried for 5 minutes, and bound cDNA eluted into 12 $\mu$ L 5mM Tris HCl pH 8. 10 $\mu$ L of this product (cDNA) or 80ng transfected plasmid DNA were then used to generate double-stranded, target-specific product from the cDNA using the same antisense primer as reverse primer and a primer in the 3'UTR WPRE element of the reporter as the forward primer for 15 cycles with Phusion HF 2x Master Mix. Product was size-selected using Ampure XP beads by bringing the 20 $\mu$ L PCR reaction to 100 $\mu$ L with water, adding 80 $\mu$ L beads, magnetization, capture of 160 $\mu$ L supernatant and adding 40 $\mu$ L fresh beads to it, magnetization, 3x washes with 80% ethanol, and air drying, followed by elution into 12 $\mu$ L 5 mM Tris HCl Ph 8. (This procedure is henceforth referenced as “80/40 Ampure XP cleanup”).

The double stranded product was then digested with HindIII-HF (NEB) and NheI-HF (NEB) in CutSmart buffer at 37°C for 1 hour, to create sticky ends flanking the barcode sequences to which sequencing adapters could be attached. The digested product was purified as above, but using 100 $\mu$ L beads and 50 $\mu$ L beads, with 180 $\mu$ L supernatant recovery added to the second (50 $\mu$ L) aliquots of beads. Product was eluted as above, using 10 $\mu$ L eluate for adapter ligation.

Adapters were ligated using Enzymatics T4 ligase (1 $\mu$ L), Enzymatics T4 10x buffer (2 $\mu$ L), 1 $\mu$ M HindIII-compatible, staggered-length (to improve read heterogeneity for sequencer compatibility) Illumina adapters for read 1 and 1 $\mu$ M single-length NheI-compatible adapter for read 2. Ligation ‘cycling’ was used, alternating cycles of 30s at 25°C (annealing) and 30s at 16°C (ligase activity)

for 560 cycles, a procedure previously shown to improve ligation yields twofold in restriction cloning<sup>139</sup>. Product was purified and eluted with the 80/40 Ampure XP cleanup; 10 $\mu$ L eluate was used for 50 $\mu$ L index PCRs with sample-specific read 1 indices, user-specific read 2 indices, 2x Q5 Ultra II Mastermix (NEB) for 9 cycles. Product was purified using 80/40 Ampure XP cleanup, eluted into 19 $\mu$ L of 5mM Tris HCl pH 8. 17 $\mu$ L were recovered (2 used for verifying proper product at ~330bp was produced by using the D1000 kit with the Tapestation 4200 system). Up to 15 $\mu$ L of each library was mixed at an equimolar amount, with approximately 1.5 times the molar amount of the sample corresponding to input DNA for maximally accurate quantification of this single DNA sample.

***MPRA Analysis.*** Raw fastq files for each RNA sample are initially processed by using the bash terminal and the string-matching command, sed. The string matching requires an exact match of 6bp upstream sequence and 8bp downstream sequence corresponding to the restriction-ligation sites wherein the reporter was ligated to the barcode, and the barcode to the plasmid, respectively. For each string matching this pattern exactly, the 10 bases of barcode sequence between are extracted and put in a text file for the corresponding fastq. To then get barcode counts, sequences in the text file of extracted barcode-positioned 10mers are subsetted to those containing a perfect 10 base match to a barcode sequence in the library and are counted using the R function *table* , which tabulates the occurrences of each barcode sequence specified in a reference metadata file (containing the designed barcodes only and related information) for each sample's string-extraction output /text file. This process is extraordinarily efficient, requiring only 1-2 minutes per 20 million reads on a personal computer.

Following the tabulation of barcode counts, barcode counts are converted to counts per million (mapped) before filtering. Several filtering steps are then used to ensure robust representation of barcodes, their paired sequences, and samples before testing expression effects. These steps proceed in the following order. 1) Barcodes with a DNA count under a specified threshold (75) are struck from the counts table from the DNA *and* RNA samples. 2) For sequences with 4 or fewer remaining DNA barcodes represented across all samples after this step, all other barcodes for the sequence are struck from all of the samples in the table to avoid analysis of sequences with inadequate barcoding depth. 3) Counts of RNA barcodes are then removed on a per-barcode-per-sample basis if they fall below a separate minimum read threshold (30 counts), set below the DNA threshold to allow for detection of repressive effects). At this point, preliminary expression values for barcodes are calculated ( $\log_2$  of the ratio of RNA barcode cpm to DNA barcode cpm) in each replicate and collapsed into a mean per sample. Single barcodes with expression values  $\geq 2$  standard deviations apart from other barcodes for a given sequence in a given sample are dropped only from that sample, as this suggests the barcode exerted a cryptic regulatory effect (e.g. as a 3'UTR element) in the sample, or that either mutations in other barcode(s) on delivered plasmid or during sequencing preparation are contributing to spurious counts for that barcode. Penultimately, all barcodes for a sequence are dropped from an individual sample if the filtering steps up until now have resulted in 4 or fewer barcodes remaining for a given sequence in that sample.

After these filtering steps, expression is calculated on a per-barcode basis in each sample by taking  $\log_2$  of the ratio of RNA cpm to DNA cpm. These values are then averaged for each sequence in each sample, resulting in a per-sequence expression value in each sequence. For the first assay,

where the sole comparison is between alleles of a sequence, a student's t-test is applied to compare the vector of samplewise mean expression values for one allele to those of the other allele. In the case of linear mixed modeling, as in the ATRA assay, the barcode-level expression values for the two respective alleles are used as model input, along with the corresponding allele for each barcode and drug/vehicle status of the corresponding samples to model both as  $\text{expression} \sim \text{allele} + \text{drug} + \text{allele}:\text{drug}$ . As a failsafe step, sequences with a calculable mean expression value in  $< 2$  of the samples in a condition (of  $n=6$  per condition in both experiments) were excluded from t-testing in the first experiment; in the ATRA experiment utilizing an LMM of barcode expression values, SNPs for which  $\geq 60\%$  of samplewise barcode expression measurements were missing (i.e.,  $\geq (12 \text{ replicates} * 10 \text{ barcodes per allele} * 2 \text{ alleles} * 0.6 \implies) 144$  sample-barcode expression measurements) were excluded from analysis.

To perform statistical correction informed by the level of cross-sample and cross-barcode noise, the vector of Z values (t test) or F values (LMM) were stored for each SNP. For single-condition t-test analysis of the ATRA samples alone, Benjamini-Hochberg FDR correction was used, as my statistics of interest regarding these samples came almost entirely from the LMM (except for dual interaction and allele main effect cases, see below). To correct test statistics in the first assay, second assay vehicle-only, and LMM analyses, the same statistical comparison (t-test or LMM) as used for the sequences of interest was performed utilizing a vector of 10,000 (first assay), 20,000 (second assay LMM), or 5,000 (second assay vehicle samples alone) Z or F test statistics random "allelic" comparisons between elements that should have a known null difference in transcription—that is, 5,000-10,000 comparisons between two sets of 6 randomly selected barcodes—representing the median number of barcodes analyzed for a given sequence—from

among the 110 barcodes assigned to the hsp68 promoter without human genomic sequences added upstream. This process effectively models the level of noise in expression calculations for the experiment due to e.g., variation in the cultures, across barcodes, and via biases during sequencing preparation. This method derives from an analogous statistical correction approach using non-targeting CRISPR gRNA sets from within a CRISPR screening library<sup>68</sup>.

For SNPs with significant allele and interaction coefficients in the LMM, a meaningful allele main effect was considered present if the single-condition vehicle and ATRA analyses showed the same allelic direction of effect, with a vehicle  $p_{\text{emp}} < 0.1$  and ATRA FDR  $< 0.1$  (*i.e.*, near-significant within each condition of  $n=6$ , thus reasonably capable of achieving significance in the LMM analysis of the two conditions combined).

***MotifbreakR analysis.*** In order to assess transcription factor binding motif perturbations corresponding to SNPs, I utilized the R package motifbreakR<sup>69</sup> and its built-in database of motif position-weight matrices (PWMs) from multiple public repositories. The tool by design identifies PWM matches overlapping input SNPs from dbSNP (here, version 151) for which at least one of the SNP alleles results in a genomic sequence significantly matching a given motif sequence. I used the default significance cutoff of  $p < 10^{-4}$  for calling motif matches in all analyses and identified changes in motif match score using the tool's default algorithm, which takes a weighted sum based on the position weights of each base in the motif sequence and considers these for the two alleles of the query SNP. As the use of dbSNP 151 required use of hg38, my queries thus took place in hg38 reference genome sequence space for the subset of SNPs with the same rsID in both the MPRA design and dbSNP 151. Computationally, all SNPs from a given set of positive and



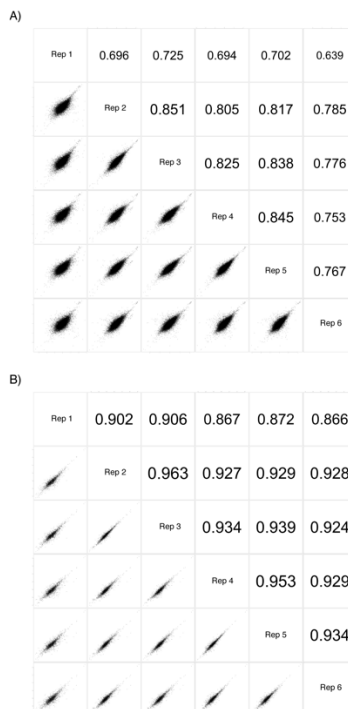
negative comparators were run through motifbreakR once with these settings to identify motif perturbations. For analysis of SNPs with identified in the first assay in N2a cells, SNPs with an allelic effect  $p_{\text{emp}} < 0.05$  were compared to those with an allelic effect  $p_{\text{emp}} > 0.05$ ; from the second assay, allele-ATRA interaction SNPs with interaction  $p_{\text{emp}} < 0.05$  were compared to allele-main-effect-only SNPs (allele  $p_{\text{emp}} < 0.1$  and interaction  $p_{\text{emp}} > 0.1$ ), and allele main effect SNPs at allele  $p_{\text{emp}} < 0.1$  were compared to the SNPs that were subject to neither main effects of allele or drug, nor their interaction (all  $p_{\text{emp}} > 0.1$ ).

Frequencies at the level of TF (which can include several motifs) were considered as the number of SNPs matched to a given TF, regardless of the number or identity of motifs to which that SNP matched. Null distributions of frequency were determined by 10,000 random selections of motif perturbations corresponding to  $n$  SNPs in the negative comparator SNP set, where  $n$  was the number of positive set SNPs analyzed based on the presence of an rsid in dbSNP 151. The p-value of frequency was then calculated from the empirical percentile of the positive SNP frequency count vs the distribution of frequencies in the negative set selections. Concordance was determined in a similar manner by permuting 10,000 random subsets of data covering  $n$  SNPs, but drawing motif assignments from the full set of SNPs analyzed and assigning a random “MPRA” allelic-directional effect to each SNP to compare observed rates of concordant allelic effects on motif match and reporter expression to the rates obtained by chance.

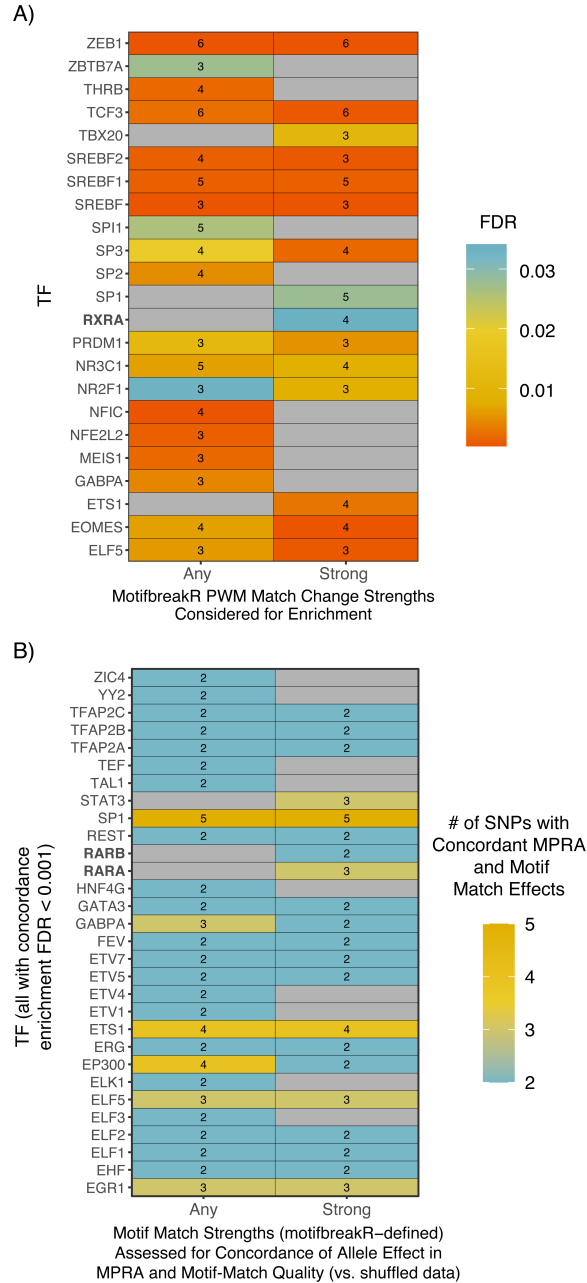
***ABAEnrichment Analysis of TF Sets.*** To examine enrichment of motifbreakR-identified TFs in spatiotemporal human brain gene expression, I utilized my four sets of TFs identified from the ATRA experiment (TFs enriched at functional variants with allele•ATRA interaction effects over

allele main-effect-only SNPs, and those enriched at allele-main-effect-only variants vs no-effect variants, considering only “strong” motif perturbations or all motif perturbations in each set). I constrained these TF sets to those considered enriched on the basis of motif perturbation at  $\geq 2$  functional SNPs. Enrichment was then tested in the 5-developmental-epoch dataset (27 brain regions per epoch) and Allen Atlas adult microarray data of over 500 brain regions using a Wilcoxon approach, weighting TFs by their “score” (number of SNPs assigned by motifbreakR at their set’s strength level). Enrichment was tested using 8 percentile thresholds for whole-dataset genes to be considered expressed in brain: 25, 37.5, 50, 62.5, 75, 82.5, 90, 95. The ABAEnrichment tool then gives a Family Wide Error Rate (FWER) by performing 1,000 shuffled analyses at each cutoff for each age•region•threshold (8 thresholds x 27 regions x 5 stages = 1080 enrichment tests).

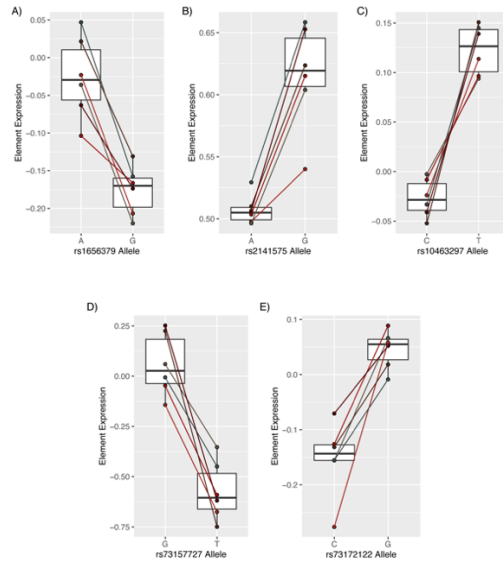
#### 4.6.2 Supplementary figures and tables



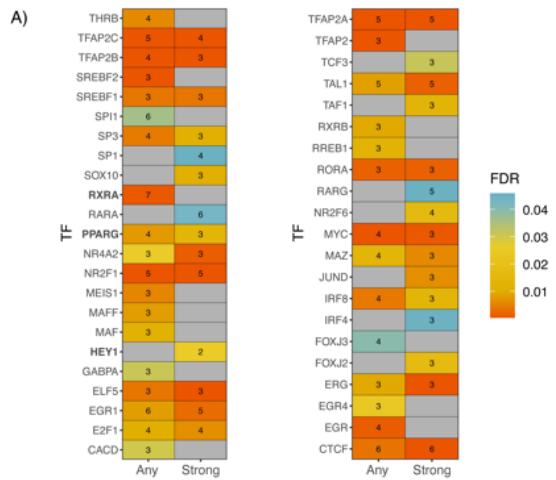
**Supplementary Figure 4.1. Cross-replicate correlations from initial N2A MPRA. A)** Cross-replicate comparison of individual barcode counts per million (CPM). **B)** Cross-replicate comparison of element-level expression (mean barcode RNA/DNA ratio). Pearson correlation coefficients are shown above the diagonal in both panels.



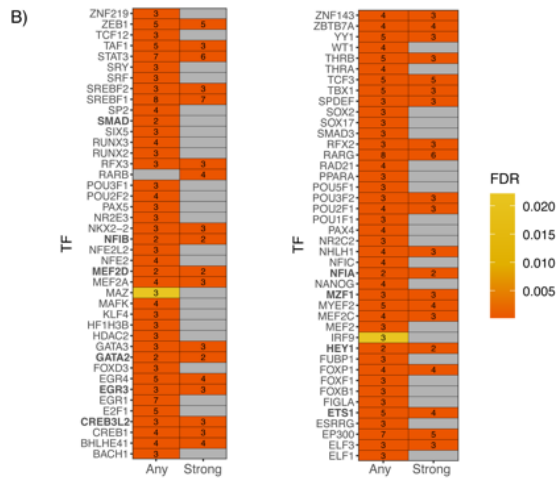
**Supplementary Figure 4.2. Extended motifbreakR results from the first MPRA.** Enrichment analysis results agnostic to the motifbreakR-defined “strength” of motif change between alleles (left-hand columns). **A)** TFs with enriched frequency of motifs among functional SNPs. The corresponding number of functional SNPs matched to each TF for a given strength are shown. **B)** TFs with greater than predicted concordant motif and MPRA effects among functional SNPs. Concordant effects were defined by greater MPRA expression driven by the allele better-matched to the corresponding TF motif and vice versa—the expected behavior of strictly activating TFs.



**Supplementary Figure 4.3. Additional functional SNPs corresponding to an enriched TF motif group (NR3C1) from the first assay.** Replicates are indicated by individual points and their connecting lines.

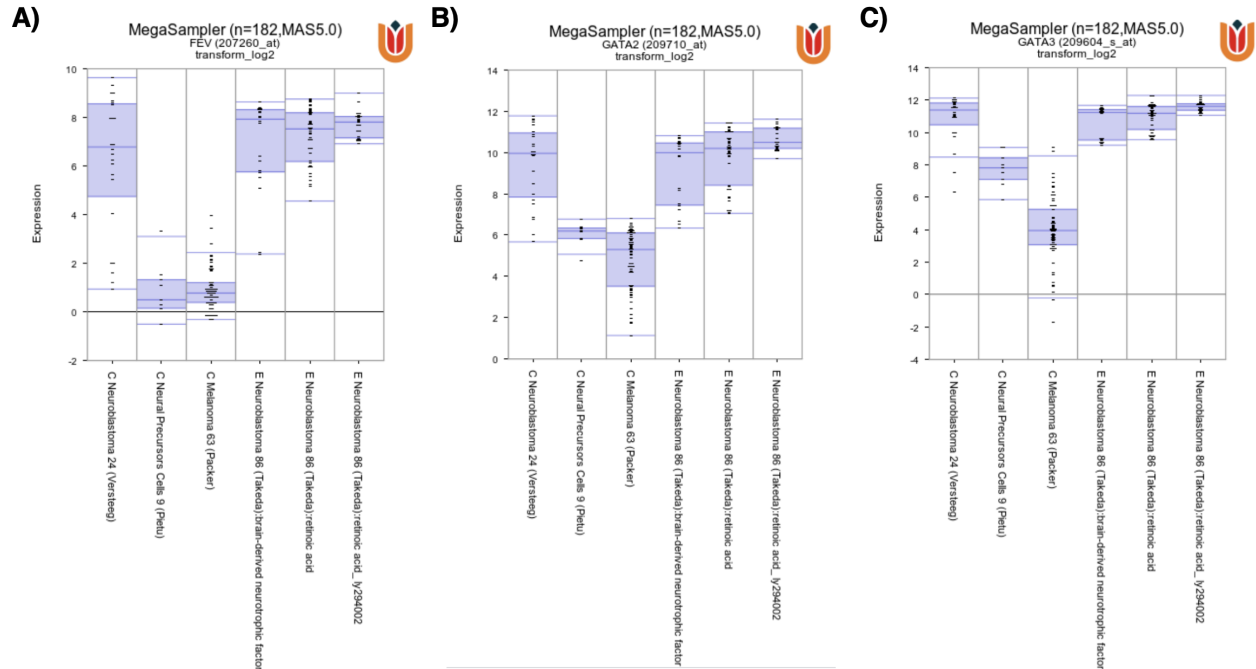


MotifbreakR PWM Match Change Strengths Considered for Enrichment (Allele Main Effect-Only SNPs vs. No-Effect SNPs)



MotifbreakR PWM Match Change Strengths Considered for Enrichment (Allele-Drug Interaction SNPs vs. Allele Main Effect-Only SNPs)

**Supplementary Figure 4.4. Extended motifbreakR results from the ATRA treatment MPRA.** Enrichment analysis results agnostic to the motifbreakR-defined “strength” of motif change between alleles (left-hand columns). **A)** TFs with enriched frequency of motifs among retinoid-independent functional SNPs compared to SNPs with no detected allelic effects. The corresponding number of functional SNPs matched to each TF for a given strength are shown. Heatmap is split into two vertical slices for visibility. **B)** TFs with motifs overrepresented among ATRA-dependent (interaction) functional SNPs compared to ATRA-independent functional SNPs. The heatmap is shown in halves for visibility. TFs identified as ATRA-upregulated in human neuroblastoma lines<sup>105</sup> are in bold font.



**Supplementary Figure 4.5. Expression of TFs *GATA2*, *GATA3*, and *FEV* in a variety of human neuroblastoma cell lines at baseline, with retinoic acid treatment, and in melanoma cells and neural precursors for comparison.** Expression values are log2 normalized microarray expression values. Each dash within the a given bar represents a single sample or replicate from the corresponding GEO dataset. Results were visualized using the R2 browser (<http://r2.amc.nl>). **A)** Expression of *FEV*. **B)** Expression of *GATA2*. **C)** Expression of *GATA3*.

**Supplementary Table 4.1. Power analysis for allelic effect size detection.**

Bonferroni $p$	Power	SD Expression	log2 Allelic Fold Change
0.000	0.000	0.040	1.000
0.050	0.838	0.040	1.051
<b>0.100</b>	<b>1.000</b>	<b>0.040</b>	<b>1.105</b>
0.000	0.000	0.056	1.000
0.050	0.320	0.056	1.051
<b>0.100</b>	<b>0.999</b>	<b>0.056</b>	<b>1.105</b>
0.000	0.000	0.080	1.000
0.050	0.059	0.080	1.051
<b>0.100</b>	<b>0.824</b>	<b>0.080</b>	<b>1.105</b>

**Supplementary Table 4.2. Number of functional SNPs, by effect type(s), overlapping at least 1 ChIP peak.** Functional SNPs were subsetted based on unique patterns of which effect(s) were significant, and filtered to those overlapping at least one of the below features. Column names indicate motifs/TF binding sites: DR\_5 retinoic acid response element predictions<sup>74</sup>; ChIP peak(s) for RARA, RARB, RARG RXRA, RXRB, THRA, or THRB<sup>135,140</sup>, and/or a VDR consensus ChIP region from lymphoblast cell lines<sup>73</sup>.

Significant Effect(s)	DR0-4 RARE hg19 FIMO	DR5 RARE	RARA	RARA (RA- treated liver)	RARB	RARB (RA- treated liver)	RARG	RXRA	RXRA (RA- treated liver)	RXRB	THRA	THRB	VDR
First Assay Hit, First and VEH Hit, Second Assay Allele Effect (n=4 SNPs)	0	0	3	1	0	2	0	2	2	2	1	0	1
<b>Second Assay Interaction Effect (n=16 SNPs)</b>	<b>0</b>	<b>1</b>	<b>6</b>	<b>5</b>	<b>1</b>	<b>6</b>	<b>1</b>	<b>14</b>	<b>6</b>	<b>7</b>	<b>4</b>	<b>2</b>	<b>2</b>
First Assay Hit (n=5 SNPs)	0	0	3	2	1	3	1	3	1	1	0	0	0
Second Assay Allele Effect (n=8 SNPs)	0	0	2	2	0	1	1	3	1	3	0	0	1
First Assay Hit, Second Assay Allele Effect (n=1 SNP)	0	0	1	0	0	0	0	0	0	0	0	0	0
<b>Second Assay Allele Effect, Second Assay Interaction Effect (n=1 SNP)</b>	<b>0</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
First Assay Hit, First and VEH Hit (n=1 SNP)	0	0	1	1	1	1	0	1	1	0	0	0	0
<b>Totals</b>													
All Interaction SNPs ± other effects (n=17 SNPs)	0	1	7	5	2	7	1	15	7	7	4	2	2
Non-interaction SNPs (n=19 SNPs)	0	0	10	6	2	7	2	9	5	6	1	0	2

## 4.7 References

1. Smith, K. Mental health: A world of depression. *Nature* **515**, 180–181 (2014).
2. Üstün, T. B., Ayuso-Mateos, J. L., Chatterji, S., Mathers, C. & Murray, C. J. L. Global burden of depressive disorders in the year 2000. *Brit J Psychiat* **184**, 386–392 (2004).
3. Greenberg, P. E., Fournier, A.-A., Sisitsky, T., Pike, C. T. & Kessler, R. C. The Economic Burden of Adults With Major Depressive Disorder in the United States (2005 and 2010). *J Clin Psychiatry* **76**, 155–162 (2015).
4. Papakostas, G. I. & Fava, M. Does the probability of receiving placebo influence clinical trial outcome? A meta-regression of double-blind, randomized clinical trials in MDD. *European Neuropsychopharmacol J European Coll Neuropsychopharmacol* **19**, 34–40 (2009).
5. Flint, J. & Kendler, K. S. The genetics of major depression. *Neuron* **81**, 484–503 (2014).
6. Corfield, E. C., Yang, Y., Martin, N. G. & Nyholt, D. R. A continuum of genetic liability for minor and major depression. *Transl Psychiat* **7**, e1131–e1131 (2017).
7. Clements, C. C. *et al.* Genome-wide association study of patients with a severe major depressive episode treated with electroconvulsive therapy. *Mol Psychiatr* 1–11 (2021) doi:10.1038/s41380-020-00984-0.
8. Wray, N. R. *et al.* Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nat Genet* **50**, 668–681 (2018).
9. Cai, N. *et al.* Sparse whole-genome sequencing identifies two loci for major depressive disorder. *Nature* **523**, 588–591 (2015).
10. Howard, D. M. *et al.* Genome-wide meta-analysis of depression identifies 102 independent variants and highlights the importance of the prefrontal brain regions. *Nat Neurosci* **22**, 343–352 (2019).
11. Levey, D. F. *et al.* Bi-ancestral depression GWAS in the Million Veteran Program and meta-analysis in >1.2 million individuals highlight new therapeutic directions. *Nat Neurosci* 1–10 (2021) doi:10.1038/s41593-021-00860-2.
12. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Sci New York N Y* **337**, 1190–5 (2012).
13. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–30 (2015).



14. Schaub, M. A., Boyle, A. P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease associations with regulatory information in the human genome. *Genome Res* **22**, 1748–1759 (2012).
15. Chen, J. *et al.* A uniform survey of allele-specific binding and expression over 1000-Genomes-Project individuals. *Nat Commun* **7**, 11101 (2016).
16. Shi, H., Kichaev, G. & Pasaniuc, B. Contrasting the Genetic Architecture of 30 Complex Traits from Summary Association Data. *Am J Hum Genetics* **99**, 139–153 (2016).
17. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits: From Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).
18. Corradin, O. *et al.* Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome Res* **24**, 1–13 (2014).
19. Consortium, B. W. G. of the P. G. *et al.* Conditional GWAS analysis to identify disorder-specific SNPs for psychiatric disorders. *Mol Psychiatr* 1–12 (2020) doi:10.1038/s41380-020-0705-9.
20. Li, S. *et al.* Regulatory mechanisms of major depressive disorder risk variants. *Mol Psychiatr* 1–20 (2020) doi:10.1038/s41380-020-0715-7.
21. Furlong, L. I. Human diseases through the lens of network biology. *Trends Genet* **29**, 150–159 (2013).
22. Liu, X., Li, Y. I. & Pritchard, J. K. Trans Effects on Gene Expression Can Drive Omnigenic Inheritance. *Cell* **177**, 1022-1034.e6 (2019).
23. Leeuw, C. A. de, Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: Generalized Gene-Set Analysis of GWAS Data. *Plos Comput Biol* **11**, e1004219 (2015).
24. Sey, N. Y. A. *et al.* A computational tool (H-MAGMA) for improved prediction of brain-disorder risk genes by incorporating brain chromatin interaction profiles. *Nat Neurosci* 1–11 (2020) doi:10.1038/s41593-020-0603-0.
25. Won, H. *et al.* Chromosome conformation elucidates regulatory relationships in developing human brain. *Nature* **538**, 523–527 (2016).
26. Torre-Ubieta, L. de la *et al.* The Dynamic Landscape of Open Chromatin during Human Cortical Neurogenesis. *Cell* **172**, 289-304.e18 (2018).
27. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat Genet* **48**, 245–52 (2016).

