

Small RNA dynamics in cholinergic systems

DISSERTATION
ZUR ERLANGUNG DES DOKTORGRADES
DER NATURWISSENSCHAFTEN

VORGELEGT BEIM FACHBEREICH I4
DER JOHANN WOLFGANG VON GOETHE-UNIVERSITÄT
IN FRANKFURT AM MAIN

VON
SEBASTIAN LOBENTANZER
AUS SCHLÜCHTERN

FRANKFURT 2019
(D30)

©2019 – SEBASTIAN LOBENTANZER
ALL RIGHTS RESERVED.

Small RNA dynamics in cholinergic systems

ABSTRACT

Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetuer erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit, vel lacinia enim metus eu nunc.

Structure of the manuscript: methods in sans serif, boxes?, E notation

Contents

I	INTRODUCTION	1
1.1	Cholinergic Systems	1
1.2	Cholinergic Aspects of Disease	3
1.2.1	Alzheimer's Disease	3
1.2.2	Schizophrenia and Bipolar Disorder	5
1.2.3	Immunity	6
1.2.4	Neuroinflammation	7
1.2.5	Stroke	9
1.2.6	Circadian Aspects of Cholinergic Systems	10
1.2.7	Neurokines	11
1.3	Transcriptional Connectomics	13
1.3.1	Transcription Factors	14
1.3.2	microRNAs	15
1.3.3	Transfer RNA Fragments	17
1.4	Nested Multimodal Transcriptional Interactions - The Need for Connectomics	18
2	miRNetDB: CREATION OF A COMPREHENSIVE CONNECTOMICS DATABASE	21
2.1	Implementation	22
2.1.1	Neo4j: A Graph-Based Infrastructure	23
2.1.2	High-throughput Database Generation	24
2.1.3	Maintenance and Quality Control	24
2.2	Materials	25
2.2.1	Gene Annotation	25
2.2.2	microRNA Annotation	26
2.2.3	Transcription Factor Targeting	26
2.2.4	microRNA Interactions	27
2.2.5	Filtering of Aggregated Prediction Scores	29
2.2.6	De-novo Prediction of tRF Targeting	29
2.2.7	microRNA Primate Specificity	30
2.3	miRNetDB Usage	32
2.4	Statistical Approach to Transcriptional Connectomics	36
2.4.1	Permutation	36
2.4.2	Gene Set Enrichment Analysis	37
3	MICRORNA DYNAMICS IN CHOLINERGIC DIFFERENTIATION OF HUMAN NEURONAL CELLS	39
3.1	Neuronal Transcriptomes - Background	39
3.2	Cortical Single-Cell RNA Sequencing	41
3.3	microRNA and Transcription Factor Targeting Predictions	42
3.4	Gene Clustering Based On Expression	43
3.4.1	Co-Expression of Functional Groups of Cholinergic Genes	43
3.5	The Cellular Model	44
3.5.1	The SH-SY5Y Neuroblastoma Cell Line	44
3.5.2	The LA-N Neuroblastoma Cell Lines	45
3.5.3	Culture	46