28. Gamazon, E. R., Zwinderman, A. H., Cox, N. J., Denys, D. & Derks, E. M. Multi-tissue transcriptome analyses identify genetic mechanisms underlying neuropsychiatric traits. *Nat Genet* **51**, 933–940 (2019).
29. Gerring, Z. F., Mina-Vargas, A., Gamazon, E. R. & Derks, E. M. E-MAGMA: an eQTL-informed method to identify risk genes using genome-wide association study summary statistics. *Bioinformatics* btab115- (2021) doi:10.1093/bioinformatics/btab115.
30. Gerring, Z. F., Gamazon, E. R., Derks, E. M. & Consortium, M. D. D. W. G. of the P. G. A gene co-expression network-based analysis of multiple brain tissues reveals novel genes and molecular pathways underlying major depression. *Plos Genet* **15**, e1008245 (2019).
31. Xu, X., Wells, A. B., O'Brien, D. R., Nehorai, A. & Dougherty, J. D. Cell type-specific expression analysis to identify putative cellular mechanisms for neurogenetic disorders. *J Neurosci Official J Soc Neurosci* **34**, 1420–31 (2014).
32. Ghosh, J. C. *et al.* Interactions that determine the assembly of a retinoid X receptor/corepressor complex. *Proc National Acad Sci* **99**, 5842–5847 (2002).
33. Mangelsdorf, D. J. & Evans, R. M. The RXR heterodimers and orphan receptors. *Cell* **83**, 841–850 (1995).
34. Liao, W.-L. *et al.* Modular patterning of structure and function of the striatum by retinoid receptor signaling. *Proc National Acad Sci* **105**, 6765–6770 (2008).
35. Bremner, J. D., Shearer, K. D. & McCaffery, P. J. Retinoic acid and affective disorders: the evidence for an association. *J Clin Psychiatry* **73**, 37–50 (2011).
36. Grøntved, L. *et al.* Transcriptional activation by the thyroid hormone receptor through ligand-dependent receptor recruitment and chromatin remodelling. *Nat Commun* **6**, 7048 (2015).
37. Kim, S.-H. *et al.* Anterior insula-associated social novelty recognition: orchestrated regulation by a local retinoic acid cascade and oxytocin signaling. *Biorxiv* 2021.01.15.426848 (2021) doi:10.1101/2021.01.15.426848.
38. Hu, P. *et al.* Chronic retinoic acid treatment suppresses adult hippocampal neurogenesis, in close correlation with depressive-like behavior. *Hippocampus* **26**, 911–923 (2016).
39. Chen, X.-N. *et al.* The Involvement of Retinoic Acid Receptor- $\alpha$  in Corticotropin-Releasing Hormone Gene Expression and Affective Disorders. *Biol Psychiat* **66**, 832–839 (2009).
40. Hu, P. *et al.* All-trans retinoic acid-induced hypothalamus–pituitary–adrenal hyperactivity involves glucocorticoid receptor dysregulation. *Transl Psychiat* **3**, e336–e336 (2013).
41. Consortium, C.-D. G. of the P. G. *et al.* Genomic Relationships, Novel Loci, and Pleiotropic Mechanisms across Eight Psychiatric Disorders. *Cell* **179**, 1469–1482.e11 (2019).

42. Reay, W. R. *et al.* Polygenic disruption of retinoid signalling in schizophrenia and a severe cognitive deficit subtype. *Mol Psychiatr* **25**, 719–731 (2018).
43. Regen, F. *et al.* Clozapine modulates retinoid homeostasis in human brain and normalizes serum retinoic acid deficit in patients with schizophrenia. *Mol Psychiatr* 1–12 (2020) doi:10.1038/s41380-020-0791-8.
44. Gandal, M. J. *et al.* Shared molecular neuropathology across major psychiatric disorders parallels polygenic overlap. *Sci New York N Y* **359**, 693–697 (2018).
45. Guo, M. *et al.* Vitamin A improves the symptoms of autism spectrum disorders and decreases 5-hydroxytryptamine (5-HT): A pilot study. *Brain Res Bull* **137**, 35–40 (2018).
46. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat Biotechnol* **30**, 271–277 (2012).
47. Kwasnieski, J. C., Mogno, I., Myers, C. A., Corbo, J. C. & Cohen, B. A. Complex effects of nucleotide variants in a mammalian cis-regulatory element. *P Natl Acad Sci Usa* **109**, 19498–503 (2012).
48. Patwardhan, R. P. *et al.* High-resolution analysis of DNA regulatory elements by synthetic saturation mutagenesis. *Nat Biotechnol* **27**, 1173–1175 (2009).
49. Mulvey, B., Lagunas, T. & Dougherty, J. D. Massively Parallel Reporter Assays: Defining Functional Psychiatric Genetic Variants across Biological Contexts. *Biol Psychiat* (2020) doi:10.1016/j.biopsych.2020.06.011.
50. Choi, J. *et al.* Massively parallel reporter assays of melanoma risk variants identify MX2 as a gene promoting melanoma. *Nat Commun* **11**, 2718 (2020).
51. Bourges, C. *et al.* Resolving mechanisms of immune-mediated disease in primary CD4 T cells. *Embo Mol Med* **12**, e12112 (2020).
52. Ulirsch, J. C. *et al.* Systematic Functional Dissection of Common Genetic Variation Affecting Red Blood Cell Traits. *Cell* **165**, 1530–1545 (2016).
53. Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants using a Multiplexed Reporter Assay. *Cell* **165**, 1519–1529 (2016).
54. Myint, L. *et al.* A screen of 1,049 schizophrenia and 30 Alzheimer’s-associated variants for regulatory potential. *Am J Medical Genetics Part B Neuropsychiatric Genetics* **183**, 61–73 (2020).
55. Shen, S. Q. *et al.* Massively parallel cis -regulatory analysis in the mammalian central nervous system. *Genome Res* **26**, 238–255 (2016).

56. Shen, S. Q. *et al.* A candidate causal variant underlying both higher intelligence and increased risk of bipolar disorder. *Biorxiv* 580258 (2019) doi:10.1101/580258.
57. Vockley, C. M. *et al.* Massively parallel quantification of the regulatory effects of noncoding genetic variation in a human cohort. *Genome Res* **25**, 1206–1214 (2015).
58. Lu, X. *et al.* Global discovery of lupus genetic risk variant allelic enhancer activity. *Nat Commun* **12**, 1611 (2021).
59. Vockley, C. M. *et al.* Direct GR Binding Sites Potentiate Clusters of TF Binding across the Human Genome. *Cell* **166**, 1269-1281.e19 (2016).
60. Johnson, G. D. *et al.* Human genome-wide measurement of drug-responsive regulatory activity. *Nat Commun* **9**, 5317 (2018).
61. Shea, T. B., Fischer, I. & Sapirstein, V. S. Effect of retinoic acid on growth and morphological differentiation of mouse NB2a neuroblastoma cells in culture. *Dev Brain Res* **21**, 307–314 (1985).
62. Evangelopoulos, M. E., Weis, J. & Krüttgen, A. Signalling pathways leading to neuroblastoma differentiation after serum withdrawal: HDL blocks neuroblastoma differentiation by inhibition of EGFR. *Oncogene* **24**, 3309–3318 (2005).
63. Ochoa, D. *et al.* Open Targets Platform: supporting systematic drug–target identification and prioritisation. *Nucleic Acids Res* **49**, gkaa1027- (2020).
64. White, M. A., Myers, C. A., Corbo, J. C. & Cohen, B. A. Massively parallel in vivo enhancer assay reveals that highly local features determine the cis-regulatory function of ChIP-seq peaks. *Proc National Acad Sci* **110**, 11952–11957 (2013).
65. Donello, J. E., Loeb, J. E. & Hope, T. J. Woodchuck hepatitis virus contains a tripartite posttranscriptional regulatory element. *J Virol* **72**, 5085–92 (1998).
66. Wu, P.-Y. *et al.* Functional decreases in P2X7 receptors are associated with retinoic acid-induced neuronal differentiation of Neuro-2a neuroblastoma cells. *Cell Signal* **21**, 881–891 (2009).
67. Chohanadisai, W., Graham, D. M., Keen, C. L., Rucker, R. B. & Messerli, M. A. Neurulation and neurite extension require the zinc transporter ZIP12 (slc39a12). *Proc National Acad Sci* **110**, 9903–9908 (2013).
68. Tian, R. *et al.* CRISPR Interference-Based Platform for Multimodal Genetic Screens in Human iPSC-Derived Neurons. *Neuron* **104**, 239-255.e12 (2019).

69. Coetzee, S. G., Coetzee, G. A. & Hazelett, D. J. motifbreakR: an R/Bioconductor package for predicting variant effects at transcription factor binding sites. *Bioinform Oxf Engl* **31**, 3847–9 (2015).
70. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
71. Davis, C. A. *et al.* The Encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res* **46**, gkx1081- (2017).
72. Hinrichs, A. S. *et al.* The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res* **34**, D590–D598 (2006).
73. Tuoresmäki, P., Väisänen, S., Neme, A., Heikkinen, S. & Carlberg, C. Patterns of Genome-Wide VDR Locations. *Plos One* **9**, e96105 (2014).
74. Lalevée, S. *et al.* Genome-wide in Silico Identification of New Conserved and Functional Retinoic Acid Receptor Response Elements (Direct Repeats Separated by 5 bp). *J Biol Chem* **286**, 33322–33334 (2011).
75. Song, M. *et al.* Mapping cis-regulatory chromatin contacts in neural cells links neuropsychiatric disorder risk variants to target genes. *Nat Genet* **51**, 1252–1262 (2019).
76. Hauberg, M. E. *et al.* Large-Scale Identification of Common Trait and Disease Variants Affecting Gene Expression. *Am J Hum Genetics* **100**, 885–894 (2017).
77. Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the human brain. *Science* **362**, eaat8464 (2018).
78. Aguet, F. *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).
79. Ng, B. *et al.* An xQTL map integrates the genetic architecture of the human brain’s transcriptome and epigenome. *Nat Neurosci* **20**, 1418–1426 (2017).
80. Xie, Z. *et al.* Gene Set Knowledge Discovery with Enrichr. *Curr Protoc* **1**, e90 (2021).
81. Mi, H. *et al.* PANTHER version 16: a revised family classification, tree-based classification tool, enhancer regions and extensive API. *Nucleic Acids Res* **49**, gkaa1106- (2020).
82. Grote, S., Prüfer, K., Kelso, J. & Dannemann, M. ABAEnrichment: an R package to test for gene set expression enrichment in the adult and developing human brain. *Bioinformatics* **32**, 3201–3203 (2016).
83. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat Genet* **50**, 1335–1341 (2018).

84. Weissbrod, O. *et al.* Functionally informed fine-mapping and polygenic localization of complex trait heritability. *Nat Genet* **52**, 1355–1363 (2020).
85. Dobbyn, A. *et al.* Landscape of Conditional eQTL in Dorsolateral Prefrontal Cortex and Co-localization with Schizophrenia GWAS. *Am J Hum Genet* **102**, 1169–1184 (2018).
86. Walker, R. L. *et al.* Genetic Control of Expression and Splicing in Developing Human Brain Informs Disease Mechanisms. *Cell* **179**, 750–771.e22 (2019).
87. Hyde, C. L. *et al.* Identification of 15 genetic loci associated with risk of major depression in individuals of European descent. *Nat Genet* **48**, 1031–1036 (2016).
88. Ciuculete, D. M. *et al.* meQTL and ncRNA functional analyses of 102 GWAS-SNPs associated with depression implicate HACE1 and SHANK2 genes. *Clin Epigenetics* **12**, 99 (2020).
89. Cimadamore, F., Amador-Arjona, A., Chen, C., Huang, C.-T. & Terskikh, A. V. SOX2–LIN28/let-7 pathway regulates proliferation and neurogenesis in neural precursors. *Proc National Acad Sci* **110**, E3017–E3026 (2013).
90. Ong, K. K. *et al.* Genetic variation in LIN28B is associated with the timing of puberty. *Nat Genet* **41**, 729–733 (2009).
91. Perry, J. R. B. *et al.* Parent-of-origin-specific allelic associations among 106 genomic loci for age at menarche. *Nature* **514**, 92–97 (2014).
92. Corre, C. *et al.* Sex-specific regulation of weight and puberty by the Lin28/let-7 axis. *J Endocrinol* **228**, 179–191 (2016).
93. Mulvey, B. *et al.* Molecular and Functional Sex Differences of Noradrenergic Neurons in the Mouse Locus Coeruleus. *Cell Reports* **23**, 2225–2235 (2018).
94. Marcus, S. M. *et al.* Gender differences in depression: Findings from the STAR\*D study. *J Affect Disorders* **87**, 141–150 (2005).
95. Brody, D. J., Pratt, L. A. & Hughes, J. P. Prevalence of Depression Among Adults Aged 20 and Over: United States, 2013–2016. *Nchs Data Brief* 1–8 (2018).
96. Pan, Q. *et al.* VARAdb: a comprehensive variation annotation database for human. *Nucleic Acids Res* **49**, gkaa922- (2020).
97. Teixeira, J. R., Szeto, R. A., Carvalho, V. M. A., Muotri, A. R. & Papes, F. Transcription factor 4 and its association with psychiatric disorders. *Transl Psychiat* **11**, 19 (2021).
98. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24–26 (2011).

99. Liu, M. *et al.* Transcriptional Profiling Reveals a Common Metabolic Program in High-Risk Human Neuroblastoma and Mouse Neuroblastoma Sphere-Forming Cells. *Cell Reports* **17**, 609–623 (2016).
100. Korade, Ž., Kenworthy, A. K. & Mirnics, K. Molecular consequences of altered neuronal cholesterol biosynthesis. *J Neurosci Res* **87**, 866–875 (2009).
101. Rao, C., Malaguti, M., Mason, J. O. & Lowell, S. The transcription factor E2A drives neural differentiation in pluripotent cells. *Development* **147**, dev.184093 (2020).
102. Imayoshi, I. & Kageyama, R. bHLH Factors in Self-Renewal, Multipotency, and Fate Choice of Neural Progenitor Cells. *Neuron* **82**, 9–23 (2014).
103. Ypsilanti, A. R. *et al.* Transcriptional Network Orchestrating Regional Patterning of Cortical Progenitors. *Biorxiv* 2020.11.03.366914 (2020) doi:10.1101/2020.11.03.366914.
104. Vierstra, J. *et al.* Global reference mapping of human transcription factor footprints. *Nature* **583**, 729–736 (2020).
105. Banerjee, D. *et al.* Lineage specific transcription factor waves reprogram neuroblastoma from self-renewal to differentiation. *Biorxiv* 2020.07.23.218503 (2020) doi:10.1101/2020.07.23.218503.
106. Tremblay, R. G. *et al.* Differentiation of mouse Neuro 2A cells into dopamine neurons. *J Neurosci Meth* **186**, 60–67 (2010).
107. Alonso, A. *et al.* Development of the serotonergic cells in murine raphe nuclei and their relations with rhombomeric domains. *Brain Struct Funct* **218**, 1229–1277 (2013).
108. Nishida, Y. *et al.* Identification and classification of genes regulated by phosphatidylinositol 3-kinase- and TRKB-mediated signalling pathways during neuronal differentiation in two subtypes of the human neuroblastoma cell line SH-SY5Y. *Bmc Res Notes* **1**, 95 (2008).
109. Marteyn, A. *et al.* Mutant Human Embryonic Stem Cells Reveal Neurite and Synapse Formation Defects in Type 1 Myotonic Dystrophy. *Cell Stem Cell* **8**, 434–444 (2011).
110. Johansson, P., Pavey, S. & Hayward, N. Confirmation of a BRAF mutation-associated gene expression signature in melanoma. *Pigm Cell Res* **20**, 216–221 (2007).
111. Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell* **174**, 999–1014.e22 (2018).
112. Ishikawa, J. *et al.* 13-cis-retinoic acid alters the cellular morphology of slice-cultured serotonergic neurons in the rat. *Eur J Neurosci* **27**, 2363–2372 (2008).

113. Schrode, N. *et al.* Synergistic effects of common schizophrenia risk variants. *Nat Genet* **51**, 1475–1485 (2019).
114. Mich, J. K. *et al.* Functional enhancer elements drive subclass-selective expression from mouse to primate neocortex. *Cell Reports* **34**, 108754 (2021).
115. Fee, C., Banasr, M. & Sibille, E. Somatostatin-Positive Gamma-Aminobutyric Acid Interneuron Deficits in Depression: Cortical Microcircuit and Therapeutic Perspectives. *Biol Psychiat* **82**, 549–559 (2017).
116. Smazal, A. L. & Schalinske, K. L. Oral Administration of Retinoic Acid Lowers Brain Serotonin Concentration in Rats. *Faseb J* **27**, 635.6-635.6 (2013).
117. Ren, H. *et al.* Genes associated with anhedonia: a new analysis in a large clinical trial (GENDEP). *Transl Psychiat* **8**, 150 (2018).
118. Li, X. *et al.* Common variants on 6q16.2, 12q24.31 and 16p13.3 are associated with major depressive disorder. *Neuropsychopharmacol* **43**, 2146–2153 (2018).
119. Power, R. A. *et al.* Genome-wide Association for Major Depression Through Age at Onset Stratification: Major Depressive Disorder Working Group of the Psychiatric Genomics Consortium. *Biol Psychiat* **81**, 325–335 (2017).
120. Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder. *Nat Genet* **51**, 431–444 (2019).
121. Smith, D. J. *et al.* Genome-wide analysis of over 106 000 individuals identifies 9 neuroticism-associated loci. *Mol Psychiatr* **21**, 1644 (2016).
122. Luciano, M. *et al.* Association analysis in over 329,000 individuals identifies 116 independent variants influencing neuroticism. *Nat Genet* **50**, 6–11 (2017).
123. Ward, J. *et al.* The genomic basis of mood instability: identification of 46 loci in 363,705 UK Biobank participants, genetic correlation with psychiatric disorders, and association with gene expression and function. *Mol Psychiatr* 1–9 (2019) doi:10.1038/s41380-019-0439-8.
124. Meier, S. M. *et al.* Genetic Variants Associated With Anxiety and Stress-Related Disorders. *Jama Psychiat* **76**, 924–932 (2019).
125. Demontis, D. *et al.* Discovery of the first genome-wide significant risk loci for attention deficit/hyperactivity disorder. *Nat Genet* **51**, 63–75 (2018).
126. Hill, W. D. *et al.* A combined analysis of genetically correlated traits identifies 187 loci and a role for neurogenesis and myelination in intelligence. *Mol Psychiatr* **24**, 169–181 (2018).



127. Lee, J. J. *et al.* Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nat Genet* **50**, 1112–1121 (2018).
128. Carvalho-Silva, D. *et al.* Open Targets Platform: new developments and updates two years on. *Nucleic Acids Res* **47**, gky1133 (2018).
129. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using RegulomeDB. *Genome Res* **22**, 1790–1797 (2012).
130. Machiela, M. J. & Chanock, S. J. LDlink: a web-based application for exploring population-specific haplotype structure and linking correlated alleles of possible functional variants. *Bioinformatics* **31**, 3555–3557 (2015).
131. Jaffe, A. E. *et al.* Profiling gene expression in the human dentate gyrus granule cell layer reveals insights into schizophrenia and its genetic risk. *Nat Neurosci* **23**, 510–519 (2020).
132. Consortium, T. B. *et al.* Developmental and genetic regulation of the human cortex transcriptome illuminate schizophrenia pathogenesis. *Nat Neurosci* **21**, 1117–1125 (2018).
133. Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455–61 (2014).
134. (DGT), F. C. and the R. P. and C. *et al.* A promoter-level mammalian expression atlas. *Nature* **507**, 462–70 (2014).
135. Consortium, T. E. P. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699–710 (2020).
136. Amiri, A. *et al.* Transcriptome and epigenome landscape of human cortical development modeled in organoids. *Sci New York N Y* **362**, eaat6720 (2018).
137. Kreimer, A., Yan, Z., Ahituv, N. & Yosef, N. Meta-analysis of massively parallel reporter assays enables prediction of regulatory function across cell types. *Hum Mutat* **40**, 1299–1313 (2019).
138. Ghazi, A. R. *et al.* Design tools for MPRA experiments. *Bioinform Oxf Engl* **34**, 2682–2683 (2018).
139. Lund, A. H., Duch, M. & Pedersen, F. S. Increased Cloning Efficiency by Temperature-Cycle Ligation. *Nucleic Acids Res* **24**, 800–801 (1996).
140. He, Y., Tsuei, J. & Wan, Y.-J. Y. Biological functional annotation of retinoic acid alpha and beta in mouse liver based on genome-wide binding. *Am J Physiol-gastr L* **307**, G205–G218 (2014).

## **Chapter 5: Sex significantly impacts the function of major depression-linked variants *in vivo***

This chapter has been preprinted on *Biorxiv* and is undergoing peer review at *Nature Neuroscience*.

The citation for the former is:

Mulvey, B., Selmanovic, D. & Dougherty, J. D. Sex significantly impacts the function of major depression-linked variants *in vivo*. *Biorxiv* 2021.11.01.466849 (2021) doi:10.1101/2021.11.01.466849.

Genome-wide association studies have discovered blocks of common variants—likely transcriptional-regulatory—associated with major depressive disorder (MDD), though the functional subset and their biological impacts remain unknown. Likewise, why depression occurs ‘associated functional variants interact with sex and produce greater impact in female brains. I developed methods to directly measure regulatory variant activity and sex interactions using massively parallel reporter assays (MPRAs) in the mouse brain *in vivo*, in a cell type-specific manner. I measured activity of >1,000 variants from >30 MDD loci, identifying extensive sex-by-allele effects in mature hippocampal neurons and suggesting sex-differentiated impacts of genetic risk may underlie sex bias in disease. Unbiased informatics approaches indicated that functional MDD variants recurrently disrupt sex hormone receptor binding sequences. I confirmed this with MPRAs in neonatal brains, comparing brains undergoing a sex-differentiating hormone surge to hormonally-quiescent juveniles. my study provides novel insights into the influence of age, biological sex, and cell type on regulatory-variant function, and provides a framework for *in vivo* parallel assays to functionally define interactions between organismal variables like sex and regulatory variation.

## 5.1 Introduction

Major depressive disorder (MDD) is a profoundly disruptive and sometimes lethal disorder, affecting women 2-3 times more frequently than men across countries and cultures<sup>1</sup>. Sex differences are present across multiple levels of the disease, from symptom profiles<sup>2</sup> and effective drug classes<sup>3</sup> to brain-wide gene expression<sup>4,5</sup>. Genome-wide association studies (GWASes) have identified dozens of linkage regions each containing numerous single-nucleotide polymorphisms (SNPs) associated with MDD, demonstrating its heritability<sup>6-8</sup>. More recently, sex-by-genotype (SxG) analyses of large GWAS cohorts have revealed that MDD risk loci are the same for men and women, yet these loci explain up to 4-fold greater MDD heritability in females<sup>9,10</sup>. These findings suggest that sex interacts with a common pool of SNPs to attenuate or amplify the MDD risk they confer. However, disease-associated SNPs are seldom found in protein-coding space, complicating prediction of their molecular consequences. Instead, these SNPs are found in probable regulatory elements (REs), including transcriptional-regulatory sequences predicted from measures such as chromatin marks, accessibility, and conformation. Specific brain regions and cell types are enriched for such measures at—and in putative regulatory target genes of—MDD-associated loci, including the hippocampus<sup>7,11</sup> and excitatory neurons<sup>12-17</sup>, suggesting sites of action for these REs. In particular, there has been long-standing interest in the hippocampus regarding both MDD pathology and sex differences in the brain. Hippocampal volume reductions in MDD patients have been widely reported<sup>18</sup>. Moreover, the hippocampus is subject to influences of sex from perinatal<sup>19,20</sup> to adult life, presenting in MDD as sex differences in hippocampal volume loss<sup>21</sup> and gene expression<sup>4</sup>.

Determining the identity of functional SNPs from MDD-associated regions is the first key step toward understanding the biological perturbations resulting from risk genotypes, which can in turn enable inference of dysregulated target genes and shared regulatory programs involved across loci. However, studies connecting MDD-associated SNPs to gene expression in brain tissue, even those that considered sex effects<sup>22</sup>, have been limited to indirect indicators of function (e.g., chromatin state), or are confounded by linkage disequilibrium (e.g., for expression quantitative trait loci (eQTLs)). In contrast, *direct measurement* of regulatory output of common variants associated with disease has largely been restricted to the *in vitro* setting. Large-scale *in vitro* identification of functional regulatory variants has been made possible by massively parallel reporter assays (MPRAs), a method for functionally detecting activity from thousands of REs (and their variants) simultaneously. In brief, MPRAs adapt a traditional reporter assay paradigm—placing REs upstream of an optically measured reporter (e.g., luciferase)—but adds a unique, RE-identifying “barcode” sequence to the reporter’s 3’ untranslated region, enabling quantification of activity for thousands of REs simultaneously by RNA barcode sequencing. MPRAs have enabled identification of trait- and disease-associated SNPs affecting REs in culturable, disease-relevant cell types *in vitro*<sup>23–26</sup>. However, the complexities of cell types interacting in the brain and of the sex hormonal milieu cannot readily be emulated *ex vivo*.

I overcome these prior limitations in functional regulatory SNP (rSNP) identification to interrogate the biological contexts under which MDD rSNPs act by delivering an MPRA library of MDD-associated variants<sup>27</sup> into the adult mouse hippocampus *in vivo*. Building on prior brain MPRAs of enhancers<sup>28</sup> and an RE variant<sup>29</sup>, my approach greatly extends *in vivo* MPRA methods to identify rSNPs and their sex interactions, including those which are cell type-specific. First, I combined

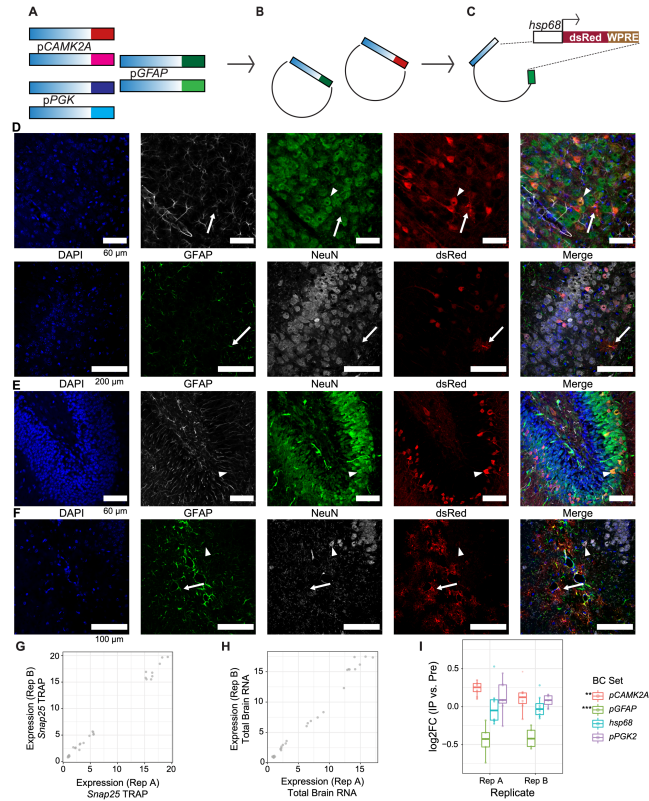
MPRA with translating ribosome affinity purification (TRAP) to simultaneously identify MDD rSNPs in both excitatory neurons and the broader hippocampus. Further, these experiments utilized mice of both sexes, enabling us to additionally test the hypothesis that rSNPs are subject to sex-by-genotype (SxG) interactions. Finally, to characterize the potential role of circulating hormones in sex-differentiated rSNP activity, and to functionally replicate predicted fetal brain RE enrichments suggesting a role for MDD SNPs during circuit organization<sup>15,30-33</sup>, I likewise delivered the library to the mouse brain *in utero*. This allowed us to identify rSNPs neonatally, coinciding with a testosterone surge and critical period for establishing sex-specific brain circuitry<sup>34</sup>, and test for loss of SxG effects in juveniles, when hormonal influences are quiescent. In sum, I illustrate that MPRA can be leveraged *in vivo* to directly identify not only functional variants, but their context-dependence on age, sex, and cell type, while demonstrating that all three of these factors have substantial impacts on MDD-associated regulatory variation.

## 5.2 Results

### 5.2.1 Combining Translating Ribosome Affinity Purification (TRAP) with MPRA

Adeno-associated viruses (AAVs) have long been used to evaluate the activity of single REs in the brain, most recently in MPRA-like designs to screen multiple REs in parallel with RNA-seq-based quantification<sup>28,35,36</sup>. This suggests it may also be possible to adapt AAV-MPRAs to study functional consequences of RE variants associated with disease. As MDD genetic risk is enriched in neuronal REs, I tested the feasibility of combining a cell type-specific profiling method, TRAP, with MPRA to attain measurements of RE activity specifically in neurons. I generated 4 small AAV9 libraries (**Figure 5.1A-C**) expressing dsRed under the control of the *hsp68* minimal promoter, with full-length human promoters (denoted *pGene*) with documented expression in

neurons (human *pCAMK2A*), astrocytes (*pGFAP*), or all cells (*pPGK2*), each carrying unique 3' untranslated region (UTR) barcodes for quantification by RNA-seq. I first individually confirmed cell-type specificity by immunofluorescence (IF) (**Figure 5.1D-F**), and then delivered<sup>37</sup> a mixed pool of the four barcoded AAV9 libraries into the brain of postnatal day 2 (P2) neuron-specific TRAP mice (*Snap25-Rpl10a-eGFP*)<sup>38</sup>. TRAP was then used to compare total brain and neuronal activity levels of the three promoters and the *hsp68* promoter alone. I prepared MPRA sequencing libraries from 1) the delivered AAV pool (DNA), 2) total brain (input) RNA, and 3) TRAP (neuronal) RNA and assessed expression of each barcode by calculating a ratio of RNA counts to DNA counts. Results were highly replicable at the level of barcodes (**Figure 5.1G-H**) in both RNA fractions. Moreover, neuronal TRAP fractions demonstrated increased expression of barcodes driven by *pCAMK2A* and lower expression of *pGFAP*-paired barcodes (**Figure 5.1I**). These results demonstrated that cell type-specific effects, even of relatively small magnitudes, can be detected using a combined MPRA-TRAP approach. I then turned to applying this method to test the effects of human variants associated with psychiatric disease.



**Figure 5.1. Proof-of-principle for cell type-specific MPRA *in vivo*.** **A)** Cell type-specific promoters, *pGFAP* (astrocytic) and *pCAMK2A* (neuronal) were barcoded by PCR. **B)** Amplicons were cloned into an AAV plasmid. **C)** Further restriction cloning added a reporter cassette containing a minimal *hsp68* promoter, dsRed, and an RNA-stabilizing 3' UTR hepatitis E "woodchuck" (WPRE) element. Barcoded pools with each promoter were packaged into AAV9 separately. **D-F)** IF of P27 mouse brain after P2 injection with a single AAV9 barcode pool. (D) *hsp68* promoter alone preferentially drove dsRed expression in neurons, while (E) *pCAMK2A* drove reporter expression solely in neurons, and (F) *pGFAP* drove predominantly astrocytic dsRed expression. **G)** Replicability of barcode expression for the SNAP25-TRAP RNA fraction (Pearson's  $r=0.9975$ ) and **H)** total brain tissue (input) ( $r=0.9927$ ). **I)** TRAP expression, compared to total brain expression, was higher for *pCAMK2A* and lower for *pGFAP*, as expected.  $**p \leq 5 \times 10^{-4}$ ,  $***p \leq 5 \times 10^{-8}$

## 5.2.2 Identification of rSNPs and their sex-allele interactions in total hippocampus and excitatory neurons

Given the association of MDD with sex, hippocampal pathology, and neuronal genetics, I sought to identify regulatory SNPs among variants selected from broad linkage disequilibrium (LD) regions associated with MDD in total hippocampus. To further investigate the roles of SNPs and

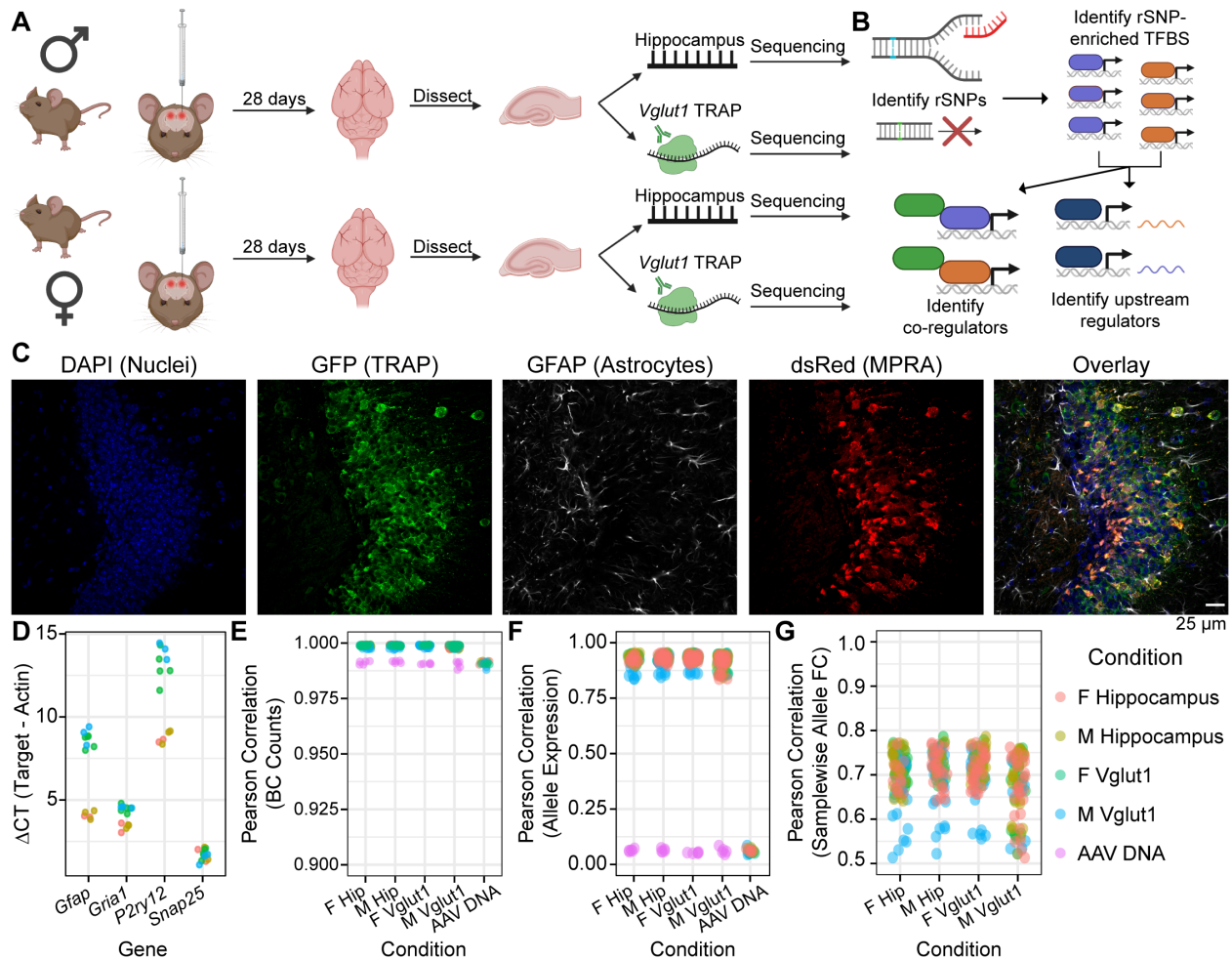
sex in hippocampal excitatory neurons in particular, I performed these experiments in a cross of *Slc17a7* (or, *Vglut1*)-*Cre* mice<sup>39</sup> to a *Cre* recombinase-dependent TRAP mouse line<sup>38</sup>.

The analyzed MPRA library covered 40 GWAS loci spanning ~1,000 SNPs in LD  $R^2 > 0.1$  with MDD-associated tag variants. SNPs were prioritized by their overlap with human brain and neural cell type eQTLs, histone marks, enhancer RNA overlap, and chromatin contacts (see *Data Availability*). Of these SNPs, 926 were from 29 MDD GWAS loci<sup>6,7,33,40-44</sup>, 19 were from 2 loci identified by meta-analysis of MDD and autism spectrum disorders<sup>45</sup>, and 21 were from 4 loci for MDD-correlated traits (mood instability, anxiety, and neuroticism)<sup>46-49</sup>. 126bp human genomic sequences centered on the SNP were generated for each allele and paired to 10 unique 10bp barcodes per allelic sequence for internal replication (**Supplementary Figure 5.1**). These were inserted into an AAV plasmid, then cloned to contain an *hsp68* minimal promoter<sup>50</sup> driving dsRed along with the “woodchuck” hepatitis 3’UTR element to improve recovery of reporter RNAs<sup>51</sup>, and packaged into AAV9.

I delivered the AAV9 library bilaterally into the hippocampus of *Vglut1*-TRAP mice ages P60-P80 (n=6 per sex), followed by hippocampal dissection and TRAP (**Figure 5.2A**) to identify rSNPs and their shared regulatory features (**Figure 5.2B**). IF of hippocampi from additional mice confirmed robust hippocampal expression of the dsRed reporter 28 days after injection (**Figure 5.2C**). I first confirmed by qPCR that TRAP RNA was depleted for glial marker genes several-fold as expected (**Figure 5.2D**), then conducted MPRA sequencing (**Supplementary Table 5.1**) and analysis of both input RNA (hippocampus) and TRAP RNA (*Vglut1*<sup>+</sup>) for each sample.



After quality control,  $n=5$  male and  $n=5$  female hippocampal samples (both the input and TRAP RNAs) were retained for MPRA analysis covering  $\sim 1,000$  SNPs. These showed replicability at the level of barcode counts (pairwise Pearson's  $r$  0.7-0.89) (**Figure 5.2E**), RE-level expression ( $r$  0.85-0.96) (**Figure 5.2F**), and mean allelic fold changes between conditions ( $r$  0.90-0.94) (sample-pairwise shown in **Figure 5.2G**), confirming I was able to reliably measure SNP-mediated regulatory effects and differences from a defined cell type *in vivo*.



**Figure 5.2. Experimental design, analysis plan, and quality control: adult mouse hippocampus and its excitatory neurons. A)** Adult male and female C57BL/6J mice received bilateral stereotaxic injections into the hippocampus delivering the AAV9-packaged MPRA. TRAP yielded two RNA fractions per sample: “input” (total hippocampal) and TRAP (Vglut1<sup>+</sup>). **B)** Analyses identified regulatory SNPs (rSNPs), transcription factor (TF) binding sites (TFBSes) enriched at rSNPs, shared protein interactors among these TFs, and shared regulators of these TFs’ expression. **C)** IF of Vglut1-TRAP mouse hippocampus 28 days after MPRA-AAV9 delivery, illustrating strong TRAP (GFP) co-expression with dsRed reporter, confirming RNA from the

latter is present in the cell type of interest. **D)** qPCR confirmed depletion of glial genes (*Gfap*, *P2ry12*) and modest enrichment of excitatory neuron marker *Gria1* in *Vglut1*<sup>+</sup> RNA. **E)** Barcode count correlations between replicates. Each point represents the cross-correlation between one sample of the type on the x-axis and one of the color-coded type. **F)** Correlation of mean barcode expression between replicates. **G)** Correlation of samplewise allelic differences in expression. PCC: Pearson correlation coefficient (or Pearson's *r*).

### 5.2.3 rSNPs in the adult hippocampus and hippocampal *Vglut1*<sup>+</sup> neurons

I first assessed ~1,000 SNPs for allelic effects in each individual sex and RNA fraction, using linear mixed models (LMMs) fitting barcode expression as a function of allele with random barcode effects (**Supplementary Figures 5.2-5.4**). I calculated empirical *p* ( $p_{\text{emp}}$ ) values using 50,000 simulated 'allelic' comparisons between subsets of random barcodes coupled to the minimal promoter alone<sup>27,52</sup> to account for technical and barcode-mediated noise. Significance was called at FDR-corrected  $p_{\text{emp}} < 0.2$ , a stringency comparable to a recent study of sex-interacting eQTLs<sup>22</sup>. Data provided at the Bitbucket link in *Acknowledgments* provides allelic beta (log2FC) values and significance status for variants at 10 different significance thresholds, while **Supplementary Figure 5.5** compares FDRs from standard and empirical *p* values in each experiment. For female mice, I additionally analyzed a separate cohort of *n*=3 *Vglut1*<sup>+</sup> TRAP mice delivered the same library into the same hippocampal coordinates.

In total hippocampus, I identified 36 (male) and 31 (female) rSNPs, 34 and 31 of which were from MDD loci, respectively. While male and female total hippocampi had similar numbers of rSNPs, I observed a striking sex difference in the number of rSNPs in *Vglut1*<sup>+</sup> cells specifically—only 7 (male) compared to 58 (female), indicating that within excitatory neurons, a higher proportion of MDD SNPs have discernible allelic effects in females. Moreover, all 7 male rSNPs were also functional in females. Female *Vglut1*<sup>+</sup> and total hippocampal SNP effects were consistent in magnitude and direction between this sample set and the additional (*n*=3) TRAP cohort (Pearson's

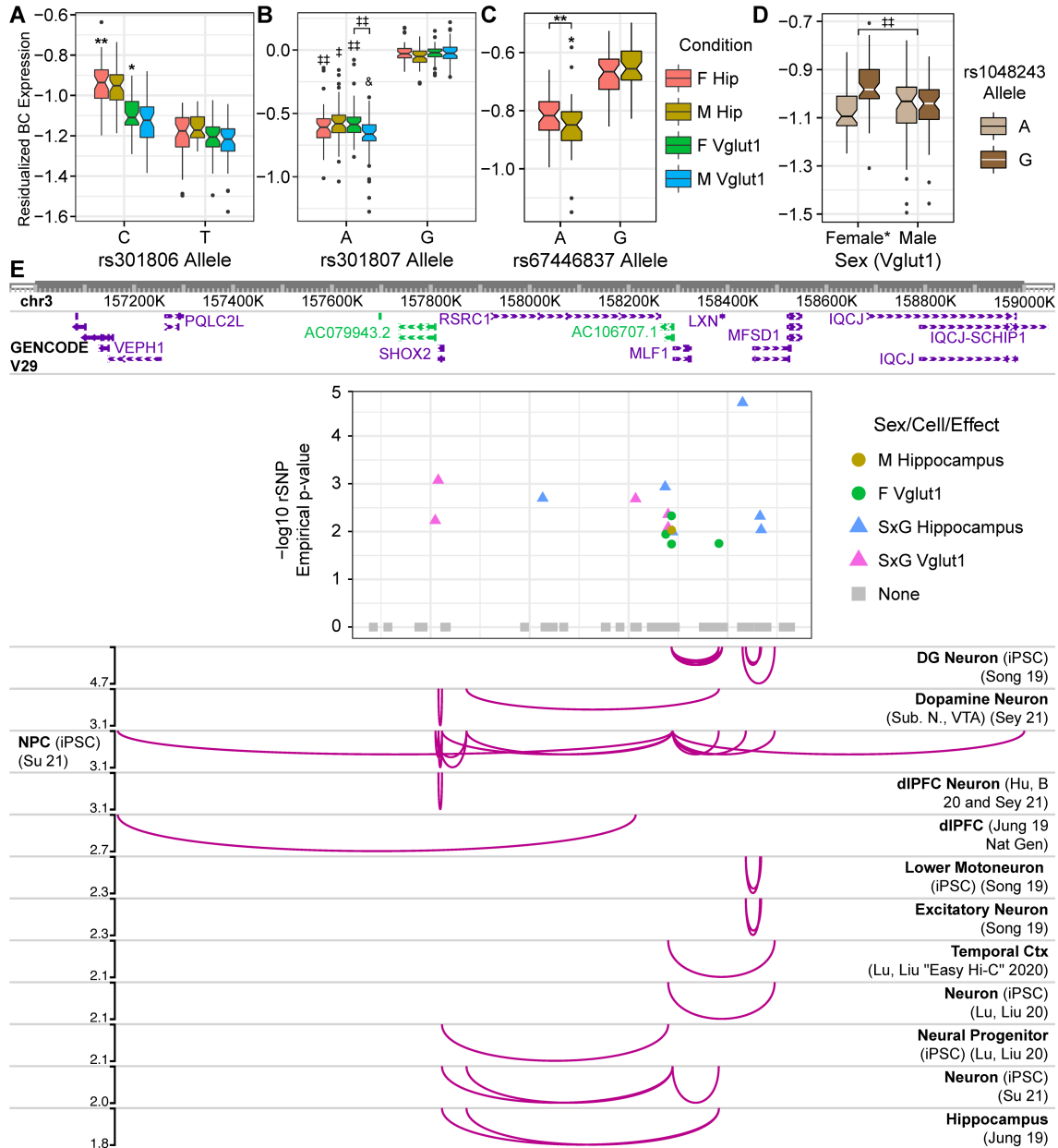
$r$  0.61-0.77) (**Supplementary Figure 5.6**). Notable rSNPs from  $\geq 1$  condition included rs2563323 and rs250427, putative brain and hippocampal<sup>53</sup> eQTLs for *SRAI*, a noncoding RNA which activates nuclear receptors even in the absence of their ligand<sup>54</sup>. Also notable were rs301806—an MDD GWAS index SNP<sup>41</sup>—and rs301807 (**Figure 5.3A-B**), both of which likely regulate the nearby gene *RERE*<sup>15,55</sup>.

#### 5.2.4 Sex-interacting rSNPs in hippocampus

Given the role of sex in MDD risk and the observed differences in rSNP activity in the adult hippocampus, I sought to investigate whether sex interacts with MDD risk genotypes. To ensure there were not confounding sex differences in minimal promoter activity, I compared minimal promoter-only barcode expression between sexes for the two tissue fractions, finding no sex difference in activity ( $t$ -test of barcode expression,  $p > 0.5$  in both comparisons; **Supplementary Figure 5.7**).

I therefore performed combined-sex LMM analysis of the *Vglut1*<sup>+</sup> and hippocampus results, identifying 41 sex-allele interaction rSNPs in each. Notably, while only 1 SxG rSNP was shared between total hippocampus and *Vglut1*<sup>+</sup>, sex-interacting SNPs originated from the same GWAS loci across both analyses: 16 of the 18 *Vglut1*<sup>+</sup> SxG loci also contained hippocampal SxG rSNPs. The tag locus rs1193510<sup>40</sup>, for example, contained 11 sex-interacting rSNPs, including several unique to hippocampus (**Figure 5.3C**) or *Vglut1*<sup>+</sup> (**Figure 5.3D**). This region is rich in neural chromatin contacts<sup>56-61</sup>, implicating several target genes of the identified rSNPs. Interestingly, SxG effects in hippocampus and *Vglut1*<sup>+</sup> segregated into distinct portions of this LD region (**Figure 5.3E**). I additionally examined three SNPs in my assay that were recently reported significant in sex-genotype interaction GWASes of over 500 traits<sup>62</sup>. One of these reported

variants, rs2400075, was found to have several significant sex-interacting associations to body traits; this variant is located near *LIN28B*—which shows sex-differential expression in mouse norepinephrine neurons<sup>63</sup>—and was a near-significant SxG rSNP in *Vglut1*<sup>+</sup> ( $0.20 < \text{FDR} < 0.25$ ).



**Figure 5.3. Adult hippocampus rSNPs and complex context-dependent, polygenic architecture of the *RSRC1* locus.** Boxplots show single-barcode (BC) expression levels adjusted for random effects across analyzed replicates. Center bars: median; boxes: 25-75% quantile; whiskers: observations spanning box edges to  $\pm 1.5 \times \text{interquantile range (IQR)}$ ; single points: observed values outside whisker range. Notched regions span  $\pm 1.58 \times \text{IQR} / \sqrt{n}$  measurements, approximating the 95% confidence interval for comparing median BC expression<sup>64</sup>. rSNPs corresponding to *RERE* locus<sup>33</sup> tag variants **A**) rs301806 and **B**) rs301807 are

shown. **C-D**) Sex-interacting SNPs from the *RSRC1* locus<sup>40</sup> included (C) rs67446837 and (D) rs1048243. **E**) rSNPs identified in the *RSRC1* locus and their regulatory target genes in human tissues ascertained by Hi-C. A plot of rSNP effects, colored by most significant condition, is embedded with its x-axis in human genomic (hg19) coordinates; chromatin contacts between the SNPs and distal gene promoters are illustrated below<sup>56-61</sup>. Curve heights correspond to  $-\log_{10}p_{emp}$  for the plotted rSNP. iPSC: Induced pluripotent stem cell; DG: dentate gyrus; Ctx: cortex. \*:  $p_{emp}$ -derived FDR < 0.25; \*\*: <0.2; ‡: < 0.15; ‡‡: < 0.1; &: < 0.05.

### 5.2.5 Transcriptional-regulatory systems shared across hippocampal rSNPs

I next asked whether there were any shared transcriptional-regulatory mechanisms underlying MDD rSNP effects in the hippocampus. I tested whether rSNPs perturbed specific transcription factor (TF) binding motifs more frequently than expected by chance (defined by their rates in rSNPs vs. non-effect SNPs in the assay)<sup>27</sup>. I assessed motif disruptions using motifbreakR<sup>65</sup> and RSAT var-tools<sup>66</sup> defining rSNPs at a nominal LMM  $p_{emp}$  of 0.05. This resulted in sets of 80-110 MPRA-identified rSNPs per condition, ensuring adequate depth for enrichment analysis. To refine these results, I filtered the enriched TFs (FDR<0.05) to those with altered putative binding sites  $\geq 4$  rSNPs and expressed in Genotype-Tissue Expression atlas (v8) hippocampus in the corresponding sex.

Altogether, I identified 38 enriched TFs in male total hippocampus rSNPs and 19 in female. The hippocampal TFs identified were largely distinctive between sexes; for example, KLF family TFs were unique to male hippocampal rSNPs, while nuclear receptor (NR) TFs were mostly unique to female rSNPs (**Figure 5.4A**). Among *Vglut1*<sup>+</sup> rSNPs, I identified 8 TFs in female and 16 in male, many of which were shared (e.g., DLX1, POU3F1/2/3). *POU3F2* has been previously shown to be a highly centralized, cross-disorder hub gene in postmortem brain co-expression analysis by PsychENCODE<sup>67</sup>.

To understand integrative biological functions of these TF sets, I utilized the tool Enrichr<sup>68</sup> for each TF set, including both tissues per sex, to identify putative *upstream* regulators, co-interacting TFs (protein-protein interactions (PPIs)), enriched disease gene sets, ontologies, and brain regions expressing the TFs (full tables in the Bitbucket link under *Acknowledgments*). The most striking enrichments (> 25% of query TFs) for male glutamatergic and hippocampal TFs combined were for regulators of these TFs' expression, including *CREB1* (14/48), *BRCA1* (19/48), and *ZBTB7A* (12/48). Male hippocampal TFs were likewise enriched for several upstream regulators, including *TCF3* (8/38), *HDAC2* (8/38), and *ZMIZ1* (9/38). *ZMIZ1* has roles in coactivation of androgen receptor (*AR*)<sup>69</sup> as well as *SMAD3*<sup>70</sup>, consistent with enrichment of these TFs in *SMAD3* (6/38) and *AR* (7/38) PPIs. TFs from male glutamatergic neurons were enriched for four PPIs: *SMAD3* and *SMAD4* (both 5/16), consistent with MPRA signal from neuronal enrichment by TRAP, and more surprisingly, sex hormone receptors: estrogen receptor  $\alpha$  (*ESR1*; 4/16) and *AR* (3/16).

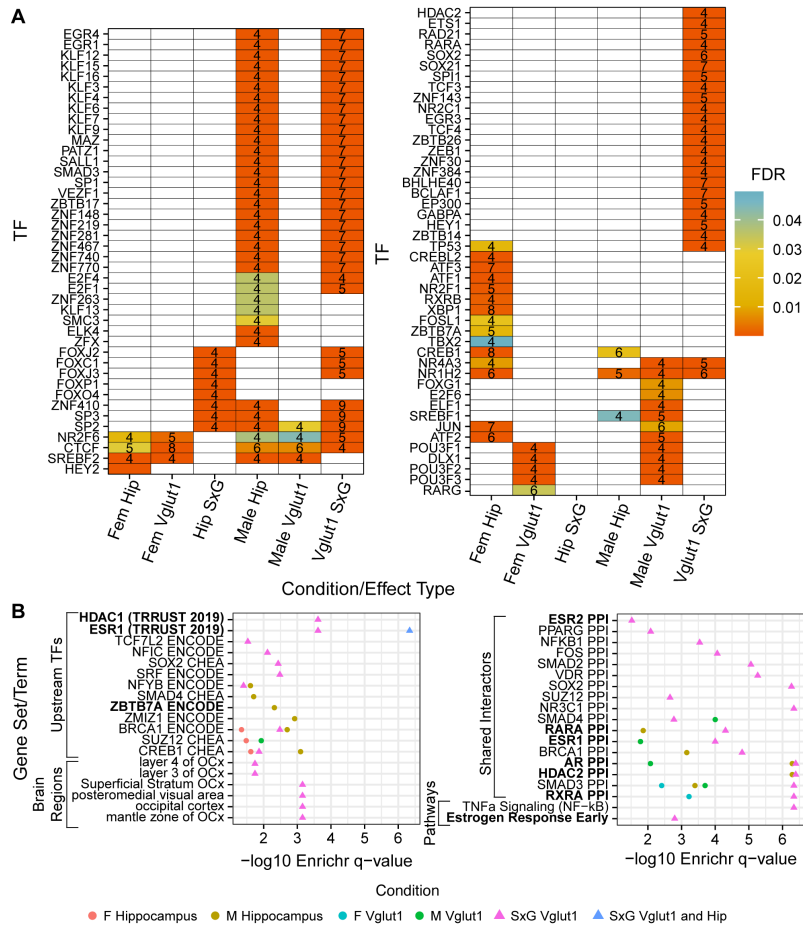
Likewise, female glutamatergic TFs were enriched for PPIs with *SMAD3* (3/8), along with *RXRA* (3/8), replicating *in vivo* my recent *in vitro* findings of retinoid-interacting rSNPs in MDD-associated loci<sup>27</sup>. Female total hippocampal TFs, on the other hand, were enriched for upstream regulation by several TFs, including *ATF2* (9/19), *BRCA1* (8/19), and *TCF3* (5/19).

### **5.2.6 Transcriptional-regulatory systems implicated in SxG interactions at MDD rSNPs in hippocampus**

Using a similar approach to that above, I investigated TFs enriched at sex-allele interaction rSNPs ( $p_{\text{emp}} < 0.05$ ) relative to sex-agnostic rSNPs (combined-sexes LMM allele main effect  $p_{\text{emp}} < 0.05$ ), and separately enriched to nonfunctional SNPs, then combined the two enrichment outputs to form TF sets (**Figure 5.4A**). I identified only 8 TFs enriched among total hippocampal SxG rSNPs,

including *ZNF410* and five TFs with a shared binding motif for *FOX(C1/J2/J3/O4/P1)*. However, I identified 57 TFs enriched at glutamatergic SxG rSNPs, including six of those from total hippocampus, suggesting coherent regulatory dynamics spanning MDD loci in this cell type. SxG TFs included several *KLF* members, as well as nuclear receptors (*NR4A3*, *NR2C1*, *NR1H2*), and another retinoid TF, *RARA*.

I ran both tissue fraction SxG TF sets through Enrichr, though nothing notable appeared for total hippocampus given the small size of the set. Top glutamatergic SxG TF enrichments included upstream regulation by *BRCAl* (21/57) and *CREB1* (15/57), as well as PPIs with AR (11/57), ESR1 (11/57), and HDAC2 (12/57) (**Figure 5.4B**). Also enriched was an MsigDB functional gene set term, “estrogen response early” (5/57). These SxG regulators implicate sex hormones and histone acetylation in both establishing sex-divergent MDD risk from upstream, e.g. via *ESR1* (5/57), and actuating it downstream via rSNP-enriched TFs and their protein interactors (AR, ESR1, HDAC2, ZMIZ1, ZBTB7A). Incidentally, both sex hormones<sup>71</sup> and histone (de)acetylases<sup>72</sup> have been major areas highlighted in recent reviews of sex differences in MDD and mouse models thereof.