3.5.4	Differentiation	46
3.5.5	RNA Isolation	47
3.6	Small RNA Sequencing and Differential Expression Analysis	49
3.6.1	Sequencing	49
3.6.2	Sequence Alignment	51
3.6.3	Differential Expression Analysis - R/DESeq2	51
3.6.4	microRNA Dynamics in CNTF-mediated Cholinergic Differentiation of LA-N-2 and LA-N-5	52
3.6.5	microRNA Family Enrichment	57
3.7	microRNA Family Gene Ontology Enrichment	57
3.7.1	Creation of miRNA Family Gene Target Sets	58
3.7.2	GO Analysis of Target Sets	58
3.7.3	Large Scale GO Term Curation	59
3.8	Whole Genome miRNA→Gene Network Generation	59
3.9	Application to Schizophrenia and Bipolar Disorder	60
3.9.1	Analysed Datasets	62
3.9.2	Microarray Quality Control and Data Preparation	62
3.9.3	Differential Expression Meta-Analysis	63
3.9.4	Sexual Dimorphism in Schizophrenia and Bipolar Disorder	64
3.9.5	Combination of Disease Data and Cell Culture	66
3.9.6	miR-125b-5p Acetylcholinesterase Targeting Assays	68
3.9.7	hsa-miR-125b-5p Targets Acetylcholinesterase	69
3.9.8	Cholinergic/Neurokine Mechanisms in Web-Available RNA Sequencing Experiments	69
4	DYNAMICS BETWEEN SMALL AND LARGE RNA IN THE BLOOD OF STROKE VICTIMS	71
4.1	RNA Sequencing, Differential Expression, and Descriptive Methods	71
4.1.1	The PREDICT Cohort	72
4.1.2	Clinical Parameters Collected in the PREDICT Study	72
4.1.3	Sample Collection, RNA Isolation, and Sequencing	72
4.1.4	RNA Sequencing Alignment	72
4.1.5	Quality Control and Filtering	73
4.1.6	RNA Sequencing Differential Expression Analysis	73
4.1.7	Gene Ontology Analyses	73
4.1.8	Homology Computation Among tRNA Fragments	74
4.1.9	t-Distributed Stochastic Neighbour Embedding	74
4.1.10	Cholinergic Association of Small RNA Species	74
4.2	Descriptive Analysis of RNA Dynamics in Blood After Stroke	74
4.2.1	Differential Expression of Large RNA	74
4.2.2	Gene Ontology Analyses of Differentially Expressed Genes	75
4.2.3	Differential Expression of small RNA	76
4.2.4	Homology Among tRNA Fragments	76
4.2.5	Cholinergic Association of Small RNA Species	76
4.3	Blood Compartments of Cholinergic Systems and Small RNA Species	78
4.3.1	Large RNA Regulatory Circuits in Tissues of the Blood	78
4.3.2	An Atlas of Small RNA Expression in Cell Types of the Blood	78
4.3.3	Large RNA Expression Patterns Identify Cholinergic Systems in CD14-positive Monocytes	80
4.3.4	Identification of Functional Enrichment of smRNA Expression in Blood-Borne Cells	80

4.3.5	Expression Patterns of Differentially Expressed and Cholinergic-Associated smRNAs	81
4.4	Regulatory Circuits of Small RNA and Transcription Factors in CD14-positive Monocytes	81
4.4.1	Comprehensive Circuit Network Creation	81
4.4.2	Gene Ontology Analyses of TF→Gene Networks of CD14-positive Monocytes	83
4.4.3	Dichotomy of Small RNA Targeting of Transcription Factors in CD14-positive Monocytes	83
4.4.4	Transcriptomic Footprints of Dichotomous Transcription Factors in CD14-positive Monocytes	84
4.4.5	miRNA-targeted Transcription Factors convey XX	85
4.4.6	tRF-targeted Transcription Factors convey XX	85
4.5	Dimensionality Reduction and Correlation of Expression with Clinical Parameters	85
4.5.1	WGCNA	85
4.5.2	Co-correlation	85
5	DISCUSSION	87
5.1	Methods	87
5.2	The Cholinergic/Neurokine Interface	88
5.3	Small RNA Therapeutics and Pharmacology	88
6	CONCLUSION	89
BIBLIOGRAPHY		91
A	TRANSCRIPTION FACTOR REGULATORY CIRCUITS - TISSUE TYPES	113
B	LIST OF PRIMATE-SPECIFIC HOMOLOGUES OF HUMAN miRNAs	115
C	MICRORNA DIFFERENTIAL EXPRESSION IN LA-N-2 AND LA-N-5	117
D	LIST OF GO TERMS FROM ANALYSIS OF DIFFERENTIALLY EXPRESSED LARGE RNA IN STROKE	119

THIS IS THE DEDICATION.

»Ever tried. Ever failed. No matter.

Try again. Fail again.

Fail better.«

Simon Beckett

Acknowledgments

Lorem ipsum consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

Abbreviations

- ACh** acetylcholine
AChE acetylcholinesterase (protein)
AD Alzheimer's Disease
Ago argonaute (protein)
API application programming interface
BD Bipolar Disorder
CA cholinergic-associated
CAGE 5' cap analysis of gene expression
cAMP cyclic adenosine monophosphate
ChAT choline acetyltransferase (protein)
CNS central nervous system
DAG directed acyclic graph
DE differentially expressed
DMEM Dulbecco's modified eagle medium
FCS fetal calf serum
FDR false discovery ratio
GEO Gene Expression Omnibus (NCBI)
GO Gene Ontology
gp130 see IL6ST (gene)
HLA-DR monocyte human leukocyte antigen isotype DR
iPSC induced pluripotent stem cell
KO knockout
LA-N-2 human neuroblastoma cell line (female)
LA-N-5 human neuroblastoma cell line (male)
LBP lipopolysaccharide binding protein
LDT laterodorsal tegmentum (Ch6)
LPS lipopolysaccharide
MBL mannan-binding lectin
miRNA microRNA
mRS modified Rankin Scale (clinical score of stroke severity)
NCBI National Center for Biotechnology Information