**Figure 5.4. Shared regulatory architecture of rSNPs by sex, cell type, and sex-interacting SNP type.** **A)** TFs with binding motif disruptions by  $\geq 4$  nominally significant (empirical  $p < 0.05$ ) rSNPs or sex-interacting rSNPs, enriched relative to nonfunctional SNPs. Number of rSNPs associated to a given binding site are shown. **B)** Terms from Enrichr<sup>68</sup> analysis, identifying shared upstream regulators (TFs controlling expression of several of the TFs in panel A), brain regions enriched for expression of the rSNP-enriched TFs, protein interactors enriched among rSNP-enriched TFs, and MSigDB pathway term enrichment for rSNP-enriched TF sets. Bolded enrichments are discussed in the text.

### 5.2.7 Identification of rSNPs in developing whole mouse brain

As sex differences in brain structure and transcriptional regulation are established in part by the effects of sex hormones, including the perinatal testosterone surge<sup>34</sup>, I sought to investigate whether MDD risk variants were subject to sex-differential regulation during early development. To be able to assess the brain during this period, I delivered the AAV library intracerebroventricularly to embryonic day 15 (E15) mice, followed by whole brain collection at

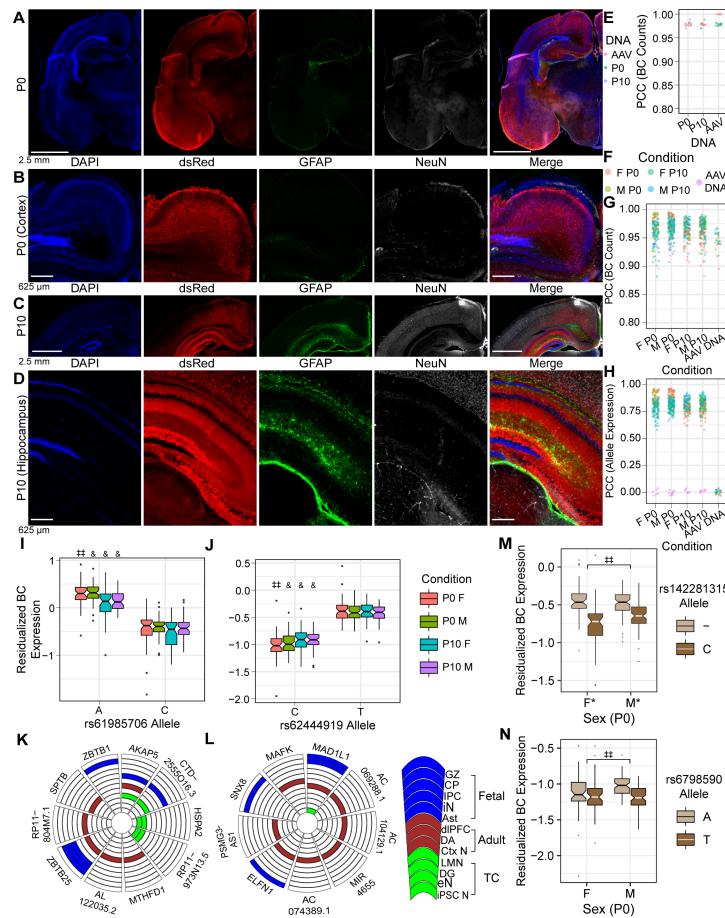


postnatal day 0 (P0) or P10. P0—while not amenable to regionally-targeted viral assays—is in the midst of the masculinizing testosterone surge, during which both acute activational effects and transcriptional-regulatory organizational effects occur; by contrast, P10 is a timepoint where sex hormones are effectively absent in normal development.

I first verified by IF that dsRed expression was detectable at P0 and P10 following *in utero* delivery. Clear, widespread reporter expression was apparent at both timepoints despite the relatively short incubation time (**Figure 5.5A-D**; **Supplementary Figure 5.8**). IF at P10 demonstrated prominent expression of the reporter in the hippocampus (**Figure 5.5D**)—a structure neither present at E15 nor well-developed at P0—consistent with prior observations of AAV9 expression ultimately occurring in hippocampus when delivered to the perinatal brain<sup>73</sup>. I subsequently collected the whole brain (except cerebellum) for RNA isolation and MPRA sequencing. Additionally, I isolated DNA from n=4 brain samples (one per age and sex) to profile the transduced library contents, verifying that the distribution of delivered MPRA barcodes was similar both between replicates ( $r$  0.86-0.94) and to input virus ( $r$  0.88-0.91) (**Figure 5.5E**). Ultimately, I analyzed 15 samples for P0 (6 female, 9 male) and 13 for P10 (6 female, 7 male). Replicates from each condition had well-correlated barcode counts (PCC 0.86-0.98) and RE expression values (PCC 0.58-0.96) (**Figure 5.5F-H**).

Within single sexes at P0, I identified 5 rSNPs in females and 12 rSNPs in males, respectively ( $p_{\text{emp}}$  FDR < 0.2), 4 of which were shared between conditions with consistent effect direction. By contrast, I identified 105 female and 72 male rSNPs at P10, with 42 rSNPs identified in both sexes with consistent direction of effect. Three of these shared rSNPs were the same rSNPs shared

between P0 sexes, with consistent effect direction in all four conditions (two of which are illustrated, **Figure 5.5I-J**). Two of these three rSNPs also have rich chromatin contact evidence supporting gene regulatory roles in fetal, adult, and cultured human neural cell types (for the illustrated rSNPs, **Figure 5.5K and 5.5L**, respectively), consistent with their detection as rSNPs in whole brain tissue, highlighting *in utero* MPRA delivery as a robust method for detecting functional variation in the developing brain.



**Figure 5.5. Validating the *in utero* MPRA delivery method, and identification of rSNPs and sex-interacting rSNPs in the developing brain. A-B)** IF of P0 brain after E15 MPRA-AAV delivery. **C-D)** IF of P10 brain after MPRA-AAV delivery. **E)** Comparability of barcode counts in recovered brain DNA and original AAV. **F)** Color legend for panels G-H. **G)** BC count correlation between samples. **(H)** Sequence expression correlation between samples. **I-J)** rSNPs rs61985706 (I) and rs62444919 (J) showed effects consistent across sexes and ages. **K-L)** Putative target genes of the respective rSNPs from Hi-C in human fetal, adult, and cultured neural tissues. **M)** Example P0 SxG SNP with comparatively small sex difference in allelic effect size. **N)** Example P0 SxG SNP with magnitude of sex difference in allelic effect comparable to smaller (female) allelic effect itself. GZ: germinal zone; CP: cortical plate<sup>74</sup>; IPC: intermediate progenitor cell; iN: inhibitory

neuron<sup>15</sup>; Ast: astrocyte<sup>58</sup>; dlPFC: dorsolateral prefrontal cortex; DA: dopamine neurons of substantia nigra and ventral tegmental area<sup>57</sup>; Ctx N: cortical neuron<sup>56</sup>; LMN: lower motor neuron; eN: excitatory neuron<sup>58</sup>; iPSC N: iPSC-derived neuron<sup>59</sup>. \*:  $p_{\text{emp}}$ -derived FDR < 0.25; \*\*: <0.2; ‡ < 0.15; ‡‡ < 0.1; & < 0.05.

### 5.2.8 Sex-allele interactions are widespread neonatally but absent during hormonal quiescence

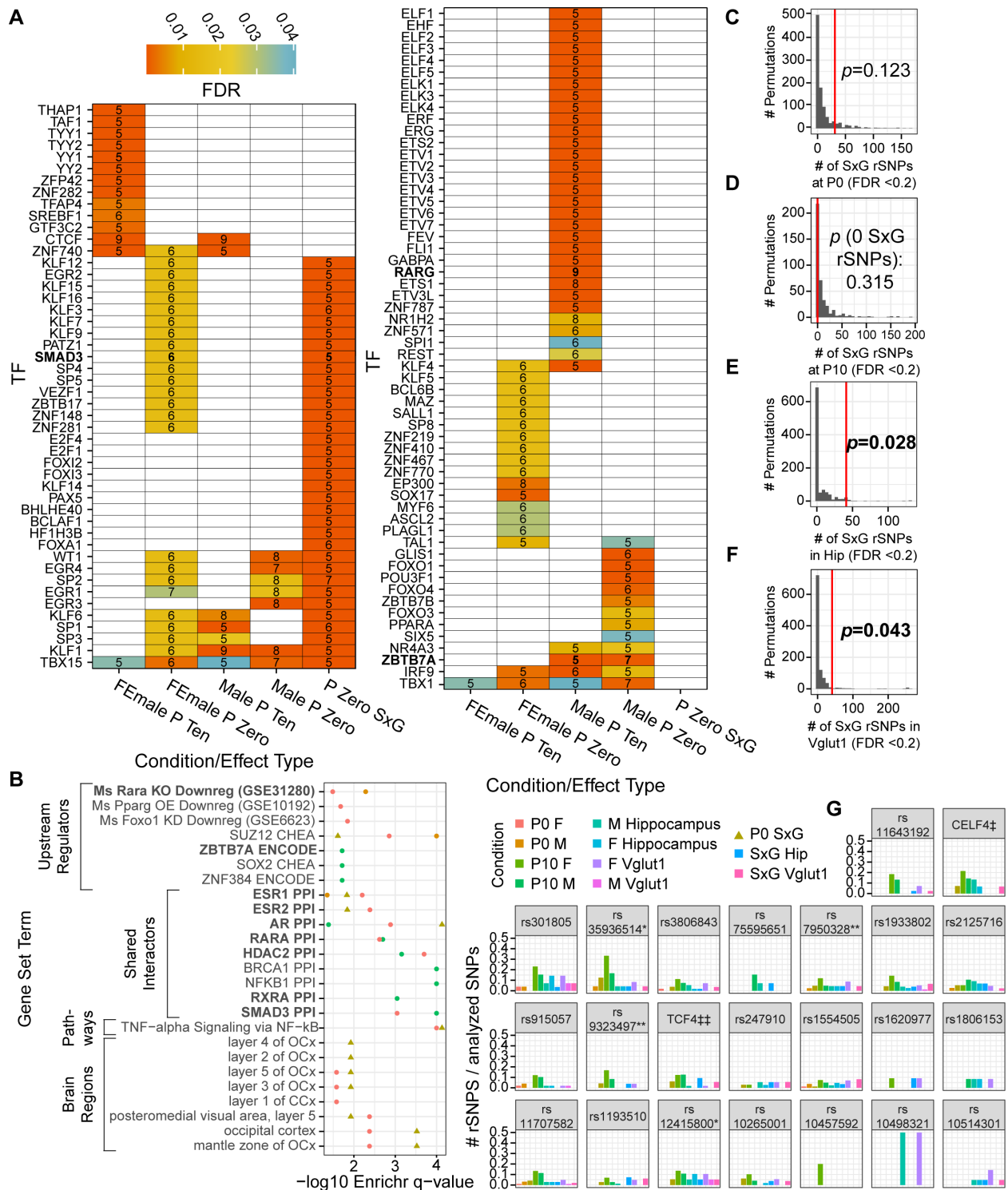
To identify SxG interactions occurring during neurodevelopment, I tested for SxG interactions as before, now within age groups. Again, the minimal promoter-only control was not sex-differentially expressed ( $t$ -test of barcode expression,  $p > 0.1$ ) (**Supplementary Figure 5.9**). At  $p_{\text{emp}}$  FDR < 0.2, I identified 31 rSNPs with sex interactions in the P0 brain (e.g., **Figure 5.5M-N**). By contrast, I identified *no* SxG interactions at  $p_{\text{emp}}$  FDR < 0.25 among the 930 analyzed SNPs at P10. I confirmed this result—despite similar  $n$  and inter-sample correlations to P0—by repeating the SxG analysis with the least variable 5 samples per sex (removing two males and one female).

### 5.2.9 Transcriptional-regulatory systems implicated in brain-wide rSNP function in postnatal development

I examined single-sex, single-age rSNP sets for enrichment of TF motif perturbations using motifbreakR and RSAT approaches as before, again using nominally significant ( $p_{\text{emp}} < 0.05$ ) rSNPs to define the sets tested for enrichment. I did not filter enriched TFs for expression in analogous human tissue, as fetal whole-brain gene expression profiles are unavailable; I instead required a more stringent minimum of 5 rSNPs to be found at significantly enriched (FDR < 0.05) TF motifs.

In P0 brain, I identified 20 TFs enriched at female rSNPs and 43 at male rSNPs, 9 of which were shared (**Figure 5.6A**). Male-specific transcription factors again included *ZBTB7A* and *NR4A3*

(**Figure 5.6A**). In P10 brain, I also found largely distinct sets of TFs in each sex (only 4 shared) (**Figure 5.6A**). Among the 42 P10 male TFs, 38 were male-specific, including *RARG*, again supporting my *in vitro* findings of retinoid-interactivity<sup>27</sup> and the cortex-specific roles of retinoid receptors across brain development<sup>75</sup>. I then looked at these TF sets as before to identify convergent regulators and functions among them (**Figure 5.6B**; full tables in the Bitbucket link under *Acknowledgments*). Female P0 TFs were enriched in Allen Atlas expression signatures for cortical layers 1, 3, and 5, and, as found in hippocampus, in PPI targets of HDAC2 and ESR1. P10 male TFs showed the greatest extent of overlap (12/42) with gene-regulatory targets of ZBTB7A. I likewise annotated P0 SxG TFs, which revealed broader extent of hormonal roles in functional variation than observed in either sex alone at P0: I found SxG TFs were again enriched for PPI targets of ESR1 (as had been P0 female TFs), but additionally enriched in PPI targets of AR and ESR2 (**Figure 5.6B**; full tables in the Bitbucket link under *Acknowledgments*).



**Figure 5.6. Regulatory architecture of rSNPs at P0 and P10, permutation analysis evaluating the number of detected SxG rSNPs, and the context-dependent landscape of MDD loci. A)** TFs with binding motif disruptions by  $\geq 5$  nominally significant ( $p_{\text{emp}} < 0.05$ ) rSNPs or SxG rSNPs, enriched relative to nonfunctional SNPs. **B)** Enrichr analysis findings for rSNP-enriched TFs in P0 and P10. Color key is shared with panel G. **C-F)** Distribution of significant ( $\text{FDR} < 0.2$ ) SxG rSNPs in 1,000 permutation analyses per condition (400 for P10); red lines: number of SxG rSNPs

identified experimentally; overlaid  $p$ -values indicate the probability of observing as many SxG rSNPs by chance. **G)** Each MDD locus, labeled by tag SNP, with bars representing the percentage of analyzed variants that were rSNPs or SxG rSNPs for each condition at  $p_{\text{emp}}$ -derived  $\text{FDR} < 0.2$ . Color key shared with panel B. \*: locus from all-female MDD GWAS<sup>43</sup>; \*\*: near-genome-wide significant locus in males with MDD developing after age 50<sup>42</sup>; ‡: collapsed results from two loci near *CELF4* with tags less than 150kb apart; ‡‡: collapsed results from LD partners of five tag SNPs, comprising two GWAS significant tags and several weaker association peaks covering a span of ~5Mb around the gene *TCF4*<sup>6</sup>.

### **5.2.10 Landscape of functional variation differs broadly across age, sex, and brain region/cell type**

If increased prevalence of MDD in females is in part caused by interactions between risk genetics and sex, then the number of SxG interactions observed in adult animals should exceed chance expectations. I therefore randomly scrambled sex labels and repeated my analyses to define a null distribution for the number of SxG rSNPs at a  $p_{\text{emp}}$  FDR of 0.2. The number of adult total hippocampal (**Figure 5.6E**) and *Vglut1*<sup>+</sup> (**Figure 5.6F**) SxG interactions I identified were both significantly greater than would be anticipated.

Overall, I found that it was the norm for loci to contain multiple functional, context-dependent SNPs, in contrast to the concept of a singular “causal variant” driving a given GWAS association. To confirm that my results reflect genetic risk effects observed in humans, I examined the relationship between GWAS and MPRA effects for over 50 SNPs directly genotyped in a large MDD GWAS<sup>7</sup>. Modest, significant Pearson correlations between MPRA and GWAS effect sizes in hippocampus and *Vglut1*<sup>+</sup> were observed, supporting the notion that functional variants beyond GWAS tags and their near neighbors influence disease risk<sup>76</sup> (**Supplementary Figure 5.10**).

Altogether, my analyses identified 280 rSNPs from 31 LD regions (28 depression-associated), with up to 13 rSNPs in a single locus found in a single condition (P10 female, tag SNP rs11707582<sup>6</sup>) (**Figure 5.6G**). In terms of age and cell type, I identified dozens of rSNPs with allele or SxG effects specific to one timepoint or tissue region: 26 P0-specific, 101 P10-specific, 34 total hippocampus-specific, and 55 specific to *Vglut1*<sup>+</sup> cells of the hippocampus. Indeed, only 64 (~23%) of rSNPs are functional in more than one developmental or cell type context. Similarly, 92 rSNPs were only identified as functional in female conditions and 37 only in male (excluding SxG interactions), while another 86 were only subject to SxG interactions. In other words, only 23% of rSNPs were sex-invariant in their *in vivo* activity.

### **5.3 Discussion and Conclusion**

I have directly measured sex-genotype interactions across MDD in the adult hippocampus and sexually differentiating brain, thus demonstrating the existence of a genetic component of sex differences in MDD and the regulatory architecture underlying these differences across space, time, and genome. I have uncovered functional differences between sexes at particular MDD-associated SNPs in the hippocampus, its excitatory neurons, and the brain during its sexual differentiation, expanding on observed genetic and clinical sex differences in MDD from general heritability to direct identification of sex-interacting variants.

One of the strengths of my study is the application of empirical *p*-values along with random-effect modeling of barcodes, controlling several aspects of experiment and design-specific technical noise in defining my results. By controlling empirically for these sources of variation in the data, I obtain results supported by a wide variety of orthogonal datasets, including reporter assays and

human brain epigenomic datasets beyond those used for variant prioritization. For example, MDD-associated SNP rs1467013 was previously demonstrated to be functional in a classical luciferase reporter assay in three different cell lines<sup>77</sup>. A prior *in vitro* MPRA identified rs301807, but not rs301806, as an rSNP<sup>27</sup>, while both were identified as functional here. This highlights the importance of *in vivo* context for obtaining relevant insights about functional variation within GWAS loci—in this case, revealing two rSNPs in close (~2kb) proximity that likely influence expression of the same target (*RERE*).

Human datasets further support the translatability of my mouse approach in identifying regulatory SNPs and sex interactions. Four rSNPs identified here, rs7244124 (*Vglut1*<sup>+</sup> SxG), rs76931017 (P10 female), rs827187 (P10 female), and rs4482931 (male hippocampus) were recently identified as chromatin accessibility QTLs in human midfetal neural progenitors and/or neurons<sup>78</sup>, consistent with the allele-differential regulatory activity I observed. Intriguingly, an early attempt to identify sex-interacting GWAS loci for MDD<sup>79</sup> found suggestive significance for a tag SNP rs1345818, near *TMEM161B* and *MEF2C*, a locus which has since been sex-agnostically associated to MDD<sup>6</sup>. My assay did not include the tag SNP, but I identified three variants in LD with rs1345818 that showed sex interactions, confirming that this risk locus indeed has sex-dependent regulatory activity: rs1814149 (hippocampal SxG, FDR<0.2, R<sup>2</sup> with rs1345818=0.67), rs5869417 (hippocampal SxG, FDR<0.1, R<sup>2</sup>=0.21), and rs6452770 (*Vglut1*<sup>+</sup> SxG, FDR<0.2, R<sup>2</sup>=0.69).

Downstream analyses aimed at identifying regulatory programs involved across rSNPs provided support for the rSNP findings themselves, while also identifying novel candidate TFs underlying sex interactions at MDD loci. My SxG-enriched TF sets were especially rich in sex hormone



receptors and interactors, consistent with expectations for an *in vivo* assay detecting sex interactions, and indicating a role of sex hormone receptors in co-regulation of MDD risk variants. My hippocampal analyses revealed male-specific roles for *AR*: male rSNPs were enriched for binding sequences of *ZMIZ1*, an *AR* co-activator, while an *AR* co-repressor, *ZBTB7A*<sup>80</sup>, was identified as a shared upstream regulator of these TFs. Notably, *ZBTB7A* also regulates human non-coding RNA *LINC00473*<sup>81</sup>, which was recently been demonstrated to have sex-differentiated effects on depressive mouse behaviors when overexpressed in cortex<sup>82</sup>. My TF analyses of P0 SxG rSNPs also identified regulatory programs consistent with the critical period for sexualization of the brain. P0 SxG rSNPs were enriched for TFs interacting with *ESR1*, *ESR2*, and *AR*, consistent with the regulatory landscape necessary for accommodation of sex hormonal signals during the perinatal testosterone surge. Additionally, *PAX5* motif disruptions were unique to P0 SxG variants; interestingly the *PAX5* motif was recently shown to be enriched in promoters of sex-differentially expressed genes in adult brain<sup>22</sup>. My neurodevelopmental rSNP-enriched TFs likewise recapitulated aspects of recent preclinical studies of sex and depression: for example, P0 female rSNPs were enriched in motifs of endothelial marker *SOX17*, consistent with demonstration of sex-differential changes in mouse brain vascular permeability after stress<sup>83</sup>.

My assay has several limitations. Notably, my P0/P10 assays could not be compared directly to my adult findings to look for sex-by-age effects, as adult experiments were limited to the hippocampus while developmental assays were brain-wide. Unfortunately, it is impractical to region-specifically deliver AAV9 *in utero*, and hippocampal morphogenesis is postnatal, precluding designs for direct comparison of hippocampus. Likewise, episomal AAV delivery may not capture all regulatory information of a genomic delivery, though benchmark studies indicate a

strong correlation<sup>84</sup>. However, I note that genome-integrated approaches for measuring variant function would require lentiviral MPRA (which failed *in vivo*, likely due to the lower titer than AAV and decreased viral spread), or base editing methods, which would be limited by efficiency and throughput *in vivo* and require mouse orthology of human SNPs. While the latter problem might be solved by using human cell lines, cell lines cannot be used to generate fully mature, differentiated neurons like those studied here (even organoids only approximate the transcriptional state of the prenatal/infant brain<sup>85</sup>), nor would they be expected to reproduce the factors that produce sex differences *in vivo* (e.g., organizational or activational hormones), precluding an *in vitro* study or replication of sex effects. Finally, results might be influenced by use of either a different minimal promoter or longer fragments if nearby elements interact with the rSNPs. However, none of the limitations above would be expected to create spurious sex effects, indicating the SxG interactions that were central to supporting my primary hypothesis are likely to be robust.

Finally, the presented *in vivo* MPRA approach indicates that critical biological and environmental factors involved in brain gene regulation and regulatory variation can be studied using a high-fidelity model of development, cell types, and biological signals. my approach provides a framework for direct, functional study of psychiatric risk genetics and their interactions with biological and environmental factors that are imperfectly modeled *in vitro*, including cell type, sex, and brain development. This same approach could be readily used in the future to directly identify variants subject to genetic-environmental interactions with other key psychiatric risk factors, such as early life adversity and chronic stress.

## 5.4 Acknowledgments and additional data

Scott Lee for assistance with immunofluorescence of pilot hippocampal injections; Michael Vasek, Ph.D. and Lexi Harris for assistance with hippocampal dissections; Ernesto Gonzalez of the Washington University in St. Louis (WU) Hope Center Animal Surgery Core for performing hippocampal AAV9 delivery; EZH Switzerland Viral Vector Facility (VVF) facility for Sanger sequencing and AAV9 packaging of the MPRA library; McDonnell Genome Institute (sequencing), funded by National Institutes of Health grant UL1TR002345; WU Center for Genome Sciences and Systems Biology (sequencing); Barak Cohen, Ph.D., Tony Fisher, Ph.D., Tomás Lagunas, Jr., and Stephen Plassmeyer for insightful discussions and methodologic guidance; WU Center for Cellular Imaging (Axioscan microscope); Biorender (hippocampus experiment schematic); WU Epigenome Browser (*RSRC1* locus visualization); and Michael White, Ph.D. and Rachel Rahn, Ph.D. for reviewing the manuscript. Funding was provided by National Institutes of Mental Health (NIMH) grant 5F30MH116654, NIMH grant 1R01MH116999, and Simons Foundation grant 571009. The annotation of library SNPs performed at the time of design are available at [https://bitbucket.org/jdlabteam/n2a\\_atra\\_mdd\\_mpra\\_paper/src/master/Library%20Design%20Epigenomic%20Annotations/](https://bitbucket.org/jdlabteam/n2a_atra_mdd_mpra_paper/src/master/Library%20Design%20Epigenomic%20Annotations/). Code is available at <https://bitbucket.org/jdlabteam/paper-resources-mdd-in-vivo-mpras/src>. Raw sequencing files and tabulated barcode counts per sample are available through GEO accession GSE186348 with reviewer token *evabumkctpsltut*. Outside datasets used for annotation are detailed in Supplementary Materials.

## 5.5 Supplementary Materials

### 5.5.1 Methods

**Animal Research.** All procedures involving animals were approved by the Institutional Animal Care and Use Committee at Washington University in St. Louis, MO.

**Design and construction of minimal promoter-reporter-WPRE 3' UTR cassette.** A previously designed MPRA reporter consisting of a minimal *hsp68* promoter and dsRed-Express2<sup>86</sup> was PCR amplified with a forward primer adding a 5' *MreI* cut site and a reverse primer with 3' overhang homologous to the 5' end of the WPRE 3'UTR element. The WPRE 3' UTR element was PCR amplified from a lentiviral plasmid encoding cyan fluorescent protein with the WPRE element (derived from plasmid FCIV)<sup>87</sup> with a forward primer containing a 5' overhang homologous to the 3' end of dsRed and a reverse primer adding a 3' *PacI* cut site. These PCR products were cut and purified from a 1% agarose gel and subjected to 10 cycles of PCR stitching (without primer) to allow overhangs to anneal and act as primers to create a contiguous sequence, followed by addition of *hsp68*-dsRed forward primer and WPRE reverse primer and 15 further PCR cycles to amplify the contiguous product. The PCR reaction was run out on a 1% agarose gel and the properly sized band (~1.5kb = ~1kb *hsp*-dsRed + ~500bp WPRE) was cut and purified. Later, Sanger sequencing (*described below*) confirmed proper assembly of the cassette.

An analogous version of this reporter was subsequently created using edited primer sequences to amplify the full cassette, replacing the 5' *MreI* site with a 5' *BsiWI* site and the 3' *PacI* site with an *AsiSI* site. This reporter was subsequently cut and cloned into single-oligo plasmids generated during the cloning optimization process. Single clones were grown in liquid culture and isolated

plasmids were verified by Sanger sequencing to provide a clonal stock of reporter for the later full-scale MPRA library of psychiatric GWAS loci.

***Design, construction, delivery, and RNA collection: cell-type-specific promoter proof-of-principle MPRA library.*** Three promoters were PCR amplified with primers adding an *MluI* cut site to the 5' end, and *MreI* and *PacI* cut sites to the 3' end (for later insertion of the hsp-dsRed-WPRE cassette). Promoters were **1)** a 2.2kb human *Gfap* promoter region<sup>88</sup> **2)** a 1.3kb mouse promoter region of the excitatory, neuron-specific gene *Camk2a*, amplified from plasmid pAAV.CamKII(1.3).eYFP.WPRE.hGH (Addgene plasmid # 105622, a gift from Karl Deisseroth to Washington University's Viral Vector Core); and **3)** a 521bp human promoter region of the constitutive, highly transcribed gene *Pgk2*, amplified from plasmid pRRLsinPGK-GFPppt<sup>87</sup>. These PCR products were used as the template for a second PCR, in which a reverse primer homologous to the added *MreI* and *PacI* cut sites was used. The reverse primer also contained an overhang consisting of a 9bp region of N's, resulting in random sequences for use as barcodes, followed by a *Sall* cut site for insertion of the product into a plasmid backbone. For the minimal-promoter only condition, a single oligonucleotide consisting of all four cut sites, intervening bases to ensure cuttability, and a barcode (*MluI-MreI-PacI*-barcode-*Sall*) was ordered, and PCR amplified. For all promoter PCR products, barcode identities were determined by Sanger sequencing after insertion of PCR product into vector (see below).