nt nucleotide
OR odds ratio
PBG parabigeminal nucleus (Ch8)
PBS phosphate buffered saline
PCA principal component analysis
RT-qPCR real-time quantitative polymerase chain reaction
PD Parkinson's Disease
PPN pedunculo-pontine nucleus (Ch5)
REM rapid eye movement
RIN RNA integrity number (RNA quality measure)
RISC RNA-induced silencing complex
RPMI1640 Roswell Park Memorial Institute medium
SCN suprachiasmatic nuclei
SCZ Schizophrenia
RNA-seq RNA sequencing
smRNA small non-coding RNA
SQL structured query language
TF transcription factor
tiRNA transfer RNA half
TPM transcripts per million
tRF transfer RNA fragment
tRNA transfer RNA
t-SNE t-distributed stochastic neighbour embedding
UTR untranslated region
vAChT vesicular acetylcholine transporter (protein; from SLC18A3 gene)
WT wild type

GENE SYMBOLS

ACHE acetylcholinesterase
ACLY ATP citrate lyase
AIF1 allograft inflammatory factor 1 (microglia marker protein)
AKT Serine/Threonine Kinase 1 (also known as Protein Kinase B)
BAD BCL-2-associated agonist of cell death
BCL-2 B cell lymphoma 2
BDNF brain-derived neurotrophic factor

BMAL1 brain and muscle ARNT-like protein 1 (also known as ARNTL)

CHAT choline acetyltransferase

CHRNA7 nicotinic acetylcholine receptor subunit $\alpha 7$

CLOCK circadian locomotor output cycles kaput

CNTF ciliary neurotrophic factor

CNTFR ciliary neurotrophic factor receptor (soluble)

CRY cryptochrome

ERK extracellular-signal-regulated kinase

GFAP glial fibrillary acidic protein (central astrocyte marker)

IFN interferon

IL interleukin

IL6R interleukin 6 receptor (soluble)

IL6ST interleukin 6 signal transducer (membrane bound; also known as gp130)

JAK janus kinase

JNK JUN N-terminal kinase

LIF leukaemia inhibitory factor

LIFR leukaemia inhibiting factor receptor (soluble)

MAPK mitogen-activated protein kinase

MCL1 myeloid leukaemia cell differentiation protein 1

MHC major histocompatibility complex

NF- κ B nuclear factor 'kappa-light-chain-enhancer' of activated B-cells

NGF nerve growth factor

NGFR nerve growth factor receptor (also known as p75)

NPAS2 neuronal PAS domain protein 2

NR1D1 nuclear receptor subfamily 1; group D; member 1 (also known as Rev-Erb α)

NR1F1 nuclear receptor subfamily 1; group F; member 1 (also known as ROR α)

NTRK1 neurotrophic receptor tyrosine kinase 1

NTRK2 neurotrophic receptor tyrosine kinase 2

OLIG1 oligodendrocyte transcription factor 1

PER period

PI3K phosphoinositide 3-kinase

RBFOX3 RNA-binding Fox-1 homolog 3 (neuronal marker gene; also known as NeuN)

RORA RAR-related orphan receptor α (see NR1F1)

SLC18A3 vesicular acetylcholine transporter (official gene symbol)

SST somatostatin

STAT signal transducer and activator of transcription

TGF transforming growth factor

TLR toll-like receptor

TNF tumour necrosis factor

TYK tyrosine kinase

VIP vasoactive intestinal peptide

1

Introduction

1.1 CHOLINERGIC SYSTEMS

NARY A PROCESS IN THE MAMMALIAN BODY CAN COMMENCE WITHOUT PARTICIPATION OF CHOLINERGIC SYSTEMS. Acetylcholine (ACh) was chemically and pharmacologically described by Henry Dale more than 100 years ago.¹ A short time later, Otto Loewi published the first proof of signal transmission by small molecules: he transferred physiological solutions from electrically stimulated frog hearts to naive hearts and observed their reactions; the solution that provoked a parasympathetic response he proposed to contain a »vagus substance«.² Finally, in 1929, Henry Dale completed the picture by isolating acetylcholine from mammalian tissue and identifying it as the molecule responsible for the parasympathetic response.³ Dale and Loewi's joint effort in »Discoveries Relating to Chemical Transmission of Nerve Impulses« was rewarded with the »Nobel Prize in Physiology or Medicine« in 1936.