Plasmid JD386<sup>89</sup>, originally encoding *mTdTomato* under the *Gfap* promoter, was digested at *MluI* and *Sall* sites (products #R3198 and #R3138, New England Biolabs, Ipswich, MA, USA) to remove the promoter and *mTdTomato*, leaving an intact human growth hormone (hGH) poly-A

signal sequence just downstream of the *Sall* site, as well as intact ampicillin resistance and inverted terminal repeats necessary for AAV packaging. The promoter PCR product was digested with the same enzymes. The plasmid digest was treated with antarctic phosphatase (AP) (NEB #M0289) and gel purified to isolate the desired backbone fragment. The digested PCR products and gel-purified, AP-treated backbone digest fragment were ligated using T4 ligase (NEB #M0202) at 16°C for 18 hours with vector:insert molar ratios ranging from 1:3 to 1:5. Ligations were then directly transformed into DH5α chemically competent cells (NEB #C2987H), outgrown for 45 minutes in a 250rpm shaker at 37°C, then plated on LB agar with ampicillin and allowed to incubate for 16-18 hours at 37°C. Individual colonies were used to incubate 1mL wells of LB-ampicillin liquid media in 96-well deep-well plates and shaken for a further 16 hours at 37°C. The liquid cultures were then used to generate two identical 96-well plates of 25% glycerol culture stocks, one of which was sent for Sanger sequencing. Sequencing in this round of cloning used primers upstream of the *MluI* site and downstream of the *Sall* site, allowing verification of the promoter sequence and identification of barcodes. Up to 10 clones per promoter were selected based on having both expected sequencing results and unique barcodes. Equal volumes of each clone's glycerol stock were then used to inoculate a 10mL LB-ampicillin liquid culture of each promoter's "library 1" and grown at 37°C for 18 hours.

The "library 1" cultures were then mini-prepped (NucleoSpin Plasmid kit #740588.250 Macherey-Nagel, Düren, Germany) to isolate plasmids for insertion of the previously mentioned reporter cassette. Each library 1 plasmid pool, as well as the reporter cassette, was digested with *MreI* and *PacI* (Thermo Fisher #ER2021, Waltham, MA, USA; NEB #R0547). As before, the plasmid digest was treated with AP and gel purified. The digested, AP-treated, purified plasmid fragment and

digested reporter cassette were combined at a vector:insert ratio of 1:5 and ligated using T4 ligase as above. Ligations were directly used for transformation, plating, and single-clone 96-well Sanger sequencing and glycerol stocking as described above. Sanger sequencing was performed with the same primers as for Library 1; this confirmed the presence of the inserted promoter (sequencing downstream from *MluI* site), as well as retention of the barcode and the presence of WPRE and the 3' end of dsRed insert (sequencing upstream from the *Sall* site). Clones were selected for generation of “library 2” based on present and intact promoter, minimal promoter, dsRed, and WPRE, and were only selected if the barcode identified was one of the barcodes from among the library 1 clones. This resulted in 5 to 9 eligible clones for each promoter condition. The clones' glycerol stock wells were again used as inoculum for a 10mL LB-ampicillin starter culture, used to inoculate 500mL LB-ampicillin cultures for maxiprep (Plasmid Maxi kit (#12163), Qiagen, Hilden, Germany). The four maxiprepped “library 2” plasmid pools (one pool per promoter) were submitted to the Washington University Viral Vectors Core for packaging into AAV9.

***Power analysis for in vivo MPRA of ~1000 variants, examining both allele and sex effects.*** In-house *in vivo* MPRA pilot data (unpublished) was used to approximate the standard deviation of MPRA expression measures from low-depth sequencing of Vglut1+ TRAP samples. Importantly, the use of low-depth sequencing (~1 million reads in the pilot experiment for ~10,000 barcodes in each of several replicates) resulted in extremely conservative power estimates. To account for measuring both sex and allele effects, I estimated 1100 \* 2 variants being analyzed for the adult hippocampal experiment. I then used the MPRA tools package from Ghazi, et. al <sup>90</sup> to approximate the variant effect sizes I would be powered to detect using  $n=8$  replicates with 4-8 barcodes per allelic measurement (our minimum requirement was 4 per allele, described in the analysis methods

below). I estimated that with 8 replicates I would be 80% powered at a Bonferroni  $p < 0.05$  ( $p < 0.1$ ) to detect log<sub>2</sub> fold changes in allele expression of 0.35-0.52 (0.33-0.49) using the standard deviation from the pilot data. When factoring in the time required for AAV delivery into hippocampus (2 hours per animal) and the intentionally hyper-conservative nature of the standard deviation estimates I used, I injected 6 animals per sex, getting 5 samples per sex with TRAP and total hippocampal RNA viable for MPRA sequencing. For P0 and P10, I did not formally perform a further power analysis; rather, I performed a pilot study using the full library of 29,000 barcodes in 4 mice, combined with immunofluorescence, to determine whether whole-brain *in utero* would be replicable animal-to-animal. Based on high percentage of barcode recovery within and barcode count correlation between neonatal pilot replicate samples, I then aimed to collect a minimum of 6 replicable sequencing samples per sex and age—again, in part due to practical considerations, as I discovered much higher rates of RNA/sequencing sample dropout with *in utero* delivery.

***Neonatal mouse transduction with cell-type promoter MPRA library(ies); subsequent TRAP and immunofluorescence (IF).*** One litter of SNAP25-RPL10a-eGFP TRAP mice (described in<sup>38</sup>), back-crossed for >10 generations to wild-type C57BL/6J mice, was genotyped at postnatal day 2 (P2). GFP-positive pups received intracranial injection of a mixture of all four viruses mixed at an equal titer; an aliquot of the mixture was set aside for later sequencing for DNA counts during MPRA analysis. GFP-negative pups (3) were injected with either the *hsp68*-only virus, *CAMK2A* promoter virus, or *GFAP* promoter virus to confirm cell-type specificity by immunofluorescence. All animals were injected with 2 $\mu$ L bilaterally of virus or viral mix into three coordinate pairs: four anterior to bregma (two anteromedial, and two ~1.5mm posterolateral to the former), and two medially midway between bregma and lambda.



At P17, the three GFP-negative pups were perfused with 4% paraformaldehyde in PBS, brains dissected and dehydrated, and sliced in the coronal plane into 40 $\mu$ m floating sections stored in PBS with 0.01% sodium azide for immunofluorescence (further IF method details below). Tissue was stained with primary antibodies as follows: mouse anti-NeuN 1:500 (MAB377), goat anti-*GFAP* 1:250 (Abcam ab53554), and rabbit anti-dsRed 1:500 (Rockland 200-301-379). Secondary antibodies were used at 1:1000 donkey anti-mouse with Alexa Fluor 647, donkey anti-goat with Alexa Fluor 488, and donkey anti-rabbit with Alexa Fluor 546. Sections were then slide-mounted and imaged on a Zeiss LSM 700 (Zeiss, Germany) using a 40x oil objective.

At P27, the three surviving GFP-positive mice were sacrificed for individual biological replicates of TRAP (methods below), performed as previously described<sup>63,91</sup>. Briefly, whole brain was separately homogenized from each mouse. 2.5% of the supernatant from the 20,000  $x$  g spin of homogenate was collected as an “input” sample of RNA, i.e. RNA from all cell types in the tissue. Both input and TRAP RNA were QC’ed using an Agilent TapeStation’s high-sensitivity RNA assay (Agilent, Santa Clara, CA). Two replicates were retained, with input RNA integrity number estimates (RINe) 7.8 and 8.2; TRAP RINe values were 6.8 and 7.4.

***Design of the neuropsychiatric GWAS MPRA library.*** The neuropsychiatric GWAS MPRA was designed as previously described<sup>27</sup>. Briefly, tag variants were selected from neuropsychiatric GWAS studies, predominantly those of MDD or meta-analyzing MDD alongside additional neuropsychiatric disorders<sup>6,7,33,40–45</sup>. Additional tag variants were selected from GWAS of anxiety disorders<sup>48</sup>, attention deficit hyperactivity disorder<sup>92</sup>, educational attainment/intelligence<sup>93,94</sup>, and

traits showing strong SNP coheritability with MDD (i.e., neuroticism<sup>46,47</sup> and mood instability<sup>49</sup>). As a negative control locus, I selected one tag related to anthropomorphic traits, rs1883640<sup>95</sup>. Loci were expanded to SNPs with LD  $R^2 > 0.1$  and minor allele frequency  $> 0.01$  with the tag variant in their appropriate 1000 genomes phase 3 population (EUR for all studies except two loci from the Han Chinese CONVERGE GWAS of MDD) using LDLink<sup>96</sup>. SNPs were manually selected for inclusion in the MPRA library on the basis of overlap with both human brain eQTL<sup>53,55,97–100</sup> and at least one additional human brain epigenomic annotation<sup>58,101–105</sup>. 11 SNPs were removed due to introduction of a nonsynonymous coding change. Two loci instead included all SNPs in LD  $R^2 > 0.65$  and SNPs with a RegulomeDb<sup>106</sup> score of  $\geq 4$  for LD  $0.1 < R^2 < 0.65$ . For loci especially sparse in overlaps with the screening annotations, a single overlap of any sort was considered adequate. I note that while this approach is not wholly empirical, recent machine learning-based approaches<sup>107</sup> have identified that regulatory variation is better predicted by considering annotations only in the variant's local genomic region (as I did manually here).

Oligonucleotides (oligos) were designed for the 1453 selected SNPs, containing up to 126bp of human genomic sequence (hg19) centered on the variant. For small insertion/deletion variants, the larger allele was used to define the 126bp sequence (i.e., the smaller allele was 126-(allele length difference) bp). Each allele of each variant was assigned 10 randomly generated 10bp barcodes from a filtered set  $\geq 2$  Hamming distances apart, comprised of 25-75% GC, and filtered for runs of  $> 3$  of any one base. Promoter-only (“basal”) oligonucleotides instead included a 126bp filler in a region of the oligonucleotide cut out during cloning (below), paired to 110 of these barcodes. The deeper barcoding of the basal oligonucleotide allows for a well-powered, within-sample normalization factor to be calculated (see analysis methods below). The remaining oligonucleotide

sequence, 200bp in total, consisted of restriction sites for cloning and primer sites for oligo amplification from the single-stranded oligo pool (**Supplementary Figure 5.1**).

A 200bp oligonucleotide pool, comprised of 29,280 unique sequences, was ordered from Twist Biosciences (San Francisco, CA). The oligonucleotides were flanked on the outside by a forward priming sequence, 5' GAGGGAAATCGTGACGCGTG 3' (forward primer same sequence), and reverse priming sequence of 5' GTCGACCAGGTCATCACTATTG 3' (reverse primer 5' CAATAGTGATGACCTGGTCGAC 3'). The internal portions of these sites are MluI (followed by a G, to prevent creation of other used restriction sites resulting from design of adjacent genomic sequence) and SallI cut sites, respectively. Proceeding from the 5' end, the remaining space comprised the  $\leq 126$ bp allelic sequence of interest, followed by BsiWI, PmeI, and AsiSI sites end-to-end, and the barcode (**Supplementary Figure 5.1**).

***Cloning of the neuropsychiatric GWAS MPRA library.*** Oligonucleotides were amplified by PCR for 12 cycles in 6, 50 $\mu$ L reactions, using Phusion High-Fidelity Polymerase 2x Mastermix (New England Biosciences) each with 500nM final primer concentrations and 10ng oligonucleotide pool as template. Reactions were thermocycled as follows: 98°C, 3 minutes  $\rightarrow$  12 cycles of 98°C/15s, 65°C/30s, 72°C/15s  $\rightarrow$  72°C/5min  $\rightarrow$  4°C. Products were loaded into an unstained, 2% tris-acetate-EDTA (TAE) agarose gel run at 90V for 105 minutes. The gel was post-stained with 15 $\mu$ L 10,000x SyBr gold (Invitrogen) with 150mL of TAE in a small pyrex dish, gently rotated for 20 minutes, and the 200bp product band cut for gel purification using the Nucleospin Gel and PCR Clean-Up kit (Takara Biosciences, Kusatsu, Shiga, Japan). 12ng of product was purified for cloning.

The purified PCR product was cut in a 100 $\mu$ L double digest reaction containing 1.5 $\mu$ L MluI-HF (NEB), 1.5 $\mu$ L SallI-HF (NEB) and supplied buffer at 37 $^{\circ}$ C / 1hr followed by an 80 $^{\circ}$ C/20 min heat kill. The AAV-compatible plasmid backbone (JD386, see above) was digested separately in four 100 $\mu$ L reactions, each containing 2 $\mu$ L MluI-HF, 2 $\mu$ L SallI-HF, and 2 $\mu$ L BamHI-HF (to resolve the target product on subsequent gel), supplied buffer, and 2 $\mu$ g plasmid, incubated at 37 $^{\circ}$ C / 1hr and heat killed at 80 $^{\circ}$ C / 10min. Subsequently, 4 $\mu$ L H<sub>2</sub>O, 4 $\mu$ L Antarctic Phosphatase, and 12 $\mu$ L of supplied buffer (NEB) were added to each reaction; the reactions were then divided into 60 $\mu$ L aliquots for further incubation at 37 $^{\circ}$ C / 1hr. The vector was then run on a 0.8% agarose gel with GelGreen dye (Biotium, Fremont CA), the target 3.4kb band cut and purified using the Nucleospin kit per manufacturer directions. The oligo digest reaction was used directly for ligation with the purified vector at a molar ratio of vector:insert 6:1. Nineteen 20 $\mu$ L ligation reactions were aliquoted from a single master mix containing 11.9ng digested oligo, 1.51 $\mu$ g digested, phosphatased, and purified vector, 10 $\mu$ L Enzymatics T4 ligase and 40 $\mu$ L supplied buffer (Qiagen). Ligations incubated at 16 $^{\circ}$ C / 14hr, then heat killed at 80 $^{\circ}$ C / 15min. Purification of the ligation product was performed using MyOne Silane magnetic beads (Thermo Fisher) from which the supplied buffer was removed and replaced with 3x the total ligation volume of Qiagen Buffer QG. The ligation reactions were recombined into a single 1.5mL microcentrifuge tube with the buffer QG•silane bead mixture, and incubated end-over-end at room temperature for 75 minutes. Beads were then washed on a magnet stand by removing supernatant and adding 80% ethanol thrice (without disrupting magnetized beads). The beads then air dried for 10 minutes and product was eluted into 196 $\mu$ L of water, with 192 $\mu$ L taken for transformation into *E. coli*.

For transformation, 6 $\mu$ L (~18 ng) of ligation product was added to each of 27 tubes of 50 $\mu$ L of DH5alpha chemical competent *E. coli* (NEB), which were transformed per manufacturer instructions. Three transformations were simultaneously performed using DH5alpha electrocompetent *E. coli* (NEB), each consisting of 8 $\mu$ L ligation product (~24ng) and 100 $\mu$ L cells in a 1mm cuvette per manufacturer instructions. All transformants were grown out for 45 minutes at 250 rpm and 37°C in their PCR tubes (chemical transformation) or after transfer into culture tubes with 2mL pre-warmed SOC (NEB; electrochemical transformation). The outgrowths were pooled (19mL total), added to a flask of 31mL luria broth (LB) with carbenicillin for a maxiprep starter which incubated with 250rpm, 37°C shaking for 6 hours. The full starter volume was then added to 450mL LB with carbenicillin and cultured with 250rpm shaking at 37°C for 10 hours. The plasmid was maxiprepped using GenElute HP Plasmid Maxiprep Kit (Sigma Aldrich, St. Louis MO). Low-depth next-generation sequencing was performed on the maxiprep with amplicons generated over the space spanning the MluI and Sall sites to confirm the presence of the library.

The Sanger-verified *hsp68*-dsRed-WPRE cassette (described above) was cut back out of its backbone in four 30 $\mu$ L digest reactions, each containing 2 $\mu$ g of the clonal miniprep plasmid, 1 $\mu$ L BsiWI-HF (NEB), 1 $\mu$ L AsiSI (NEB), supplied buffers, heat killed at 80°C/10min, then gel purified from a GelGreen-stained 0.8% agarose gel using the Nucleospin kit. The maxiprepped “first plasmid library” was digested in 20 $\mu$ L reactions of 215ng each, using 1 $\mu$ L AsiSI for 37°C/1hr, followed by direct addition of 1 $\mu$ L BsiWI-HF, 1 $\mu$ L 10x CutSmart buffer, and 8 $\mu$ L water, and digestion at 37°C for a further hour. 1 $\mu$ L rSAP phosphatase (NEB) and 9 $\mu$ L H<sub>2</sub>O was then added to the reaction and incubated at 37°C / 1hr, 80°C / 20min, and brought to 4°C. Forty 20 $\mu$ L ligation

reactions were prepared as described above, each containing ~50ng DNA in 1:6 molar ratio of the cut vector and reporter cassette. The ligations were thermocycled 25°C/30s, 16°C/30s 20 hours (~800 cycles in practice) followed by 80°C/20min heat kill and 4°C incubation. Ligations were cleaned up with Silane beads as described above.

The “second plasmid library” was transformed into 25 tubes of DH5a chemically competent cells; 20 reactions were 200µL cells, 19µL (~100ng) ligation; 2 were ~150µL cells and 5µL (~25ng) ligation; 2 were 50µL cells and 5µL ligation, and one was 50µL cells and 2.5µL (~12.5ng) ligation, all of which were transformed per manufacturer instructions. Cells were outgrown for 50 minutes at 37°C/250rpm. Outgrowth was pooled and used to inoculate a starter and subsequent maxiprep culture as described above. Sequencing (in the same manner as performed for later experimental analyses) was then performed to verify barcode coverage of the plasmid library. The resulting plasmid was sent to EZH Zurich Viral Vector Facility, which verified its sequence by Sanger sequencing and packaged the library in AAV9.

***Hippocampal stereotaxic delivery of the neuropsychiatric GWAS MPRA library.*** Hippocampal stereotaxic injections were performed in a counter-balanced fashion (alternating order male/female, with the order switched each day) over a three-day period (n=4 per day). Mice were anesthetized with continuous isoflurane up to 5% until unreactive to toe pinch, then maintained at 2-3%. The skin of the head was retracted to expose the surface of the skull. A 0.485 µM hamilton syringe was stereotactically guided into place and used with an automatic pump to deliver 2µL of AAV9 at a rate of 0.2µL per minute to each of four positions (8µL AAV9 total per animal; deeper coordinate first): A/P -2.0 mm, M/L -1.5mm and 1.5mm, and depth (from dura) of -1.75mm first,

followed by partial withdrawal of the syringe and an additional delivery at -1.25mm depth. The needle was allowed to dwell for 10 minutes after each injection was completed. Mice were then given intraperitoneal buprenorphine 0.1mg/kg and cranial subcutaneous lidocaine 25mg/kg for post-operative analgesia and recovered in the home cage under a heat lamp with cagemates. Mice were checked for ambulation and absence of distress at 1, 24, 48, and 72-hours after return to consciousness. Mice then continued to live in the home cage with same-sex cagemates (same cagemates as prior to surgery) with the same *ad libitum* access to food and water and the 12:12 light/dark cycle they had been reared under.

***In-Utero AAV injections and brain collection/lysis.*** CD1 IGS timed-pregnant female mice were purchased and delivered (Charles River Laboratories). Pregnant dams were allowed to house overnight after delivery to reduce stress and ensure optimal success during in-utero injections. Pregnant mice at 15 days of gestation were placed under isoflurane anesthesia and a midline laparotomy was done to expose the uterus. 1uL of viral MPRA-AAV was mixed with 0.025% Fast Green FCF (Sigma, St. Louis MO) and administered through a 10uL-Drummond glass micro dispenser pipette (Drummond Scientific, Broomall, PA) into pups through the uterine wall. Injections were targeted towards the anterior horn of the lateral ventricles as previously described in IUE and AAV injection studies<sup>108,109</sup>. For each dam, 2-4 pups were left non-injected due to position of pup in-utero or impaired development compared to other pups. On P0 and P10, mice were sacrificed and screened for fluorescence under a dual fluorescent protein flashlight (NightSea, Lexington MA). P10 mice brains were weighed post-removal to assess long-term changes in brain weight due to viral injection. Pups were decapitated and brains dissected out with removal of the cerebellum, and the remaining brain placed in a microcentrifuge tube on dry ice,

then stored at -80°C overnight to aid homogenization. The following day, the brains were removed from -80°C storage for addition of 500µL (P0) or 1mL (P10) Trizol reagent, homogenized thoroughly using a battery-powered hand pestle, and returned to -80°C until all samples were collected for that age group for single-batch RNA purification (*below*).

***Translating Ribosome Affinity Purification (TRAP).*** For immunoprecipitation (IP) of GFP-containing ribosomes, 60µL per sample of streptavidin MyOne T1 beads were resuspended in 17µL of 1µg/µL/sample protein L (reconstituted in 1x PBS), 36µg/µL/sample anti-eGFP 19C8 and 36µg/µL/sample anti-eGFP 19F7 (both available through Sloan-Kettering's antibody & bioresource center <https://www.mskcc.org/research/ski/core-facilities/monoclonal-antibody-core-facility>), brought up to 200µL times the number of immunoprecipitations to be run with 1x PBS. The beads in this mixture were incubated with end-over-end mixing at 4°C for two hours to bind antibodies to the magnetic beads. Beads were then separated on a magnet stand and washed by resuspension and remagnetization 5 times using 0.1% bovine serum albumin (BSA) in 1x PBS. Beads were then washed three more times using wash buffer (above), then resuspended to a final volume of 105µL/IP and kept on ice until needed.

On the day of tissue collection, DTT, RNase inhibitors, and cycloheximide were added to stock TRAP buffers to the specified concentrations detailed below. Mice were deeply anesthetized with isoflurane, rapidly decapitated, and the brain removed and dissected for TRAP. Each brain was bluntly dissected to remove anterior and posterior-most portions, then placed in a pre-chilled dish on ice containing 15-25mL of modified TRAP-compatible buffer consisting of 1x phosphate-buffered saline (PBS), 0.1 mg/mL cycloheximide (to halt translation for ribosome capture), and



1/2500 vol/vol each of rRNAsin (Promega, Madison, WI USA) and SUPERase<sup>•</sup>in (Thermo Fisher, Pittsburgh, PA USA). Hippocampus was dissected out bilaterally in a dish of buffer on ice, then both hippocampi from a single animal homogenized to constitute a sample. Samples were sequentially dissected and homogenized.

TRAP was then performed as previously described with slight modifications. A wall-powered drill, run at full speed, was fitted with a Teflon pestle was used to homogenize each sample in 1mL pre-chilled homogenization buffer (10mM HEPES pH 7.4, 150 mM KCl, 10mM MgCl<sub>2</sub>, 0.5 mM dithiothreitol (DTT), 0.1 mg/mL cycloheximide, and 1/1000 vol/vol each of rRNAsin and Superasin, and Roche EDTA-free protease inhibitor cocktail (dissolved in homogenization buffer without DTT, cycloheximide, or RNase inhibitors) to a final concentration of 1x). Homogenates were spun down at 2,000 x g for 10 minutes at 4°C. 825µL of supernatant was collected and combined with 100µL each of 10% NP40 in water and 300mM 1,2-diheptanoyl-sn-glycero-3-phosphocholine (DHPC; Avanti Polar Lipids, Alabaster, AL USA) and incubated on ice for 30 minutes. This solution was then spun at 20,000 x g for 15 minutes at 4°C. For the input (or here, “total hippocampus”) RNA fraction, 50µL of supernatant was collected and added to 200µL of wash buffer (1% vol/vol NP40, 10mM HEPES pH 7.4, 150 mM KCl, 10mM MgCl<sub>2</sub>, 0.5 mM dithiothreitol (DTT), 0.1 mg/mL cycloheximide, and 1/1000 vol/vol each of rRNAsin and Superasin, and Roche EDTA-free protease inhibitor cocktail (dissolved in homogenization buffer without DTT, cycloheximide, or RNase inhibitors)) and 750µL Trizol LS (Thermo Fisher), and stored at -80°C until all TRAP and input samples had been collected.

975 $\mu$ L of the same supernatant was taken for ribosome capture by IP and added to an aliquot of 100 $\mu$ L of the resuspended, antibody-coupled beads. The mixture then incubated end-over-end for 5.3-5.7 hours at 4°C. After incubation, beads were separated on a magnet stand, supernatant removed, and resuspended in 1mL high-salt wash buffer (10mM HEPES pH 7.4, 350 mM KCl, 10mM MgCl<sub>2</sub>, 0.5 mM dithiothreitol (DTT), 0.1 mg/mL cycloheximide, and 1/1000 vol/vol each of rRNasin and Supersasin, and Roche EDTA-free protease inhibitor cocktail (dissolved in homogenization buffer without DTT, cycloheximide, or RNase inhibitors)) before remagnetizing. This supernatant was removed, the beads suspended in a second wash of high-salt buffer, and transferred to a new microcentrifuge tube (to avoid nonspecific RNA stuck on tube walls from releasing by later addition of Trizol LS). Beads were then magnetized and washed twice more with the high-salt wash buffer. Beads were then resuspended in 250 $\mu$ L of the *wash* buffer (i.e., the 150mM KCl buffer described), and 750 $\mu$ L Trizol LS was added. These samples too were stored at -80°C until all samples were collected.

***RNA purification (all experiments) and brain DNA isolation (P0/P10).*** RNA purification was using the Zymo Clean and Concentrator-5 (Zymo, CA USA), with simultaneous processing of all samples for each condition (i.e., all hippocampal input and TRAP samples in one batch; all P0 samples in one batch; all P10 samples in one batch, see **Supplementary Table 5.1**). The samples were removed from -80°C, allowed to come to room temperature for 5-10 minutes, followed by addition of 20% volume of chloroform. Tubes were shaken vigorously by hand for 30 seconds and allowed to stand for 7 minutes at room temperature. Subsequently, tubes were spun at 12,000 x g for 20 minutes at 4°C. For hippocampal samples, 500 $\mu$ L of supernatant was collected; otherwise,

175 $\mu$ L was collected. The remaining phase-separated Trizol-chloroform sample was returned to -80°C (see DNA collection below).

For uniform RNA handling, the Zymo kit instructions were followed, but prepared one mastermix adequate for all samples being purified so as to avoid variability in volumes of buffer/ethanol added to each. This mastermix consisted of 2 supernatant volumes of Zymo RNA Binding Buffer and 3 supernatant volumes of 100% ethanol per sample. 5 supernatant volumes of this mix was added to each Trizol-chloroform supernatant, mixed thoroughly by pipetting, and applied to the kit columns, spun through at 12,000 x g at room temperature until columns were loaded. Manufacturer instructions were then followed for the remainder of cleanup. RNA quality was assessed using Agilent High-Sensitivity RNA TapeStation assay. All samples in all experiments used for sequencing preparation, including the proof-of-principle pilot, had an RNA integrity number (RIN) of, at minimum, 6.

For isolation of DNA from P0 and P10 brain to verify the presence of MPRA barcodes at the DNA level, phase-separated Trizol-chloroform mixtures were brought to room temperature and used for DNA isolation according to manufacturer instructions.

***Immunofluorescence.*** Brains not isolated for RNA analysis were removed and fixed in paraformaldehyde in PBS (4%) followed by serial sucrose in PBS solutions (15%, 30%). Following post-fixation, brains were embedded in OCT (Sakura, Torrence CA) and sectioned at 30 $\mu$ m (hippocampus) or 35  $\mu$ m (P0, P10) using a Leica CM1950 cryostat (Buffalo Grove, IL) and

processed for immunofluorescence as slide-mounted sections for P0 and free-floating sections for P10 pups and adult hippocampus.

To characterize adult hippocampal AAV9 delivery, 30 $\mu$ M coronal sections were cut 21 days after delivery, incubated in 1x PBS with 5% normal donkey serum and 1:1000 chicken anti-GFP (to identify TRAP-positive cells), 1:500 rabbit anti-RFP ((1:500, Rockland, 600-401-379; to identify AAV-transduced cells), and 1:500 goat anti-GFAP (1:500, Abcam ab53554) to visualize potential astrocytosis around the viral injection sites. Slices incubated in primary antibodies overnight at room temperature on a horizontal mixer, were rinsed three times in 3x PBS, followed by 45 minutes of incubation in 1x PBS with 5% normal donkey serum and 1:1000 each of Alexa Fluor 488 donkey-anti-chicken, Alexa Fluor 568 donkey anti-rabbit, and Alexa Fluor 647 anti-goat. Sections were rinsed again with 1x PBS, then incubated for 5 minutes in PBS with 1:20,000 DAPI for nuclear fluorescence, rinsed once more with 1x PBS, then mounted onto slides with application of Prolong Gold, followed by nail polishing of cover slips into place. Slides were stored in a foil-covered box at 4°C until imaged on the Axioscan.Z1 slide scanner (ZEISS, Germany) at 10x resolution.