Although we have learned much about cholinergic systems in these past 100 years, our understanding of the mammalian nervous system still is fairly limited. Even when disregarding peripheral nervous systems, the complexity of cholinergic transmission is immense, and a myriad functions have been attributed to cholinergic circuits in the central nervous system (CNS). Central nervous projections of cholinergic fibres were extensively mapped by Marek-Marsel Mesulam and others in the 1980s,^{4,5} with a majority of long projection neurons originating in one of the eight cholinergic nuclei, Ch1-Ch8. While many of these anatomical structures have been filled with meaning by associations with both rudimentary as well as higher brain functions, there are still as many cholinergic pathways whose function is entirely unclear (Figure 1.1, from Lobentanzer *et al.*⁶). This holds particularly

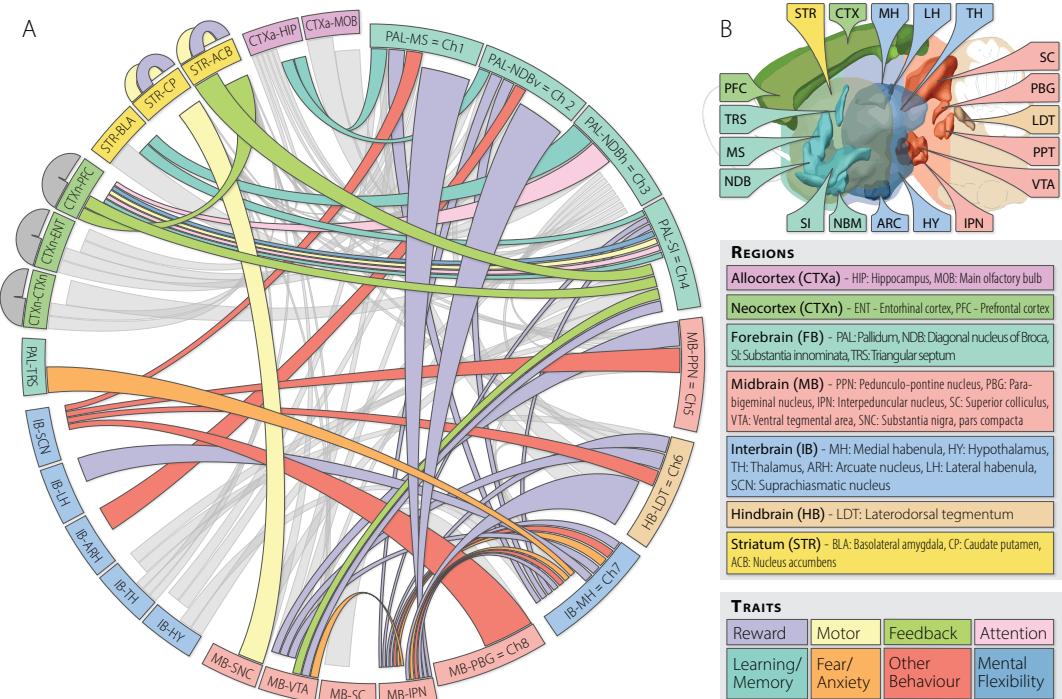


Figure 1.1: Cholinergic Projections in the CNS. Cholinergic systems are implicated in many diverse functional categories. **A)** The bulk of cholinergic projection neurons stems from one of the eight cholinergic nuclei, Ch1-Ch8 (right side of ideogram). They innervate wide areas of the mammalian CNS, and in turn receive incoming connections from all around the brain (left side of ideogram). Efferent connections are indicated by a small gap between ideogram and connector, in the first clockwise half of each ideogram component, afferent connection by a large gap, in the second half. A number of projections has been associated with specific functions, as implicated by the colours of the connectors. Two populations of cholinergic interneurons have been identified, in the striatum and the neocortex (outside of ideogram). **B)** Brain region and trait colour legend for A).

true for the only recently discovered cortical cholinergic interneurons, which, in comparison to their projecting counterparts, are very small and numerically vastly inferior to other neuron types in the cortex. Thus, their detection and analysis with current methods is challenging.

The histological definition of what constitutes a cholinergic neuron is not without debate. The staining procedures established in the 1970s utilised monoclonal antibodies against acetylcholinesterase (AChE),⁷ whose association with cholinergic neurons is not definitive, as it can be expressed post-synaptically as well (for an overview of genes of the cholinergic systems, see Box 1). Later on, developments in horseradish peroxidase systems allowed immunohistochemistry on choline acetyltransferase (ChAT), which is a more immediate marker of cholinergic neurons,⁴ albeit much more lowly expressed than AChE. However, AChE-based staining still was consistently used in addition to ChAT staining,⁵ sometimes without much differentiation. Recently, single-cell RNA sequencing (RNA-seq) allows a more detailed appreciation of the transcriptional diversity of neurons, and enables a clearer distinction between cholinergic and non-cholinergic neurons expressing AChE (see Section 3.2).

1.2 CHOLINERGIC ASPECTS OF DISEASE

CHOLINERGIC SYSTEMS ARE INTEGRAL for a myriad physiological functions, and as such they are critically involved in aetiologies and phenotypes of a number of central and peripheral diseases. Of interest to this dissertation are the cholinergic aspects of degenerative and non-degenerative central nervous diseases (such as Alzheimer's Disease, Bipolar Disorder, Schizophrenia), ischemic conditions in stroke, and peripheral modulation of immune responses, particularly in the context of the aforementioned diseases.

1.2.1 ALZHEIMER'S DISEASE

Alzheimer's Disease (AD) was characterised by Alois Alzheimer in 1906 and later named after him by his colleague and mentor, Emil Kraepelin.⁸ AD is a progressive neurodegenerative disease, its main risk factor is age, and it is estimated to make up 60-70% of all dementia cases. The disease incidence and progression distinctly differ between the sexes;⁹ generally, women are affected more often. Unlike the very rare familial form (that can affect patients in their fifties), spontaneous AD usually only begins to manifest symptomatically in the 6th to 7th life decade. As a result of the demographic change in most western countries, patient numbers, and thus, medical efforts, are expected to more than double in size until the year 2050. The cognitive decline associated with AD is progressive and ultimately leads to exhaustive care dependency; there is no cure.

The pathological hallmarks of AD are two types of atypical protein aggregates, inter-cellular amyloid β »plaques«, and intra-cellular »neurofibrillary tangles« composed of hyper-phosphorylated

Box 1: The Cholinergic Genes

Acetylcholine is synthesised from acetyl-CoA - supplied by **ATP citrate lyase** (*ACLY*) - and choline via enzymatic catalysis by **choline acetyltransferase** (ChAT from the *CHAT* gene). It is then packed into vesicles by the **vesicular acetylcholine transporter** (vAChT from the *SLC18A3* gene). After release into the synaptic cleft, it binds to a variety of **nicotinic and muscarinic receptors** (*CHRNx*, 16 subunits, and *CHRMx*, 5 subtypes). Of those, the nicotinic receptors form heteropentameric or, seldom, homopentameric ion channels, while the muscarinic receptors are monomeric G protein-coupled transmembrane receptors. The human possesses a duplicate of the nicotinic $\alpha 7$ receptor, **dup $\alpha 7$** (*CHRFAM7A*), which cannot bind ACh and supposedly acts as a dominant negative regulator of the $\alpha 7$ homomeric receptor. Termination of the signal is mainly achieved by **acetylcholinesterase** (AChE from the *ACHE* gene), one of the fastest enzymes known, with a theoretical rate of 25 000 molecules per second. AChE tetramers are usually tethered to cell membranes in the synaptic vicinity by the **proline-rich membrane anchor** (*PRIMA1*) or **collagen Q** (*COLQ*) peptides. Complementary to the mostly residual AChE is the circulatory **butyryl cholinesterase** (BChE from the *BCHE* gene), which can also nonspecifically degrade ACh. After degradation, residual choline is reimported into cells via the **high affinity choline uptake transporter** (HACU, from the *SLC5A7* gene).

tau-protein. Often, pathological aggregation of these proteins begins decades before the onset of symptoms. However, there have also been numerous cases of cognitively healthy subjects showing high amounts of protein aggregates. These inconsistencies and the unclear causality of pathology and symptoms have led to a redirection of scientific efforts to processes unrelated to amyloid and tau, such as neuroinflammation.