For P0 Primary antibodies included anti-RFP (as above), anti-GFAP (as above), and anti-NeuN (1:250, MAB377). Fluorescently conjugated secondary antibodies (AlexaFluor 488, 568, and 647) were obtained, and nuclei were labeled with a DAPI counterstain. “No primary” controls were done for both sets of time points to indicate and ensure staining was specific to primary antibody targets. Multi-channel imaging was performed at 20X using an AxioScan.Z1 slide scanner to assess both independent region and whole-brain viral transduction. Image editing was performed

using ImageJ software and only included re-scaling of resolution, brightness/contrast adjustments, and cropping.

***RNA-seq library prep (proof-of-principle and Neuropsychiatric GWAS MPRA libraries).*** RNA samples were treated with the Turbo DNA-Free kit (Ambion #AM1907, Austin, TX, USA) to remove extant DNA using the manufacturer's instructions for high-concentration DNA (2 $\mu$ L of enzyme, followed by 20% volume of Inactivation Reagent to remove the DNase). Sequencing libraries were prepared from RNA by performing a variation on the methods of *e.g.*<sup>50,86,110</sup> by using a reporter-specific primer, targeting the polyA signal sequence just 3' to the barcode sequence<sup>111</sup>, during first-strand reverse transcription with Superscript III Reverse Transcriptase (Invitrogen 18080044, Carlsbad, CA, USA). Double-stranded cDNA was then synthesized and amplified by PCR using Phusion HF (NEB #M0531) using the same reverse primer and a forward primer in the WPRE of the 3'UTR. These primers added unique cut sites allowing subsequent Illumina adapter ligation. To prevent sequencer errors due to homogenous sequence at the start of read 1 (3' end), these adapters were a mix of four different lengths to stagger the first base of the 3' end read. Digestion, clean-up, ligation, clean-up, and final PCR with primers with partial homology to the adapter ends were used to add the remaining Illumina sequences and sample indices. The full details of each sample processing step, including bead-based size selection, were as described<sup>27</sup>. Sample input mass, number of PCR cycles for the two PCR steps (single-strand cDNA to double stranded DNA and index PCR), and read depths are described in **Supplementary Table 5.1**. Sequencing of proof-of-principle samples was performed on a MiSeq instrument (Illumina, San Diego, CA); all other experiments were sequenced on an Illumina NovaSeq 6000 instrument. NovaSeq 6000 (Illumina). Three samples from a pilot of the hippocampal *Vglut1*<sup>+</sup> TRAP-MPRA

in female adults were additionally used as a validation dataset. In the pilot experiment, two technical replicates for sequencing were prepared for each input (total hippocampus) sample, using 66ng (to match the input RNA mass of the pilots' *Vglut1*<sup>+</sup> samples) or 150ng RNA.

***qPCR verification of cell-type marker changes in TRAP.*** 20 $\mu$ L reverse transcriptase reactions were prepared using Quanta Biosciences qScript with supplied 5x buffer (containing both random hexamer and poly-T primers), containing 20ng of sample RNA. Reverse transcriptase reactions were incubated at 25°C/5minutes, 42°C/30minutes, 85°C/5minutes, then kept on ice. RT reactions were diluted with 140 $\mu$ L water (final volume 160 $\mu$ L) for use as qPCR template. 10 $\mu$ L qPCR reactions (technical triplicates per sample•gene) were prepared using Sybr Green 2x Mastermix (Thermo Fisher), containing 4 $\mu$ L cDNA, 0.5 $\mu$ L each primer, and 5 $\mu$ L of qPCR mastermix. qPCR thermocycling ran for 40 cycles at 95°C/15s, 63°C/30s per cycle, followed by a melt curve on a Quantstudio 6 instrument. Before analysis, data were quality checked by a) examining melt curve product heights (all were a singular, consistent peak per gene over 80°C, corresponding to a true amplicon as opposed to primer dimers) and b) identification of outlier wells based on a cycles to threshold of detection (CT) value  $\geq 1$  cycle different from other sample•gene technical replicates. One row of technical replicates corresponding to a total hippocampal sample were removed due to large CT discrepancies relative to the other two technical replicate sets; 6 other singular wells were excluded on the same basis of outlier status, and 1 well was excluded for failure to amplify any product. All analyzed sample•gene wells contained at least two technical replicates in strong agreement (CT values within 0.5 of one another). To assess enrichment in TRAP relative to total hippocampal RNA, the technical replicate mean CT value for the internal control gene,  $\beta$ -actin

(*Actb*) was subtracted from the CT value of each other gene. qPCR primer sequences are in **Supplementary Table 5.2**.

***MPRA sequencing analysis: proof-of-principle experiment.*** Barcodes were counted from read 1 sequences, allowing up to 3 mismatches in the 20bp upstream of the barcode and 0 mismatches within the barcode sequence itself. The number of reads mapping to each barcode were totaled and normalized to counts per million (CPM) with normalization for sequencing library size (in number of reads mapped) using EdgeR<sup>112,113</sup>. Expression for a given barcode was then calculated as the ratio of (CPM RNA / CPM viral DNA) for each RNA sample (TRAP and input RNA from each brain). Expression values were normalized to the within-sample mean expression of the minimal promoter (*hsp68*) alone by taking expression (BC in sample) / expression (mean(*hsp68* BCs)). By normalizing within sample (i.e., input or TRAP), expression is thus normalized to general minimal promoter activity for that sample. Significance was calculated by performing repeated-measures ANOVA / linear mixed modeling. Each barcode group's expression was implemented as a repeated measure and modeled as a dependent variable of RNA sample type and of a random variable for source tissue, thus:  $expression \sim RNA.fraction + (I|mouse)^{114}$ . For plotting and interpretation of barcode enrichment/depletion between biologically paired input and TRAP RNA samples, log2 fold-change in expression was calculated by subtracting log2(normalized expression in Input) from log2(normalized expression in TRAP) for each barcode.

***MPRA sequencing analysis: all other experiments.*** Barcodes were counted from read 1 sequences, allowing mismatches neither within the barcode sequence nor the 6bp upstream/8bp downstream of flanking sequence. CPM were calculated using the total number of barcode-

mapping reads prior to several filtering steps; for samples with multiple sequencing runs of data (DNA technical replicates and a subset of P0 RNA samples), a single CPM value per barcode per sample was first calculated by obtaining the mean CPM across the sequencing runs. 1) Barcodes with a DNA count under a specified threshold (187 or approximate CPM equivalent) are excluded from the counts table from the DNA *and* RNA samples. 2) For sequences with 4 or fewer remaining DNA barcodes represented across all samples after this step, all other barcodes for the sequence are excluded from all of the samples in the table to avoid analysis of sequences with inadequate barcoding depth. 3) Counts of RNA barcodes are then removed on a per-barcode-per-sample basis if they fall below a separate minimum read threshold (75 counts or approximate CPM equivalent), set below the DNA threshold to allow for detection of repressive effects. 4) Preliminary expression values for barcodes ( $\log_2(\text{RNA barcode CPM} / \text{DNA barcode CPM})$ ) are calculated for each replicate and collapsed across barcodes into a mean for each Regulatory Element (RE) within sample. Single barcodes with outlier expression values ( $\geq 2$  standard deviations) apart from other barcodes for that sample are dropped only from that sample. (The expression values are not written out to a results table at this time). 5) Penultimately, all barcodes are dropped from individual samples if 4 or fewer barcodes remain for a given sequence in that sample, such that all samples analyzed for a given sequence have at least 4 barcodes represented in each sample. 6) A final check is made to ensure that each barcode remaining is represented in at least 50% of samples, and those represented in fewer samples are removed, followed by a second check that all samples have  $\geq 4$  barcode expression values remaining.

Expression for a given barcode was then calculated as the  $\log_2$  ratio of (CPM RNA / CPM viral DNA) for each RNA sample. Prior to linear modeling, within-sample barcode expression values



were normalized by subtracting the within-sample mean expression of the set of barcodes paired to the minimal promoter (*hsp68*) alone (<=110 barcodes). By normalizing within sample, expression is thus normalized to minimal promoter activity among the cell types and proportions comprising the RNA sample—while this notably does not change inter-sample correlations, it *does* alter the Euclidean distance between samples in hierarchical clustering (namely, samples with outlying barcode wise expression before normalization cluster back in with the other samples after this transformation). Shapiro tests for normality were performed on each condition-wide set of barcode expression values for each SNP (i.e., 4-10 barcodes \* N samples \* 2 alleles); all Shapiro tests were >0.05 for all analyzed SNPs in all eight conditions.

Linear mixed modeling was then applied within single sexes for each condition to analyze allelic effects alone, and with data from both sexes pooled to test for allele-by-sex interactions. Each barcode group's expression was implemented as a repeated measure and modeled as expression as a dependent variable of allele (and for interaction models, of sex and sex-by-allele), thus:  $expression \sim allele + sex + allele*sex + (1|barcode)$ . I note that the random intercept values determined for barcodes were consistent even when running the model using different samples, indicating I were detecting and removing biologically invariant effects of the barcode sequences on RNA levels (**Supplementary Figures 5.2-5.4**). Empirical null test statistics were calculated as previously described<sup>27</sup>; in brief, I applied the respective experimental model to 50,000 comparisons between two “alleles” each comprised of 6 randomly selected barcodes from among the 110 corresponding to the *hsp68* minimal promoter alone, thus controlling for the noise inherent to the assay and sample set. These values were then used in the qvalue package<sup>52</sup> to determine

empirical p values, and corresponding q-values and FDR significance were generated from the empirical test statistics and p-values (also using the qvalue package).

The pilot/validation hippocampal MPRA-TRAP dataset from n=3 adult females was sequenced with technical replicates using two starting quantities of total hippocampal RNA: 66ng or 150ng. For validation analysis of total hippocampal RNA from these three RNA samples, the two input masses were analyzed individually and jointly using the same analysis approach as described above. The correlations of allelic log<sub>2</sub> fold-change between the main female hippocampal TRAP experiment and pilot are shown comparing the respective *Vglut1*<sup>+</sup> analyses, as well as comparing each total hippocampal sequencing set (66ng, 150ng, or combined analysis), are shown in **Supplementary Figure 5.6**. The corresponding  $P_{emp}$ -derived FDR values and log<sub>2</sub> fold-changes from the main experiment and each of the four replication analysis sets are provided in the additional datasets linked on bitbucket in *Acknowledgments*.

***Analysis of functional SNPs for enrichment in TF motifs.*** To assess TF binding sites potentially disrupted by functional variants, I first generated sets of positive (functional) and negative (non-functional) variants for the analyses. For single-sex, single age/tissue analyses, I defined functional SNPs as those with an uncorrected  $P_{emp} < 0.05$ , and the remainder of the measured SNPs as non-functional. I also performed two comparisons for each age/tissue to identify TFs enriched at sex-genotype interaction SNPs. One comparison assessed interaction SNPs at  $P_{emp} < 0.05$  to “non-functional” SNPs ( $P_{emp} > 0.05$  for both allele and sex-allele interaction) to maximize the size of the negative set used in enrichment analysis. I additionally compared functional interaction SNPs to those with only a significant main effect allele effect, so as to identify TFs potentially involved

in sex-divergent variant effects. The enrichment procedure required a negative set of greater size than the positive set (a random set of negatives equal to the size of positives was drawn each iteration, see below); to meet this condition, I allowed for more lenient definition of allele-only effects from the LMM.

I performed motif perturbation analyses for all SNPs designed into the MPRA utilizing two tools: 1) the R package motifbreakR<sup>65</sup> and its built-in database of motif position-weight matrices (PWMs) from multiple public repositories, and 2) RSAT var-tool<sup>66</sup> with each of 3 motif databases; its own 2017 database, comprised of clustered motifs (some without TFs assigned) based on similarities across multiple motifs for multiple TFs<sup>115</sup>, JASPAR 2020 TF-specific motifs<sup>116</sup>, and cisBP 2017's human database<sup>117</sup>, which consists of both specific TF motifs and more general motifs not assigned to any one TF. Both tools are designed to identify PWM matches overlapping input SNPs in dbSNP (version 151 in hg38 for motifbreakR) or Ensembl (hg37, for RSAT) for which at least one of the SNP alleles results in a genomic sequence significantly matching a given motif sequence. I used the default significance cutoff of  $p < 10^{-4}$  for calling motif matches in all analyses and identified changes in motif match score using the tools' default algorithms. MotifbreakR considers a weighted sum based on the position weights of each base in the motif sequence and considers these for the two alleles of the query SNP; the magnitude of these differences is used to classify motif perturbations as "strong" or "weak". For motifbreakR, I performed separate enrichment analyses only considering those changes classified as strong, and regardless of the algorithm's classification. RSAT performs a similar analysis, identifying the strongest position-weighted p-value match to a motif for each allele of a SNP within a sequence of user-defined length (here, I used 122bp flanks to approximate the 110-126bp sequences assayed

in the MPRA) under a first-order (dinucleotide) background frequency model, and reporting the best match p-value for each allele to each motif where at least one allele exceeds the defined cutoff for the best match value.

Frequencies at the level of TF (which can include several motifs) were considered as the number of SNPs matched to a given TF, regardless of the number or identity of motifs to which that SNP matched. Null distributions of frequency were determined by 50,000 random selections of  $n$  SNPs of motif perturbations identified in the negative SNP set, where  $n$  was the number of positive set SNPs analyzed. Due to incomplete compatibility of hg37 and UK Biobank rsIDs with hg38/dbSNP 151, only 1277 SNPs were actually analyzed by motifbreakR; 1452 were analyzed by RSAT; the number of positive SNPs analyzed, and negative SNPs drawn in permutations were based on the number of SNPs actually analyzed in each respective tool. The p-value of frequency was then calculated from the empirical percentile of the positive SNP frequency count vs the distribution of frequencies in the negative sets, and these were corrected using standard FDR correction within each individual analysis, with resulting significant enrichments considered as  $FDR < 0.05$ . I additionally logged whether each TF/motif was depleted in the positive SNP set relative to the permuted negative sets.

Finally, some results from the RSAT analyses using the RSAT 2017 and cisBP human 2017 databases corresponded to motif sequences not ascribed to particular TFs (by nature of those databases). I separated these enrichment results from those for which a TF was explicitly listed in motifbreakR or RSAT. To predict corresponding transcription factors for the undefined motifs, I utilized the MEME-suite tool TomTom<sup>118</sup> to predict significant matches (using Euclidean distance)

of the database motif to TF-specific motifs across 4 databases: JASPAR CORE Vertebrates 2018 (non-redundant)<sup>119</sup>, Jolma 2013<sup>120</sup>, Mouse Uniprobe<sup>121</sup>, and HOCOMOCO v11 (human and mouse motifs)<sup>122</sup>. Corresponding TFs were assigned for all TomTom matches at  $p < 10^{-4}$ ; for cisBP/RSAT motifs where no TomTom match achieved this p-value, the single-lowest match p-value under 0.01 was retained.

***Gene set enrichment analyses of motif-enriched TFs.*** From the above analyses, I generated lists of unique TFs identified across the 3 RSAT and motifbreakR analyses for each condition (9 hippocampal TF sets total--Vglut1 sex-genotype interaction, hippocampus sex-genotype interaction, both interaction types combined, and allelic rSNPs from each sex in each of tissue fraction and in both tissue fractions for each sex). For the neurodevelopmental conditions, I generated one TF set per age•sex, and one TF set for sex-by-allele effect variants from the P0 condition.

To narrow the hippocampal TF sets down to those most likely present (expressed) and thus able to exert regulatory activity in the adult hippocampus, I collected the publicly available GTEX v8 transcripts per million (TPM) expression dataset and subsetted to hippocampal samples (see *data availability* below), averaging the genewise TPM values across all samples (for interaction TF filtering) or against single-sex sample sets (for single-sex allele-effect TF filtering). From those TFs significantly enriched ( $FDR < 0.05$ ) from the analyses above—including those identified by matching nonspecific database motifs to putative TFs with TomTom—I filtered down to those with  $\geq 3$  average TPM in the respective GTEX hippocampal sample set. I did not perform any expression filtering of the P0 or P10 TF sets as comparable whole brain datasets do not exist. I

then utilized Enrichr to identify gene sets across ontologies, pathways, and drug perturbations where the TFs were enriched as a set (only considering enrichments of reported q-value < 0.01 and driven by  $\geq 3$  input genes if either the input gene or the result gene was a retinoid receptor or sex hormone receptor).

***Permutation tests of sex-by-allele interactions.*** In order to determine the null expectation for the rate of significant sex-genotype interactions, I performed 1,000 iterations of each sex-by-genotype linear mixed model, wherein the sample labels were randomly shuffled by sex. The same linear model, including the 50,000-iteration empirical p-value calculation step as described above, was performed for each permutation to ensure that the empirical p-values were of the same granularity as the experimental analyses. The 20% FDR cutoff for the empirical p-values were determined for each iteration, and the number of SNPs significant for a sex-allele interaction at this threshold were recorded for each iteration. The end result was a vector of 1,000 numbers of FDR 20% “sex-genotype” significant rSNPs from the permutations, constituting a null distribution of the number of interaction effects for a given experiment. That was then compared to the actual number of interaction SNPs found with the true labels.

***Comparison of MPRA allele effects to MDD GWAS effects.*** Summary statistics from the Howard 2019 MDD meta-GWAS<sup>7</sup>, which included subjects from the UK Biobank, were obtained from the Psychiatric Genomics Consortium website and subsetted to variants measured in the MPRA experiments presented. The identity of genotyped SNPs in UK Biobank were obtained at <http://geneatlas.roslin.ed.ac.uk/>. MPRA results for each individual condition were then split into those SNPs genotyped or imputed (*i.e.*, not genotyped but with a summary statistic in the GWAS

results) and absolute MPRA allele effects from the condition were correlated to the GWAS summary statistic's absolute effect sizes.

### 5.5.2 Additional highlighted findings from TF analysis of rSNPs in P0 and P10 brain

In P0 brain, I identified 20 TFs enriched at female rSNPs and 43 at male rSNPs, 9 of which were shared (**Figure 5.6A**). Shared transcription factors included *EGR1/4*, *TBX1/15*, and *IRF9*. Male-specific transcription factors again included *ZBTB7A* and *NR4A3* (**Figure 5.6A**). Female-specific TFs included a variety of zinc finger TFs, Krüppel-like factors (*KLFs*), endothelial-developmental regulator *SOX17*, and the neurodevelopmental TF *SMAD3* (**Figure 5.6A**). Despite the absence of sex interactions in the P10 brain, I also found largely distinct sets of TFs in each sex at this age (only 4 shared) (**Figure 5.6A**). Among the 42 P10 male TFs, 38 were male-specific, including *RARG*—again supporting my *in vitro* findings of retinoid-interactivity<sup>27</sup> and the cortex-specific roles of retinoid receptors across brain development<sup>75</sup>— and several *KLF* members. (In contrast, *KLFs* were instead only enriched at P0 in rSNPs from *females*). The 15 female P10 TFs were predominantly core transcriptional machinery, including *YY1/2*, *GTF3C2*, *TAF1*, and in both sexes, *CTCF*. I then looked at these TF sets as before to identify convergent regulators and functions among them (**Figure 5.6B**). Of the few enriched annotations for male P0 TFs, I notably found that 4 corresponded to genes downregulated by *Rara* knockout, and 6 were putative regulatory targets of *PBX3*, which shows widespread subcortical and midbrain expression in E18.5 and P4 mouse brain<sup>123</sup>. Female P0 TFs were enriched in Allen Atlas expression signatures for cortical layers 1, 3, and 5, and, as found in hippocampus, in PPI targets of HDAC2 and ESR1. P10 male TFs showed the greatest extent of overlap (12/42) with gene-regulatory targets of *ZBTB7A*

and were enriched for PPI targets of HDAC2, RXRA and RARA. The only enrichment found for P10 female TFs was via a modest (3/15) set of *FOXA1* regulatory targets.

Given the absence of FDR-corrected sex-genotype interactions at P10, I only analyzed P0 interaction rSNPs for TF motif perturbation enrichment. TFs enriched at interaction SNPs were comprised largely of *EGR*, *KLF*, and *SP* family members, as well as *PAX5* and neurodevelopmental factor *SMAD3* (**Figure 5.6A**). Annotation of SxG TFs revealed a broader extent of hormonal roles in functional variation than observed in either sex alone at P0: I found SxG TFs were again enriched for PPI targets of ESR1 (as had been P0 female TFs), but additionally enriched in PPI targets of AR and ESR2 (**Figure 5.6B**).

### 5.5.3 Annotation datasets and other outside datasets

•**GTEX v8 TPM** [https://storage.googleapis.com/gtex\\_analysis\\_v8/rna\\_seq\\_data/GTEX\\_Analysis\\_2017-06-05\\_v8\\_RNASeQCv1.1.9\\_gene\\_tpm.gct.gz](https://storage.googleapis.com/gtex_analysis_v8/rna_seq_data/GTEX_Analysis_2017-06-05_v8_RNASeQCv1.1.9_gene_tpm.gct.gz) ; de-identified metadata used to identify hippocampal samples <https://www.ebi.ac.uk/arrayexpress/files/E-MTAB-5214/E-MTAB-5214.sdrf.txt>

•**Hi-C contacts for dopaminergic neurons and cortical neurons:** [https://github.com/thewonlab/H-MAGMA/blob/master/Input\\_Files/Midbrain\\_DA.genes.annot](https://github.com/thewonlab/H-MAGMA/blob/master/Input_Files/Midbrain_DA.genes.annot) and [https://github.com/thewonlab/H-MAGMA/blob/master/Input\\_Files/Cortical\\_Neuron.genes.annot](https://github.com/thewonlab/H-MAGMA/blob/master/Input_Files/Cortical_Neuron.genes.annot)

•**Brain Hi-C contact matrices from Jung 2019:** [ftp://ftp\\_3div:3div@ftp.kobic.re.kr](ftp://ftp_3div:3div@ftp.kobic.re.kr)

•**Song 2019 *in vitro* neural cell type and fetal primary astrocyte Hi-C:** Corresponding paper's Supplementary Table 2

•**Song 2020 fetal radial glia, intermediate progenitor cell, excitatory neuron, and inhibitory neuron chromatin contacts:** files "iN.MAPS.peaks.txt", "IPC.MAPS.peaks.txt", "RG.MAPS.peaks.txt", and "eN.MAPS.peaks.txt" thru BDbag linked at <https://assets.nemoarchive.org/dat-uidoqy8b>

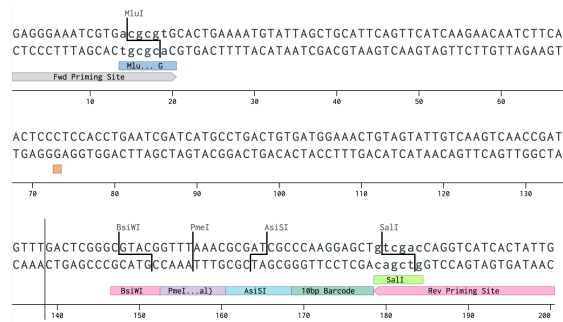
•**Su 2021 Hi-C data for neural tissues:** ebi.ac.uk with accession E-MTAB-9159

•**Fetal cortical plate and germinal zone Hi-C contacts from Won, 2016:** corresponding paper's supplementary tables S22 and S23.

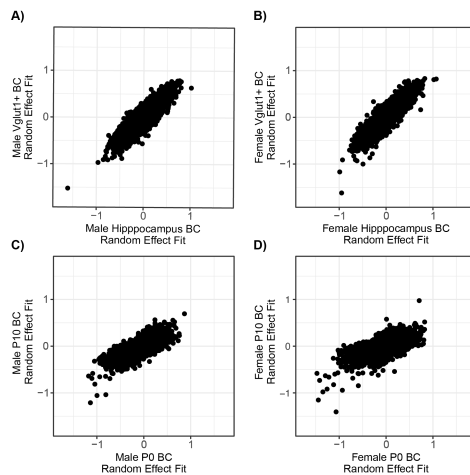
•**Howard 2019 MDD GWAS summary statistics:** entry "10.7488/ds/2458" on the Psychiatric Genomics Consortium's results download page, <https://www.med.unc.edu/pgc/download-results/>



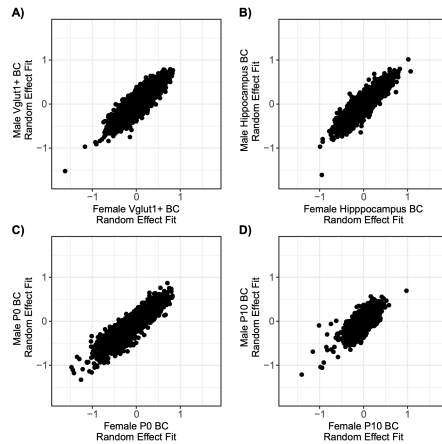
## 5.5.4 Supplementary figures and tables



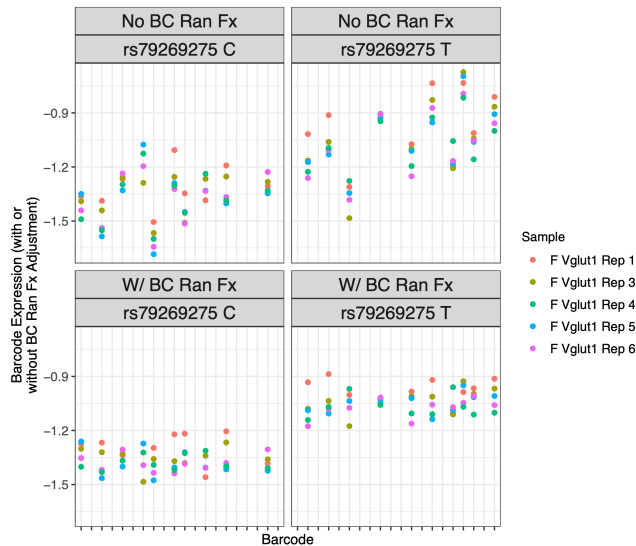
**Supplementary Figure 5.1. Template MPRA oligonucleotide.** Example MPRA oligonucleotide illustrating the features described including priming and cloning sites. The red box underlines the variant position at the center of the 126bp human genomic sequence tile. The PmeI site serves as a failsafe in the design such that, should a high fraction of plasmid not take on reporter constructs during cloning, digestion linearizes the reporter-negative plasmids and retransformation of the digested DNA results in isolation of reporter-positive plasmids. (This was not necessary for this library). Illustration captured from sequence design tools bundled with digital lab notebook, Benchling.



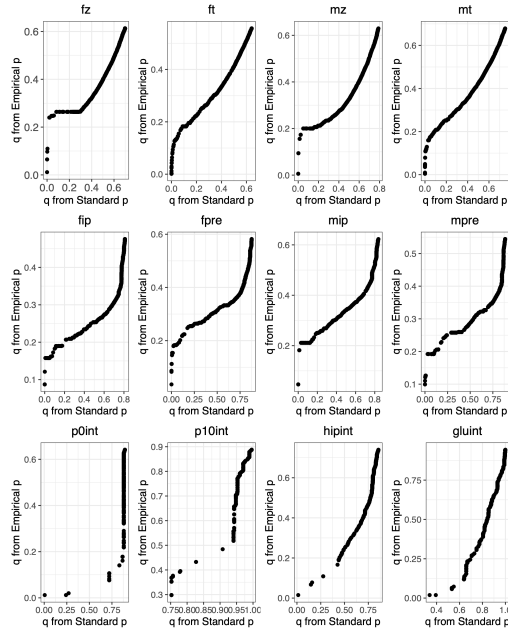
**Supplementary Figure 5.2. Barcode random effect coefficients are consistent within sex across ages/cell types.** A) Random effect coefficients from separate LMMs used to analyze male total hippocampus and male Vglut1+ data. B) *Ibid.* for female. C) Random effect coefficients from separate LMMs used to analyze male P0 and male P10 MPRA data. D) *Ibid.* for female.



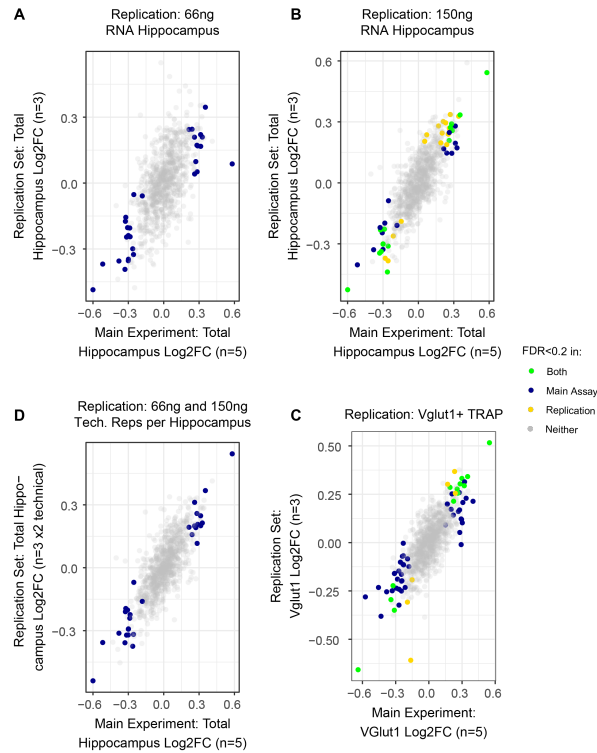
**Supplementary Figure 5.3. Barcode random effect coefficients are consistent between sexes.**  
**A)** Vglut1+ BC random effect coefficients from female vs. male. **B)** *Ibid.* for total hippocampus.  
**C)** *Ibid.* for P0. **D)** *Ibid.* for P10.



**Supplementary Figure 5.4. Example of barcode random effect fitting on expression values.**  
 Each X axis position is a barcode, with its expression level before or after adjustment for random effects shown by color for each replicate among female Vglut1+ samples.

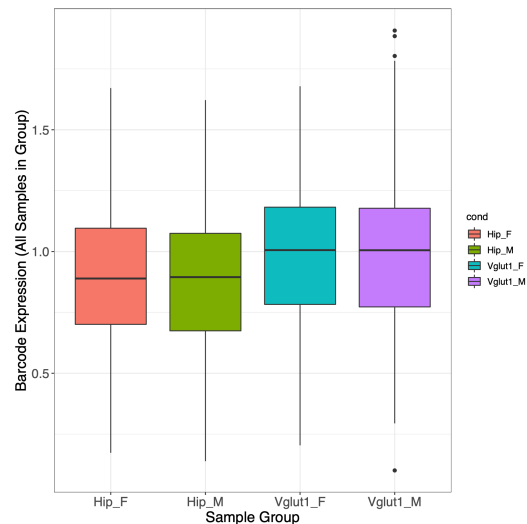


**Supplementary Figure 5.5. Comparison of Benjamini-Hochberg corrected p-values and empirical p-values of allelic differences in each single-sex analysis, and of SxG interactions in each interaction analysis.** Where X=M or F, representing male or female sex, Xz: P0; Xt: P10; Xhip: total hippocampus; Xglu: Vglut1+ hippocampal trap fraction; int: SxG interaction.

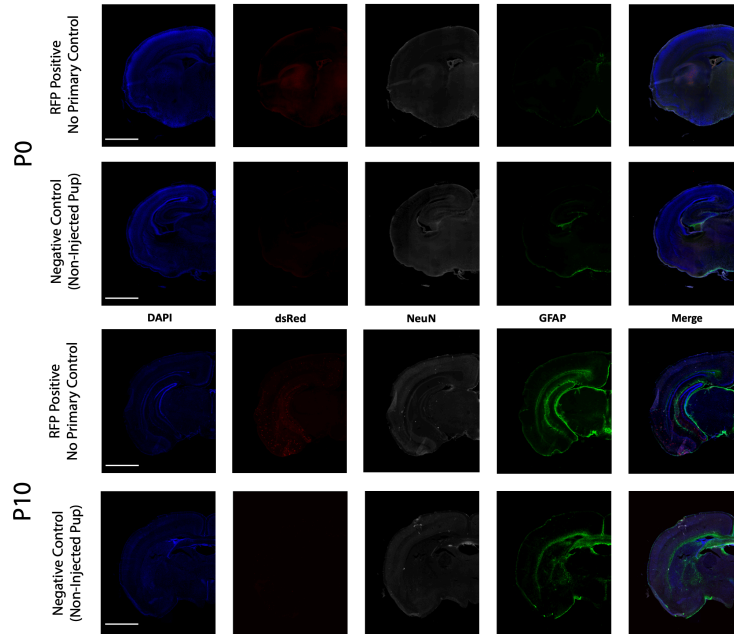


**Supplementary Figure 5.6. Log<sub>2</sub>FC between results from main adult female hippocampal-TRAP experiment and a validation set of three additional female hippocampal-TRAP sample sets.** Total hippocampal RNA was prepared for sequencing in parallel using 66ng or 150ng

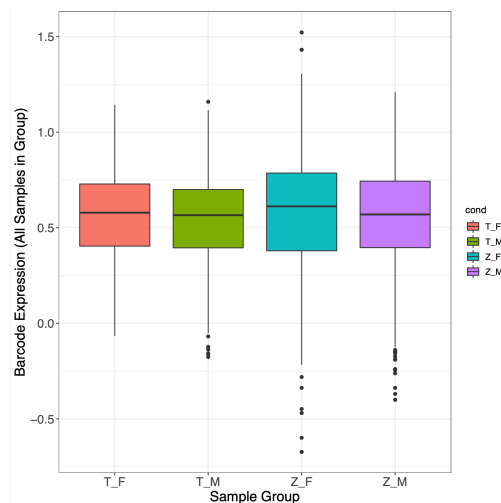
from each sample (see Sequencing analysis: All other experiments above). The sequencing results from each input mass were analyzed separately and jointly. **A)** log<sub>2</sub>FC correlations between main experimental total hippocampus and replication samples prepared using 66ng input RNA. Pearson correlation 0.6098. **B)** log<sub>2</sub>FC correlations between main experimental total hippocampus and replication samples prepared using 150ng input RNA. Pearson correlation 0.7721. **C)** log<sub>2</sub>FC correlations between main experimental total hippocampus and replication samples when analyzing the 66ng and 150ng input RNA technical replicates together. Pearson correlation 0.7699. **D)** log<sub>2</sub>FC correlations between main experimental Vglut1+ and replication Vglut1+ samples. Pearson correlation 0.7549.



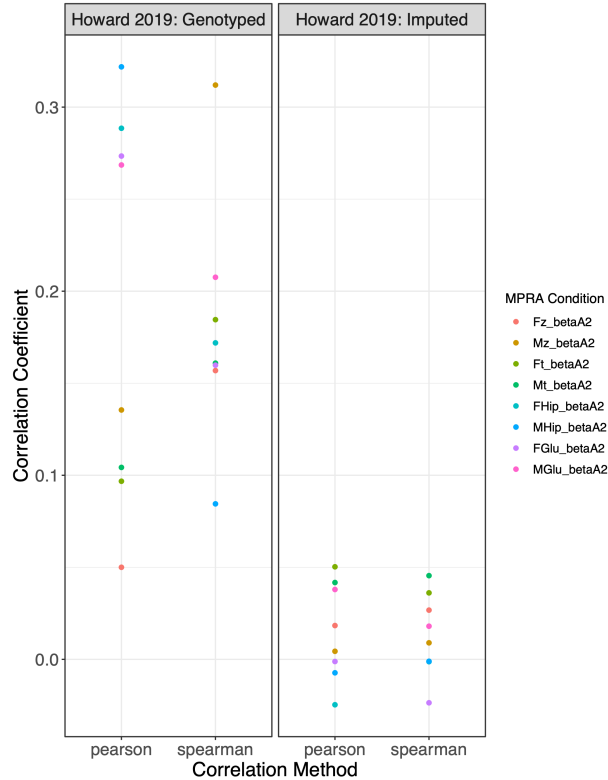
**Supplementary Figure 5.7. Basal (minimal promoter alone) barcode expression values do not vary by sex in total hippocampus or *Vglut1*<sup>+</sup> TRAP.**  $p > 0.5$  for both sex comparisons using student's  $t$ -test;  $\geq 99$  BC expression values shown per sample type.



**Supplementary Figure 5.8. IF negative controls for P0 and P10 AAV9 delivery.** **A)** An RFP-positive P0 brain without primary antibodies applied. **B)** An RFP-negative P0 brain with primary antibody staining (i.e., demonstrating that the dsRed primary antibody is selectively marking the AAV-delivered reporter per the main figure IFs). **C)** An RFP-positive P10 brain without primary antibodies applied. **D)** An RFP-negative P10 with primary antibody staining.



**Supplementary Figure 5.9. Basal (minimal promoter alone) barcode expression values do not vary by sex in P0 or P10 whole brain.** P0 sex comparison  $p=0.14$  (student's t-test); P10 sex comparison  $p=0.11$  (student's t-test).



**Supplementary Figure 5.10. Absolute MPRA SNP effects correlate more strongly to MDD GWAS summary statistic effects for SNPs directly genotyped than for those imputed in the Howard 2019 meta-GWAS of MDD.** SNPs were divided into those genotyped in UK biobank samples (which comprised a portion of the cases and controls in the Howard 2019 GWAS) and those for which genotypes were imputed. Absolute GWAS effects and absolute SNP effects from the MPRA were then correlated (as the direction of effect on disease risk and on gene expression may not be the same). All SNPs with an effect measurement in both the GWAS and this study were considered, representing 792-865 SNPs with imputed genotypes in GWAS and measured by MPRA, or 51-56 SNPs with measured genotypes in GWAS and measured by MPRA. Fz/Mz: female or male, P0; Ft/Mt: female or male, P10; FHip/MHip: female or male total hippocampus; FGlu/MGlu: female or male TRAP (Vglut1+).

**Supplementary Table 5.1. Sequencing preparation steps and data QC.** Table describing sequencing preparation parameters: starting mass of RNA or DNA, number of PCR cycles for each amplification step, and number of resulting samples for sequencing. Sequencing outcomes provided are the number of retained (QC-passing) samples for analysis, and read depth for RNA and DNA samples in each sequencing group. \*: For P0 and P10 expression calculations, the mean CPM of these 7 DNA samples was used. \*\*P0 RNA samples A8, A9, A10, A12, A16, A17, A18, and A22 were prepared and sequenced in both of these runs; barcode CPM values for twice-sequenced samples were averaged to obtain a single value per samples. Raw sequencing files and barcode counts for each run individually are included in the GEO dataset.

Preparation and Sequencing Batch	Input DNA per technical replicate in ng (# of tech. reps)	Input RNA per sample (ng)	# of cycles (cDNA / AAV PCR)	# of cycles (index PCR)	N total RNA	N analyzed RNA	RNA avg # reads (x 10 <sup>6</sup> )	DNA avg # reads (x 10 <sup>6</sup> )	RNA avg map rate (± SD)	DNA avg map rate (± SD)
Proof-of-principle TRAP-MPRA	1.44 (1)	15	17 for RNA, 12 for AAV DNA	9	6	4	0.4	0.9	41.8 (± 22) %	91.90%
Adult hippocampal TRAP	68 (4)	68.3 (3.5-6ng for 3 samples)	18 (RNA and DNA)	13	24	20	26	35	92.6 (± 0.6) %	94.4 (± 0.1) %
P0 samples alone	10 (3*)	1500	19	12	15	11**	42	14	86.3 (± 7.8) %	94.9 (± 0.5) %
P0 samples and P10 samples	9.25 (4*)	1500	18 cycles (P0), 19 cycles (P10), 16 cycles (DNA)	11	43	23**	37.5	19	84.8 (± 8) %	94.3 (± 1.3) %
3 additional P0 samples	NA	1000	19	12	4	3	94	NA	91.2 (± 4) %	NA

**Supplementary Table 5.2. qPCR primers used for *Vglut1*<sup>+</sup> TRAP validation.**

Primer Name	Primer sequence (5' to 3')
Actb R	CAATAGTGATGACCTGGCCGT
Actb F	AGAGGGAAATCGTGCGTGAC
Snap25 F	CAACTGGAACGCATTGAGGAA
Snap25 R	GGCCACTACTCCATCCTGATTAT
Gfap Fw	aaccgcatcaccattct
Gfap R	cgcattctccacagtctttacc
Gria1 F	CAAGTTTTCCCGTTGACACATC
Gria1 R	CGGCTGTATCCAAGACTCTCTG
P2ry12 F	ATGGATATGCCTGGTGTCAACA
P2ry12 R	AGCAATGGGAAGAGAACCTGG

## 5.6 References

1. Salk, R. H., Hyde, J. S. & Abramson, L. Y. Gender Differences in Depression in Representative National Samples: Meta-Analyses of Diagnoses and Symptoms. *Psychol Bull* 143, 783–822 (2017).
2. Marcus, S. M. *et al.* Gender differences in depression: Findings from the STAR\*D study. *J Affect Disorders* 87, 141–150 (2005).
3. LeGates, T. A., Kvarta, M. D. & Thompson, S. M. Sex differences in antidepressant efficacy. *Neuropsychopharmacol Official Publ Am Coll Neuropsychopharmacol* 44, 140–154 (2018).
4. Labonté, B. *et al.* Sex-specific transcriptional signatures in human depression. *Nat Med* 23, 1102–1111 (2017).
5. Nagy, C. *et al.* Single-nucleus transcriptomics of the prefrontal cortex in major depressive disorder implicates oligodendrocyte precursor cells and excitatory neurons. *Nat Neurosci* 1–11 (2020) doi:10.1038/s41593-020-0621-y.
6. Wray, N. R. *et al.* Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nat Genet* 50, 668–681 (2018).
7. Howard, D. M. *et al.* Genome-wide meta-analysis of depression identifies 102 independent variants and highlights the importance of the prefrontal brain regions. *Nat Neurosci* 22, 343–352 (2019).
8. Levey, D. F. *et al.* Bi-ancestral depression GWAS in the Million Veteran Program and meta-analysis in >1.2 million individuals highlight new therapeutic directions. *Nat Neurosci* 1–10 (2021) doi:10.1038/s41593-021-00860-2.
9. Blokland, G. A. M. *et al.* Sex-Dependent Shared and Non-Shared Genetic Architecture, Across Mood and Psychotic Disorders. *Biol Psychiat* (2021) doi:10.1016/j.biopsych.2021.02.972.
10. Trzaskowski, M. *et al.* Quantifying between-cohort and between-sex genetic heterogeneity in major depressive disorder. *Am J Medical Genetics Part B Neuropsychiatric Genetics* 180, 439–447 (2019).
11. Fullard, J. F. *et al.* An atlas of chromatin accessibility in the adult human brain. *Genome Res* 28, 1243–1252 (2018).
12. Bryois, J. *et al.* Genetic identification of cell types underlying brain complex traits yields insights into the etiology of Parkinson’s disease. *Nat Genet* 1–12 (2020) doi:10.1038/s41588-020-0610-9.



13. Hauberg, M. E. *et al.* Common schizophrenia risk variants are enriched in open chromatin regions of human glutamatergic neurons. *Nat Commun* 11, 5581 (2020).
14. Dong, P. *et al.* Transcribed enhancers in the human brain identify novel disease risk mechanisms. *Biorxiv* 2021.05.14.443421 (2021) doi:10.1101/2021.05.14.443421.
15. Song, M. *et al.* Cell-type-specific 3D epigenomes in the developing human cortex. *Nature* 1–6 (2020) doi:10.1038/s41586-020-2825-4.
16. Polioudakis, D. *et al.* A Single-Cell Transcriptomic Atlas of Human Neocortical Development during Mid-gestation. *Neuron* 103, 785-801.e8 (2019).
17. Zhang, K. *et al.* A single-cell atlas of chromatin accessibility in the human genome. *Cell* (2021) doi:10.1016/j.cell.2021.10.024.
18. Schmaal, L. *et al.* Subcortical brain alterations in major depressive disorder: findings from the ENIGMA Major Depressive Disorder working group. *Mol Psychiatr* 21, 806–812 (2015).
19. Zhang, J. M., Tonelli, L., Regenold, W. T. & McCarthy, M. M. Effects of neonatal flutamide treatment on hippocampal neurogenesis and synaptogenesis correlate with depression-like behaviors in preadolescent male rats. *Neuroscience* 169, 544–554 (2010).
20. Isgor, C. & Sengelaub, D. R. Effects of neonatal gonadal steroids on adult CA3 pyramidal neuron dendritic morphology and spatial memory in rats. *J Neurobiol* 55, 179–190 (2003).
21. Kohli, M. A. *et al.* The Neuronal Transporter Gene SLC6A15 Confers Risk to Major Depression. *Neuron* 70, 252–265 (2011).
22. Oliva, M. *et al.* The impact of sex on gene expression across human tissues. *Science* 369, eaba3066 (2020).
23. Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants using a Multiplexed Reporter Assay. *Cell* 165, 1519–1529 (2016).
24. Lu, X. *et al.* Global discovery of lupus genetic risk variant allelic enhancer activity. *Nat Commun* 12, 1611 (2021).
25. Choi, J. *et al.* Massively parallel reporter assays of melanoma risk variants identify MX2 as a gene promoting melanoma. *Nat Commun* 11, 2718 (2020).
26. Doan, R. N. *et al.* Mutations in Human Accelerated Regions Disrupt Cognition and Social Behavior. *Cell* 167, 341-354.e12 (2016).
27. Mulvey, B. & Dougherty, J. D. Transcriptional-regulatory convergence across functional MDD risk variants identified by massively parallel reporter assays. *Transl Psychiat* 11, 403 (2021).

28. Lambert, J. T. *et al.* Parallel functional testing identifies enhancers active in early postnatal mouse brain. *Biorxiv* 2021.01.15.426772 (2021) doi:10.1101/2021.01.15.426772.
29. Shen, S. Q. *et al.* A candidate causal variant underlying both higher intelligence and increased risk of bipolar disorder. *Biorxiv* 580258 (2019) doi:10.1101/580258.
30. Werling, D. M. *et al.* Whole-Genome and RNA Sequencing Reveal Variation and Transcriptomic Coordination in the Developing Human Prefrontal Cortex. *Cell Reports* 31, 107489 (2020).
31. Kouakou, M. R. *et al.* Sites of active gene regulation in the prenatal frontal cortex and their role in neuropsychiatric disorders. *Am J Medical Genetics Part B Neuropsychiatric Genetics* (2021) doi:10.1002/ajmg.b.32877.
32. Schork, A. J. *et al.* A genome-wide association study of shared risk across psychiatric disorders implicates gene regulation during fetal neurodevelopment. *Nat Neurosci* 22, 353–361 (2019).
33. Consortium, C.-D. G. of the P. G. *et al.* Genomic Relationships, Novel Loci, and Pleiotropic Mechanisms across Eight Psychiatric Disorders. *Cell* 179, 1469-1482.e11 (2019).
34. Arnold, A. P. & Breedlove, S. M. Organizational and activational effects of sex steroids on brain and behavior: A reanalysis. *Horm Behav* 19, 469–498 (1985).
35. Shen, S. Q. *et al.* Massively parallel cis -regulatory analysis in the mammalian central nervous system. *Genome Res* 26, 238–255 (2016).
36. Hrvatin, S. *et al.* A scalable platform for the development of cell-type-specific viral drivers. *Elife* 8, e48089 (2019).
37. Kim, J.-Y., Grunke, S. D., Levites, Y., Golde, T. E. & Jankowsky, J. L. Intracerebroventricular viral injection of the neonatal mouse brain for persistent and widespread neuronal transduction. *J Vis Exp Jove* 51863 (2014) doi:10.3791/51863.
38. Dougherty, J. D. *et al.* Candidate pathways for promoting differentiation or quiescence of oligodendrocyte progenitor-like cells in glioma. *Cancer Res* 72, 4856–68 (2012).
39. Harris, J. A. *et al.* Anatomical characterization of Cre driver mice for neural circuit mapping and manipulation. *Front Neural Circuit* 8, 76 (2014).
40. Li, X. *et al.* Common variants on 6q16.2, 12q24.31 and 16p13.3 are associated with major depressive disorder. *Neuropsychopharmacol* 43, 2146–2153 (2018).
41. Hyde, C. L. *et al.* Identification of 15 genetic loci associated with risk of major depression in individuals of European descent. *Nat Genet* 48, 1031–1036 (2016).

42. Power, R. A. *et al.* Genome-wide Association for Major Depression Through Age at Onset Stratification: Major Depressive Disorder Working Group of the Psychiatric Genomics Consortium. *Biol Psychiat* 81, 325–335 (2017).
43. Cai, N. *et al.* Sparse whole-genome sequencing identifies two loci for major depressive disorder. *Nature* 523, 588–591 (2015).
44. Ren, H. *et al.* Genes associated with anhedonia: a new analysis in a large clinical trial (GENDEP). *Transl Psychiat* 8, 150 (2018).
45. Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder. *Nat Genet* 51, 431–444 (2019).
46. Smith, D. J. *et al.* Genome-wide analysis of over 106 000 individuals identifies 9 neuroticism-associated loci. *Mol Psychiatr* 21, 1644 (2016).
47. Luciano, M. *et al.* Association analysis in over 329,000 individuals identifies 116 independent variants influencing neuroticism. *Nat Genet* 50, 6–11 (2017).
48. Meier, S. M. *et al.* Genetic Variants Associated With Anxiety and Stress-Related Disorders. *Jama Psychiat* 76, 924–932 (2019).
49. Ward, J. *et al.* The genomic basis of mood instability: identification of 46 loci in 363,705 UK Biobank participants, genetic correlation with psychiatric disorders, and association with gene expression and function. *Mol Psychiatr* 1–9 (2019) doi:10.1038/s41380-019-0439-8.
50. White, M. A., Myers, C. A., Corbo, J. C. & Cohen, B. A. Massively parallel in vivo enhancer assay reveals that highly local features determine the cis-regulatory function of ChIP-seq peaks. *Proc National Acad Sci* 110, 11952–11957 (2013).
51. Inoue, F., Kreimer, A., Ashuach, T., Ahituv, N. & Yosef, N. Identification and Massively Parallel Characterization of Regulatory Elements Driving Neural Induction. *Cell Stem Cell* (2019) doi:10.1016/j.stem.2019.09.010.
52. Storey, J. D. & Tibshirani, R. Statistical significance for genomewide studies. *Proc National Acad Sci* 100, 9440–9445 (2003).
53. Jaffe, A. E. *et al.* Profiling gene expression in the human dentate gyrus granule cell layer reveals insights into schizophrenia and its genetic risk. *Nat Neurosci* 23, 510–519 (2020).
54. Coleman, K. M., Lam, V., Jaber, B. M., Lanz, R. B. & Smith, C. L. SRA coactivation of estrogen receptor-alpha is phosphorylation-independent, and enhances 4-hydroxytamoxifen agonist activity. *Biochem Bioph Res Co* 323, 332–8 (2004).
55. Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the human brain. *Science* 362, eaat8464 (2018).

56. Hu, B. *et al.* Neuronal and glial 3D chromatin architecture informs the cellular etiology of brain disorders. *Nat Commun* 12, 3968 (2021).
57. Sey, N. Y. A. *et al.* Chromatin architecture in addiction circuitry elucidates biological mechanisms underlying cigarette smoking and alcohol use traits. *Biorxiv* 2021.03.18.436046 (2021) doi:10.1101/2021.03.18.436046.
58. Song, M. *et al.* Mapping cis-regulatory chromatin contacts in neural cells links neuropsychiatric disorder risk variants to target genes. *Nat Genet* 51, 1252–1262 (2019).
59. Su, C. *et al.* 3D promoter architecture re-organization during iPSC-derived neuronal cell differentiation implicates target genes for neurodevelopmental disorders. *Prog Neurobiol* 201, 102000 (2021).
60. Jung, I. *et al.* A compendium of promoter-centered long-range chromatin interactions in the human genome. *Nat Genet* 51, 1442–1449 (2019).
61. Lu, L. *et al.* Robust Hi-C Maps of Enhancer-Promoter Interactions Reveal the Function of Non-coding Genome in Neural Development and Diseases. *Mol Cell* (2020) doi:10.1016/j.molcel.2020.06.007.
62. Bernabeu, E. *et al.* Sex differences in genetic architecture in the UK Biobank. *Nat Genet* 53, 1283–1289 (2021).
63. Mulvey, B. *et al.* Molecular and Functional Sex Differences of Noradrenergic Neurons in the Mouse Locus Coeruleus. *Cell Reports* 23, 2225–2235 (2018).
64. McGill, R., Tukey, J. W. & Larsen, W. A. Variations of Box Plots. *Am Statistician* 32, 12 (1978).
65. Coetzee, S. G., Coetzee, G. A. & Hazelett, D. J. motifbreakR: an R/Bioconductor package for predicting variant effects at transcription factor binding sites. *Bioinform Oxf Engl* 31, 3847–9 (2015).
66. Santana-Garcia, W. *et al.* RSAT variation-tools: An accessible and flexible framework to predict the impact of regulatory variants on transcription factor binding. *Comput Struct Biotechnology J* 17, 1415–1428 (2019).
67. Chen, C. *et al.* The transcription factor POU3F2 regulates a gene coexpression network in brain tissue from patients with psychiatric disorders. *Sci Transl Med* 10, eaat8178 (2018).
68. Xie, Z. *et al.* Gene Set Knowledge Discovery with Enrichr. *Curr Protoc* 1, e90 (2021).
69. Sharma, M. *et al.* hZimp10 is an androgen receptor co-activator and forms a complex with SUMO-1 at replication foci. *Embo J* 22, 6101–6114 (2003).

70. Li, X., Thyssen, G., Beliakoff, J. & Sun, Z. The Novel PIAS-like Protein hZimp10 Enhances Smad Transcriptional Activity. *J Biol Chem* 281, 23748–23756 (2006).
71. Williams, E. S., Mazei-Robison, M. & Robison, A. J. Sex Differences in Major Depressive Disorder (MDD) and Preclinical Animal Models for the Study of Depression. *Csh Perspect Biol* a039198 (2021) doi:10.1101/cshperspect.a039198.
72. Kawatake-Kuno, A., Murai, T. & Uchida, S. The Molecular Basis of Depression: Implications of Sex-Related Differences in Epigenetic Regulation. *Front Mol Neurosci* 14, 708004 (2021).
73. Chakrabarty, P. *et al.* Capsid Serotype and Timing of Injection Determines AAV Transduction in the Neonatal Mice Brain. *Plos One* 8, e67680 (2013).
74. Won, H. *et al.* Chromosome conformation elucidates regulatory relationships in developing human brain. *Nature* 538, 523–527 (2016).
75. Shibata, M. *et al.* Regulation of prefrontal patterning and connectivity by retinoic acid. *Nature* 598, 483–488 (2021).
76. Corradin, O. *et al.* Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome Res* 24, 1–13 (2014).
77. Li, S. *et al.* Regulatory mechanisms of major depressive disorder risk variants. *Mol Psychiatr* 1–20 (2020) doi:10.1038/s41380-020-0715-7.
78. Liang, D. *et al.* Cell-type-specific effects of genetic variation on chromatin accessibility during human neuronal differentiation. *Nat Neurosci* 24, 941–953 (2021).
79. Aragam, N., Wang, K.-S. & Pan, Y. Genome-wide association analysis of gender differences in major depressive disorder in the Netherlands NESDA and NTR population-based samples. *J Affect Disorders* 133, 516–521 (2011).
80. Cui, J. *et al.* FBI-1 functions as a novel AR co-repressor in prostate cancer cells. *Cell Mol Life Sci* 68, 1091–1103 (2011).
81. Zhang, L. *et al.* ZBTB7A Enhances Osteosarcoma Chemoresistance by Transcriptionally Repressing lncRNA LINC00473-IL24 Activity. *Neoplasia New York N Y* 19, 908–918 (2017).
82. Issler, O. *et al.* Sex-Specific Role for the Long Non-coding RNA LINC00473 in Depression. *Neuron* (2020) doi:10.1016/j.neuron.2020.03.023.
83. Dion-Albert, L. *et al.* Sex-specific blood-brain barrier alterations and vascular biomarkers underlie chronic stress responses in mice and human depression. *Biorxiv* 2021.04.23.441142 (2021) doi:10.1101/2021.04.23.441142.