Cholinergic systems have long been associated with AD, as evidenced by the cholinergic hypothesis that was posed in the 1980s. To cite from Bartus *et al.*, 1982:¹⁰

»We have been guided by three deductive requirements that must be satisfied if the cholinergic hypothesis is to deserve continued attention: (i) specific dysfunctions in cholinergic markers should be found in the brains of subjects suffering from age-related memory loss, (ii) artificial disruption of central cholinergic function in young subjects should induce behavioural impairments that mimic the cognitive loss found naturally in aged subjects, and (iii) appropriately enhancing central cholinergic activity in aged subjects should significantly reduce age-related cognitive deficits.«

Although many reports substantiate all three deductive prerequisites, the cholinergic hypothesis has in the last decades been overshadowed by alternative hypotheses, particularly, amyloid-related theories. However, the therapeutic approaches developed along the lines of preventing amyloid β aggregation or otherwise clearing the plaques or soluble aggregates have not been successful in alleviating the cognitive decline in patients(cite). Thus, pro-cholinergic intervention by means of AChE inhibition still makes up the majority of approved drugs. The monotherapeutic approach of AChE inhibition is based on a premise that seems simple in light of the enormous complexity surrounding the interplay of the billions of neurons in the process of memory formation and recall, and in fact, pro-cholinergic therapy has only been shown to alleviate symptoms or delay their onset; a reversal of cognitive losses has so far not been achieved by any drug regime. As such, even in regard only to the cholinergic systems, novel approaches are sorely needed.

On the other hand, the cholinergic deficit in AD is not purely symptomatic. There is considerable debate whether the pathology originates from the entorhinal cortex or the basal forebrain. As Heiko and Eva Braak have shown,¹¹ AD pathology follows a characteristic brain region distribution process that can be stratified into stages and starts in the entorhinal region. The early stages of pathology by far precede the onset of symptoms. Taylor Schmitz and colleagues substantiate the cholinergic origin of neurodegeneration in their longitudinal *in vivo* imaging studies:^{12,13} In cognitively normal and impaired human subjects, basal forebrain volume predicted longitudinal entorhinal degeneration, but not *vice versa*. As such, the cholinergic basal forebrain dysfunction precedes as well as predicts pathology in other affected areas and cognitive deficits.¹² Additionally, the spread of Alzheimer's pathology in the longitudinal progression of the disease reflects the spatial topography of basal forebrain cholinergic projections.¹³

1.2.2 SCHIZOPHRENIA AND BIPOLAR DISORDER

Cognitive deficits can also occur without neuron death. The earliest description of what we today call Schizophrenia (SCZ) was coined by Emil Kraepelin: »*dementia praecox*«, premature dementia.¹⁴ Cognitive deficits are an integral but often overlooked part of the clinical picture of SCZ, which is often dominated by the more impressive positive symptoms such as hallucinations and paranoia. Likewise, Bipolar Disorder (BD) can present with cognitive impairments. The cognitive impairments affecting both SCZ and BD patients involve diminished problem solving capabilities and reduced intelligence, and are more pronounced in SCZ than in BD.¹⁵ They have been connected to cholinergic dysfunction^{16,17} and the sum of anticholinergic medications.^{18,19} A human polymorphism in the $\alpha 5$ nicotinic receptor subunit predicts a higher propensity for smoking and SCZ, showing parallel manifestations in engineered mice²⁰ and rats.²¹ Correspondingly, cholinergic stimulation can improve cognition^{22,23,24} and mood,²⁵ but can on the other hand provoke schizotypic behaviour in AD patients.²⁶

SCZ and BD clinically present with clear sexual dimorphisms. Compared to women, men have a higher SCZ prevalence with an odds ratio (OR) of 1.4, are affected earlier (at 15-25 as compared to 25-35 years of age), and face a worse prognosis.²⁷ Cholinergic participation also appears sex-dependent: Male SCZ patients more often self-medicate by smoking (7.2 versus 3.3 weighted average OR with 90% lifetime prevalence).²⁸ BD, on the other hand, affects men and women almost equally, with an OR of ~1. However, women make up 80-90% of so-called »rapid cyclers«, a subgroup of patients showing short intervals between manic and depressive phases which is associated with a worse prognosis.²⁹ Additionally, major depressive disorder, which is a prerequisite for BD diagnosis, more often affects women²⁹ (OR = 2).