84. Klein, J. C. *et al.* A systematic evaluation of the design and context dependencies of massively parallel reporter assays. *Nat Methods* 1–9 (2020) doi:10.1038/s41592-020-0965-y.
85. Bhaduri, A. *et al.* Cell stress in cortical organoids impairs molecular subtype specification. *Nature* 578, 142–148 (2020).
86. Kwasniewski, J. C., Fiore, C., Chaudhari, H. G. & Cohen, B. A. High-throughput functional testing of ENCODE segmentation predictions. *Genome Res* 24, 1595–1602 (2014).
87. Li, M. *et al.* Optimal promoter usage for lentiviral vector-mediated transduction of cultured central nervous system cells. *J Neurosci Meth* 189, 56–64 (2010).
88. Lee, Y., Messing, A., Su, M. & Brenner, M. GFAP promoter elements required for region-specific and astrocyte-specific expression. *Glia* 56, 481–493 (2008).
89. Sakers, K. *et al.* Astrocytes locally translate transcripts in their peripheral processes. *Proc National Acad Sci* 114, E3830–E3838 (2017).
90. Ghazi, A. R. *et al.* Design tools for MPRA experiments. *Bioinform Oxf Engl* 34, 2682–2683 (2018).
91. Doyle, J. P. *et al.* Application of a translational profiling approach for the comparative analysis of CNS cell types. *Cell* 135, 749–62 (2008).
92. Demontis, D. *et al.* Discovery of the first genome-wide significant risk loci for attention deficit/hyperactivity disorder. *Nat Genet* 51, 63–75 (2018).
93. Hill, W. D. *et al.* A combined analysis of genetically correlated traits identifies 187 loci and a role for neurogenesis and myelination in intelligence. *Mol Psychiatr* 24, 169–181 (2018).
94. Lee, J. J. *et al.* Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nat Genet* 50, 1112–1121 (2018).
95. Ghossaini, M. *et al.* Open Targets Genetics: systematic identification of trait-associated genes using large-scale genetics and functional genomics. *Nucleic Acids Res* 49, gkaa840- (2020).
96. Machiela, M. J. & Chanock, S. J. LDlink: a web-based application for exploring population-specific haplotype structure and linking correlated alleles of possible functional variants. *Bioinformatics* 31, 3555–3557 (2015).
97. Aguet, F. *et al.* Genetic effects on gene expression across human tissues. *Nature* 550, 204–213 (2017).
98. Hauberg, M. E. *et al.* Large-Scale Identification of Common Trait and Disease Variants Affecting Gene Expression. *Am J Hum Genetics* 101, 157 (2017).

99. Consortium, T. B. *et al.* Developmental and genetic regulation of the human cortex transcriptome illuminate schizophrenia pathogenesis. *Nat Neurosci* 21, 1117–1125 (2018).
100. Ng, B. *et al.* An xQTL map integrates the genetic architecture of the human brain's transcriptome and epigenome. *Nat Neurosci* 20, 1418–1426 (2017).
101. Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* 507, 455–61 (2014).
102. (DGT), F. C. and the R. P. and C. *et al.* A promoter-level mammalian expression atlas. *Nature* 507, 462–70 (2014).
103. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature* 489, 57–74 (2012).
104. Consortium, T. E. P. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* 583, 699–710 (2020).
105. Amiri, A. *et al.* Transcriptome and epigenome landscape of human cortical development modeled in organoids. *Sci New York N Y* 362, eaat6720 (2018).
106. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using RegulomeDB. *Genome Res* 22, 1790–1797 (2012).
107. Kreimer, A., Yan, Z., Ahituv, N. & Yosef, N. Meta-analysis of massively parallel reporter assays enables prediction of regulatory function across cell types. *Hum Mutat* 40, 1299–1313 (2019).
108. Matsui, A., Yoshida, A. C., Kubota, M., Ogawa, M. & Shimogori, T. Mouse *in Utero* Electroporation: Controlled Spatiotemporal Gene Transfection. *J Vis Exp* 3024 (2011) doi:10.3791/3024.
109. Hu, S., Yang, T. & Wang, Y. Widespread labeling and genomic editing of the fetal central nervous system by in utero CRISPR AAV9-PHP.eB administration. *Development* 148, dev195586 (2020).
110. Kwasnieski, J. C., Mogno, I., Myers, C. A., Corbo, J. C. & Cohen, B. A. Complex effects of nucleotide variants in a mammalian cis-regulatory element. *P Natl Acad Sci Usa* 109, 19498–503 (2012).
111. Rabani, M., Pieper, L., Chew, G.-L. & Schier, A. F. A Massively Parallel Reporter Assay of 3' UTR Sequences Identifies In Vivo Rules for mRNA Degradation. *Mol Cell* 68, 1083-1094.e5 (2017).

112. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26, 139–140 (2009).
113. McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res* 40, 4288–97 (2012).
114. Rieger, M. A. *et al.* CLIP and Massively Parallel Functional Analysis of CELF6 Reveal a Role in Destabilizing Synaptic Gene mRNAs through Interaction with 3' UTR Elements. *Cell Reports* 33, 108531 (2020).
115. Castro-Mondragon, J. A., Jaeger, S., Thieffry, D., Thomas-Chollier, M. & van Helden, J. RSAT matrix-clustering: dynamic exploration and redundancy reduction of transcription factor binding motif collections. *Nucleic Acids Res* 45, gkx314- (2017).
116. Fornes, O. *et al.* JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* 48, D87–D92 (2019).
117. Weirauch, M. T. *et al.* Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* 158, 1431–1443 (2014).
118. Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. Quantifying similarity between motifs. *Genome Biol* 8, R24 (2007).
119. Khan, A. *et al.* JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res* 46, D260–D266 (2017).
120. Jolma, A. *et al.* DNA-Binding Specificities of Human Transcription Factors. *Cell* 152, 327–339 (2013).
121. Newburger, D. E. & Bulyk, M. L. UniPROBE: an online database of protein binding microarray data on protein–DNA interactions. *Nucleic Acids Res* 37, D77–D82 (2009).
122. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of transcription factor binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids Res* 46, D252–D259 (2018).
123. Thompson, C. L. *et al.* A High-Resolution Spatiotemporal Atlas of Gene Expression of the Developing Mouse Brain. *Neuron* 83, 309–323 (2014).



## Chapter 6: Discussion

### 6.1 Overview of findings

This body of work has highlighted sex as a critical variable that shapes gene expression, behavior, and genetic risk effects relating to MDD and potential biological underpinnings thereof. Noradrenergic neurons of the locus coeruleus (LC) had previously been demonstrated to have sex-differential anatomy and responsivity to peptide signals like *CRF*, though the underlying extent of these observed differences in LC were unknown. Using TRAP, I identified that adult mouse LC shows sex differences in its transcriptional profile after social (single-housing) and physical (surgical) stress, while serotonin neurons of the raphe nuclei lack such sex differences. Importantly, behavior follow-up work demonstrated replicability of the *Ptger3* sex difference in group-housed mice solely after survival surgery, and electrophysiology indicated the same functional sex difference is present without isolation or surgical stress. The sequence motif analyses presented on LC genes with sex-differential expression had, tantalizingly, shown local enrichment of *rodentia*-specific B-family transposons that nonetheless constituted ‘conserved’ sequences across mammals up to humans (which possess an orthologous transposon family, *Alu* elements). An analysis in 2010 of rodent and primate genomes indicated that, despite their emergence after the rodent-primate evolutionary split, these orthologous transposons were ultimately inserted across the genome in very similar fashions<sup>1</sup>. As my LC motif analyses were limited to conserved bases across mammals including primates and humans, this finding suggests that orthologous *Alu* transpositions are present near the same genes in human. I will later discuss approaches that would constitute steps toward verifying human translatability of my findings on sex differences in mouse LC.

Subsequently, I identified functional variation across more than two dozen MDD loci both *in vitro* and, novelly, *in vivo* at multiple timepoints and in a cell type-specific manner. One key finding from both the *in vivo* and *in vitro* experiments was a breadth of functional variation per risk-associated linkage region. This strongly suggests that MDD risk loci contain far more than one “causal variant” each, contrasting with *a priori* assumptions of methods like genetic fine-mapping and GWAS-eQTL colocalization. Secondly, both experiments identified shared, enriched regulatory features across functional variation associated with MDD, including retinoid receptor enrichments in neuroblastoma, developing mouse, and adult hippocampal assays. This finding confirms my hypothesis of shared regulatory architecture across MDD risk loci. Finally, the majority of functional variation in MDD loci showed dependency on sex, age, retinoids, and/or cell type. Strikingly, a median of only 45% of the functional variants identified in each of the 11 *in vivo* analyses (8 single-sex-single-age/cell type analyses and 3 sex-allele interaction analyses) were functional in at least one of the three neuroblastoma analyses (allele effect, allele effect with vehicle treatment, and ATRA interaction), emphasizing the importance of context while also highlighting the impact of adapting psychiatric-genetic MPRA to the mouse brain. Finally, this work demonstrated that common regulatory variation is very much capable of interacting with biological sex, likely via acute sex hormone receptor activity, and in the case of MDD, is over-represented within hippocampus and its excitatory neurons, consistent with the hypothesis that sex-by-variant interactions underlie sex differences in genetically-mediated MDD risk.

The following sections discuss extending findings from mouse LC to the human brain and several follow-up experiments and future directions for MPRA of MDD GWAS loci.

## 6.2 LC sex differences: human validity and transcriptional-regulatory mechanisms

The clinical-translational value of my findings in the mouse LC depends on whether similar sex-differential gene expression (or regulatory mechanisms thereof) occur in human brain. There are several potential avenues for exploring this question. One is collection of male and female human LC samples by, for example, laser confocal microdissection, to characterize gene expression and test for sex differences. Lower-throughput options include use of human tissue samples for IF or *in-situ* sequencing<sup>2-4</sup> for the ~150 genes identified as LC-enriched and sex-differential in Chapter 2, along with, say, ~150 top-enriched candidate LC markers from the TRAP experiments (to verify their human applicability and identify LC neurons within the tissue), and an additional top ~500 genes in terms of sex-differential expression, both LC-enriched and not. All of these techniques are complicated, however, by the small number of LC neurons, even in human, and the LC as a structure running a long, thin, and angular path in all three dimensions, such that only a minority of the small population of LC occupies any section (of a thickness tractable for *in-situ* sequencing, laser confocal microdissection, etc.) along any canonical plane.

An indirect approach to determining whether these sex differences in gene expression/regulation are conserved in human would be to test the enriched motifs near these genes and, critically, *their orthologous human sequences*, using MPRA. These assays would: A) test for transcriptional-regulatory activity, which could be easily achieved in neuroblastoma cell lines, which are arrested in mid-differentiation to an adrenergic cell type<sup>5-7</sup>; B) contrast both the human and mouse motifs' activity in human- and mouse-derived neuroblastoma lines (explained below); and C) test for sex-differential activity (i.e., *in vivo* MPRA), which ideally would be assayed in the LC itself (and a comparator cell type— perhaps serotonergic neurons, given the close relationship at the level of

gene expression but with a lack of sex differences) by multiplexing TRAP and MPRA. If the previously identified motifs are indeed a signature of evolutionarily conserved regulatory sequences in catecholaminergic cell types, then the bare-minimum outcome of experiments A/B would be regulatory activity of both the mouse sequence as assayed in mouse neuroblastoma and likewise of the human ortholog assayed in human neuroblastoma. This would support conserved function—that is, the human sequence in the presence of human regulatory proteins serves a function similar to the mouse sequence in the presence of mouse regulatory proteins. Human ortholog sequences that additionally show similar activity in mouse neuroblastoma would then be appropriate candidates—as this finding would confirm their compatibility with mouse gene-regulatory architecture—for follow-up *in vivo* MPRA to test for sex-differential regulatory activity in mouse LC.

### **6.3 Further dissection of functional MDD risk variants**

The experiments described in Chapter 5 identify both functional and sex-interacting variants from MDD loci, but leave several key aspects of the identified functional variants unresolved. One unforeseen limitation to the study is the inability to define enhancers and repressors under an *hsp68* minimal promoter in neural cell types, as virtually all sequences assayed had lower activity compared to the *hsp68* minimal promoter alone. Additional areas for further investigation include direct demonstration of roles for functional variant-enriched regulatory sequences by perturbations of their cognate TFs in an MPRA setting. Likewise, a mechanistic demonstration of a role for sex hormones in SxG variants is warranted. Finally, recent advances in MPRA (see Section 6.3.5) are ripe for implementation in the *in vivo* brain and could constitute means to further assess the predicted roles of transcriptional regulators, identify additional environmental and organismal

variables that may shape functional variation and/or its interaction with sex, and explore variant effects across multiple brain cell types. These areas of follow-up study are each explored in-depth below.

### **6.3.1 Defining variant-perturbed enhancers and repressors and minimal promoter selection in MPRA**

One unforeseen limitation of the MPRA experiments—both in culture and mouse brain—was near-global repressive effects of sequences placed upstream of the *hsp68* minimal promoter relative to the inserts of minimal promoter alone. Given that all sequences in the MPRA were included on the basis of several (mostly enhancer-like) epigenomic signatures from neural tissues, these findings do not necessarily mean that the sequences are “repressive” in their native genomic context. Rather, it is likely that the *hsp68* promoter itself, and/or flanking regions of the DNA cassette, interact in a net-repressive manner with any sequence placed upstream of the promoter. Critically, this prevented determination of sequences with enhancer or repressor activity, which would ordinarily be identifiable by significant deviations from expression driven by the promoter alone.

The simplest approach to this problem is to replace the reporter cassette used in the presented experiments with another cassette containing a different minimal promoter. Studies using human neural progenitor cells have utilized a short “super core promoter”<sup>8–10</sup> and similar sequences to successfully discern enhancers and repressors, including in the setting of functional regulatory variants. Since the reporter is cloned in separately from the oligonucleotides to the plasmid, the final plasmid library could simply be re-digested at the reporter insertion sites, gel purified, and ligated with a new reporter cassette using the same cut sites as before.

Identifying sequences with enhancer or repressor activity in the assays presented would substantially aid interpretation and classification of results. First, MPRA of disease-associated variation in other tissues have only examined active regulatory elements as defined above for allelic effects of variants they contain<sup>11,12</sup>. This filtering approach arguably supports greater biological relevance of identified functional SNPs; that is, if the variant is in a sequence with clear *cis*-regulatory activity, it is more likely to be perturbing a *bona fide* biological role in the genome. Additionally, the ability to define enhancer and repressor sequences would improve resolution of subsequent analyses of shared regulatory mechanisms, as motifs could be separately identified in repressors and enhancers, providing information on particular regulatory mechanisms dysregulated in each direction. Likewise, such stratification would enable discernment of whether dual-functioning TFs, like retinoid receptors, are implicated in MDD by disruption of their repressive effects or their enhancing effects.

### **6.3.2 Validation/demonstration of predicted regulatory factor mechanisms**

One challenging aspect of *in vivo* MPRA—especially the sex-differential aspects of these experiments—is that the approaches used to validate MPRA findings *in vitro*, such as single-variant luciferase assays, are not nearly as readily implemented in the mouse brain. The following three sections discuss novel approaches that could be employed to maintain the *in vivo* contexts of the original assays while demonstrating roles for the predicted key regulatory systems and enabling refined parsing of variant effects.

### **6.3.3 Confirming the role for sex hormones in hippocampal SxG variant effects**

Cumulatively, the three MPRA performed *in vivo* implicated sex hormones as key mediators of functionality and sex-interactivity of MDD-associated variation in two key ways. First, sex-by-

genotype interactions were only observed at the two age points where there are sex hormones circulating (P0 and adulthood). Second, sex hormone receptors and TFs known to modulate their activity were enriched regulators and interactors among functional variants. To functionally confirm these findings, any number of well-established methods for modulating sex hormonal effects in the mouse brain could be leveraged. Chief among them is ovariectomy/gonadectomy of adults, which results in a substantial reduction of ordinarily circulating sex hormones.

Delivery of the same MPRA library in an extended experiment including intact mice of both sexes and gonadectomy and/or ovariectomy would enable re-assessment of variants in the absence of sex hormones. In a three-condition setting of intact males and females plus ovariectomized females, the roles of local testosterone conversion to estrogens (in males) could also be discerned. Given the original results, my primary prediction would be that ovariectomy reduces the number of female-specific *Vglut1*<sup>+</sup> rSNPs toward that of males. Sex-by-genotype interactions from intact animals could be altered in more complex ways, however. Sex interactions driven primarily by androgens should remain unaffected by ovariectomy, as the sex difference in androgen abundance remains essentially unaltered in this scenario. Interactions driven primarily by progestins, on the other hand, will likely be absent altogether given progestin reduction in ovariectomized females. Finally, and most interestingly, there may potentially be *novel* sex-genotype interaction variants revealed by comparison of wild-type males and ovariectomized females—specifically at loci where local androgen conversion to estrogen results in an estrogen-mediated regulatory effect. In this case, comparison of the intact sexes would not have revealed an interaction (assuming comparable levels of female systemic estrogen and male local conversion of estrogen). However, the female-specific reduction of estrogen by ovariectomy would reveal that the variant's effect

observed in the original experiments changes in magnitude or direction in the absence of estrogen, and thus result in an ovariectomy-vs-male specific sex interaction.

#### **6.3.4 Leveraging *in vivo* MPRA to explore consequences of other environmental factors on transcriptional-regulatory variants**

The adaptation of MPRA to the *in vivo* mouse brain may additionally enable assessment of the effects of other environmental factors relevant to MDD on associated variation. Preclinical studies on MDD and anxiety use, in fact, almost exclusively environmental manipulations to induce these phenotypes in rodents in the form of stressors. Delivery of the MPRA to the hippocampus or another brain region of interest—especially regions responsive to stressors—followed by stress and control conditions to induce depressive phenotypes and MPRA comparing variant effects across these two groups may serve to highlight means by which stress unmasks or potentiates risk variant effects on the way to depression. Other insights that could be gleaned from the *in vivo* context include antidepressant effects on functional variation—that is, to address whether antidepressants result in attenuated risk variant effects in the process of ameliorating the disease. (Since long-term antidepressant treatments do increase the time between depressive episodes, it is not implausible that effects of genetic risk are “suppressed” by these medications longitudinally.) One interesting experiment toward this end this end would be to study the sex difference in MDD treatment response—that is, that tricyclic antidepressants are more effective in men while SSRIs are more effective in women, though both drug classes target broadly-projecting, monoaminergic neuromodulatory nuclei. If de-functionalization of MDD risk variants is a component of successful treatment, then I would expect to see A) a general decrease in the number of functional variants (or in their comparative effect sizes) between mice receiving vehicle and either a tricyclic or SSRI. Moreover, given the human sex difference in treatment responsiveness, I would B) expect to see more



risk variants with attenuated effects or rendered non-functional in male hippocampus by chronic tricyclic intake compared to SSRI, while I would expect to see more attenuation of variant function in female hippocampus resulting from chronic SSRI.

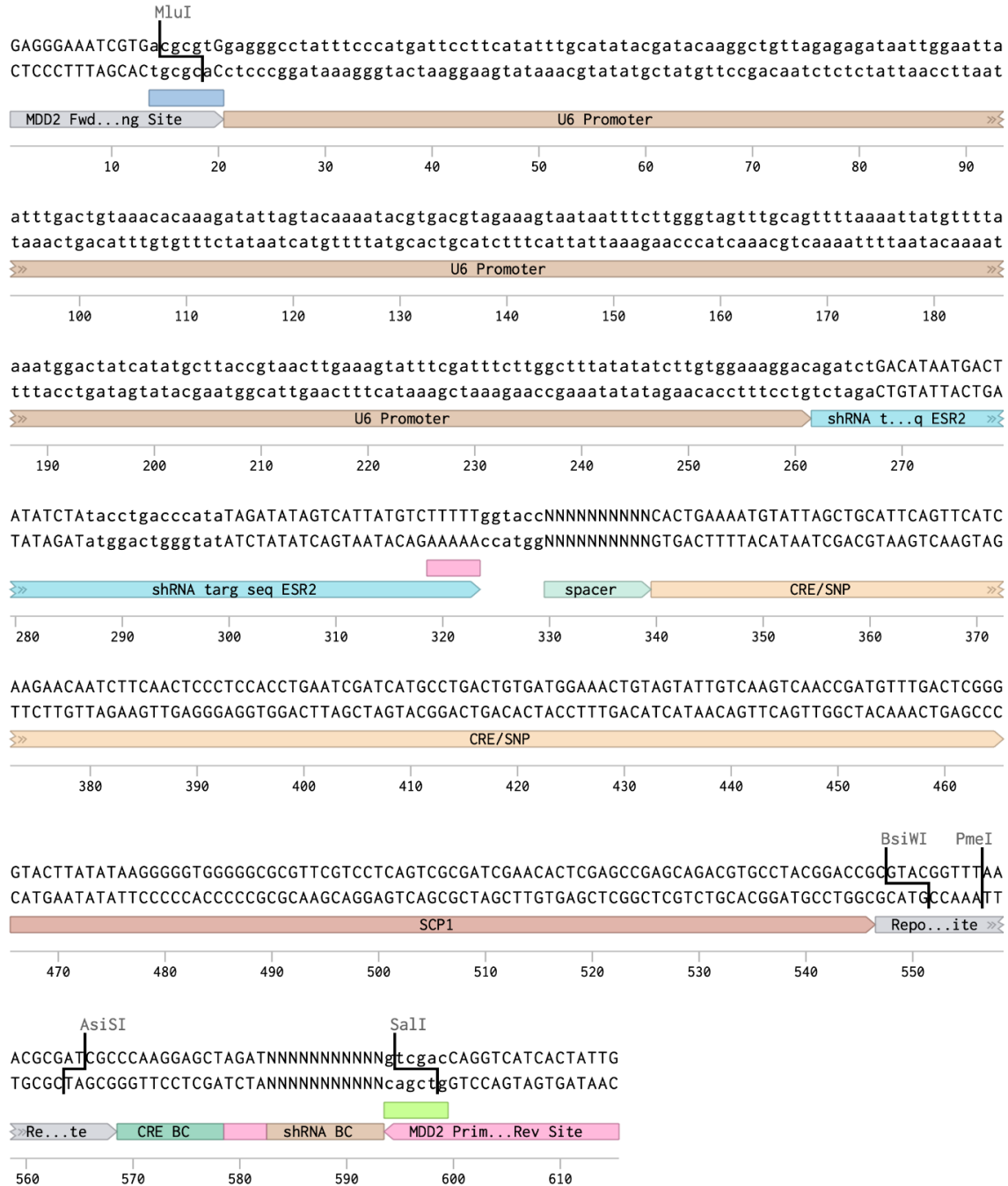
### **6.3.5 Single-cell MPRA: opportunities ahead**

A new frontier for MPRA is on the horizon with novel implementations that co-utilize single-cell RNA sequencing. This method uses a portion of each cell's RNA to identify its expressed genes, and the other portion to detect reporter barcodes<sup>13</sup>. In this way, cell typing and subsequent identification of functional variation in the same sample becomes possible, enabling dissection of functional variation in several cell types simultaneously. Performing these assays in a whole-brain manner, like the experiments presented at ages P0 and P10, would permit detection of *Vglut1*<sup>+</sup> neurons (even if not hippocampal per se) and specific detection of sex-genotype interactions within the cell type for comparison to adult *Vglut1*<sup>+</sup> neurons, for example. However, the cell type-specificity information regarding functional variants would be extended from one cell type and one tissue type per brain/brain region (*i.e.*, in TRAP) to *all* AAV-transducible cell types of reasonable abundance in the sequenced tissue, making it possible to determine virtually *all* relevant cell type contexts for functional variants from a similar number of samples.

Moreover, the single-cell MPRA approach could be leveraged to confirm roles for enriched motifs/TFs and their interactors by integrating co-cloned interfering RNAs, resulting in an MPRA library with built-in capacity to validate roles for candidate regulatory factors. For example, given ten top-enriched TF motifs or cross-TF co-regulators, a library of larger oligonucleotides could be constructed to contain the same sequences as assayed in the initial library here, but with an upstream RNA polymerase III transcription start site and a small hairpin RNA (shRNA) targeting

one of the ten candidate regulators (along with non-targeting controls). As a result, the library contains co-transcribed machinery to determine whether each of the candidate regulators is necessary for an observed allelic effect. Importantly, this type of multiplexing would require single-cell resolution in order to discern which shRNA(s) was (were) present in a given reporter-expressing cell. Excitingly, this could also facilitate observation of combinatorial regulator effects at candidate variants: for example, several functional variants directly perturbed predicted binding sites of *ZBTB7A*, which is in turn bound by the androgen receptor. A library containing shRNAs targeting *ZBTB7A* and *AR* in separate constructs across all assayed variants would result in cells with each TF repressed individually, as well as cells with both repressed, enabling examination of each TF's singular and combined effects at *ZBTB7A* site variants. Whether this would be feasible in one brain per biological replicate or require several would likely require empirical testing to determine the transduction depth per cell type per regulatory sequence and shRNA of interest (i.e. a minimum number of cells for a full-fledged experiment being,  $a$  MPRA barcodes  $\times$   $b$  shRNAs  $\times$   $c$  cells per cell type  $\times$   $n$  unique cell types to get strongly correlated results between replicates for each cell type of interest with low noise between barcodes).

Oligonucleotides for such a library could be produced by ordering oligo pools of 176bp, containing each intended 614bp product broken into four overlapping oligos with 30bp shared ends, enabling self-priming and extension into full-length target products using the multiplex pairwise assembly (MPA) method<sup>14</sup>. A schematic of such a hypothetical 614bp final product is shown in **Figure 6.1**.



**Figure 6.1 Example MPA oligo product for a hypothetical MPRA library with co-delivered shRNAs.** In this example, an shRNA targeting mouse estrogen receptor *Esr2* is present. The pink block at the 3' end of the shRNA indicates the RNA polymerase III termination site, a spacer, and the allelic sequence to be assayed. The remainder of the oligo contains the super core promoter 1 minimal promoter (SCP1) and the original internal cloning sites for reporter gene insertion. Adjacent to the genomic sequence-tagging barcode is a second barcode specific to the shRNA separately encoded upstream.

## 6.4 Conclusions

This body of work has illustrated several themes regarding MDD-pertinent genetic risk and sex differences. First, molecular sex differences with neurophysiological and behavioral consequences are present in the noradrenergic locus coeruleus of adult mice. Second, genetic risk loci for MDD are virtually all characterized by more than one functional variant per locus both *in vitro* and *in vivo*. Third, the function of MDD-associated variation is in no sense constitutive: which variants are functional and the nature of their effect varies with sex, age, retinoids, and cell type. In all, this work demonstrates that the physiologic environment—including cell types, signals, and sex—impact gene regulation, shaping transcriptomes and modulating the functional genetics of MDD risk.

## 6.5 References

1. Tsirigos, A. & Rigoutsos, I. Alu and B1 Repeats Have Been Selectively Retained in the Upstream and Intronic Regions of Genes of Specific Functional Classes. *Plos Comput Biol* 5, e1000610 (2009).
2. Chen, K. H., Boettiger, A. N., Moffitt, J. R., Wang, S. & Zhuang, X. Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* 348, aaa6090–aaa6090 (2015).
3. Maynard, K. R. *et al.* Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat Neurosci* 24, 425–436 (2021).
4. Moffitt, J. R. *et al.* Molecular, spatial and functional single-cell profiling of the hypothalamic preoptic region. *Science* 362, eaau5324 (2018).
5. Tao, T. *et al.* LIN28B regulates transcription and potentiates MYCN-induced neuroblastoma through binding to ZNF143 at target gene promoters. *P Natl Acad Sci Usa* 201922692 (2020) doi:10.1073/pnas.1922692117.
6. Banerjee, D. *et al.* Lineage specific transcription factor waves reprogram neuroblastoma from self-renewal to differentiation. *Biorxiv* 2020.07.23.218503 (2020) doi:10.1101/2020.07.23.218503.
7. Zimmerman, M. W. *et al.* Retinoic acid rewires the adrenergic core regulatory circuitry of neuroblastoma but can be subverted by enhancer hijacking of MYC or MYCN. *Biorxiv* 2020.07.23.218834 (2020) doi:10.1101/2020.07.23.218834.
8. Juven-Gershon, T., Cheng, S. & Kadonaga, J. T. Rational design of a super core promoter that enhances gene expression. *Nat Methods* 3, 917–22 (2006).
9. Ryu, H. *et al.* Massively parallel dissection of human accelerated regions in human and chimpanzee neural progenitors. *Biorxiv* 256313 (2018) doi:10.1101/256313.
10. Uebbing, S. *et al.* Massively parallel discovery of human-specific substitutions that alter enhancer activity. *Proc National Acad Sci* 118, e2007049118 (2021).
11. Lu, X. *et al.* Global discovery of lupus genetic risk variant allelic enhancer activity. *Nat Commun* 12, 1611 (2021).
12. Choi, J. *et al.* Massively parallel reporter assays of melanoma risk variants identify MX2 as a gene promoting melanoma. *Nat Commun* 11, 2718 (2020).
13. Zhao, S., Hong, C. K., Granas, D. M. & Cohen, B. A. A single-cell massively parallel reporter assay detects cell type specific cis-regulatory activity. *Biorxiv* 2021.11.11.468308 (2021) doi:10.1101/2021.11.11.468308.

14. Klein, J. C. *et al.* A systematic evaluation of the design and context dependencies of massively parallel reporter assays. *Nat Methods* 1–9 (2020) doi:10.1038/s41592-020-0965-y.