Psychiatric genomics has recently identified a high amount of shared heritability between SCZ and BD.³⁰ Likewise, transcriptomic analyses have shown a high correlation (71%) between the transcriptional perturbations in the two diseases.³¹ Clinical as well as molecular pathology intensifies from BD to SCZ, suggesting the two lie on different points of a shared spectrum. However, their genetic origins are tremendously complex. Multiple disease-relevant markers have been identified by GWAS (genome-wide association studies), even able to distinguish between several sub-phenotypes of each disease.³² These markers are found in neurotransmitter receptors (e.g., dopaminergic, glutamatergic, cholinergic), scaffolding proteins (DISC: »disrupted in schizophrenia«), multiple transcription factors (TFs), microRNAs (miRNAs), and non-coding regions without known function.^{33,34,35}

Considering this complex disease aetiology, it is not surprising that there are no »designer« drugs available against SCZ and BD. All available therapeutic options have been empirically found, starting with the first antipsychotic, chlorpromazine, synthesised in 1950 by the French pharmaceutical company Rhône-Poulenc. Originally developed in a series dedicated to the search for antihistaminics, it was soon recognised for its antipsychotic potential, and widely prescribed only few years later. Many

other neuroleptic compounds have been derived from chlorpromazine, and through binding affinity assays, their receptor profiles were established. Most compounds with antipsychotic properties have a wide spectrum of different receptor activities, but most early drugs were strong antagonists of the D₂ dopamine receptor. Thus, the »dopaminergic hypothesis« of SCZ was formulated. However, aetiology as well as therapeutic principles are unclear to this day, and most antipsychotic substances still are very »dirty drugs«. In fact, newer developments leading to the discovery of the second generation (»atypical«) neuroleptic substances, starting with clozapine, have created molecules with an even wider spectrum of interactions and thus less specificity towards a single therapeutic mechanism of action. Similarly, the archetypal »mood stabiliser« Lithium, that has been found to ameliorate depressive as well as manic phases of BD, likely influences a wide variety of neuronal functions via mechanisms yet unclear.³⁶ This general development is contrary to pharmaceutical research practice, where most developments aim for a higher specificity. It is thus very likely that pharmacological therapy of SCZ and BD requires an approach consisting of multiple pharmacodynamic angles, to account for the multigenic disruption.

1.2.3 IMMUNITY

Aside from its vast neuronal functions, ACh also is highly relevant in immune cells, recently reviewed by Fujii and colleagues.³⁷ The first to isolate ACh from an animal organ, Dale and Dudley,³ used the spleens of oxen and horses. The spleen receives sympathetic, but not parasympathetic innervation, and as such, the ACh found by Dale and Dudley had to have come from immune cells. Indeed, nearly all mammalian immune cells express cholinergic components, most importantly, B- and T-cells, macrophages, and dendritic cells. While ACh is physico-chemically rather stable, it is extremely susceptible to enzymatic degradation, and cholinesterases are ubiquitarily distributed and cleave ACh with stunning efficiency, reducing its diffusion range to few millimetres. As a result, ACh has to be supplied synaptically or, at most, in paracrine fashion. ChAT activity has been confirmed in B- and T-cells, which both contain significant amounts of ACh, although T-cells generally possess higher amounts. Additionally, in peripheral cells ACh can be synthesised by the mitochondrial carnitine acetyltransferase. In addition to B- and T-cells, *CHAT* mRNA has been found in macrophages and dendritic cells. ChAT expression and ACh synthesis can be induced by various immune mediators, such as lipopolysaccharide (LPS) and toll-like receptor (TLR) agonists.

In addition to ACh synthesis, all of the aforementioned cell types can receive cholinergic signals. They express all muscarinic receptors as well as a selection of nicotinic receptor subunits, and the signal-terminating esterases. Although it is not completely clear how the parasympathetic signal reaches the immune cells, cholinergic activation as a result of inflammation can dampen the immune response in what is described as the »cholinergic anti-inflammatory reflex«.³⁸ This reflex loop is designed to protect the body from pathogens and inflammation, but also from the harmful effects of immune stimulation. Upon afferent signalling through the afferent vagus nerve and humoral compo-

nents, the brain releases humoral (via the hypothalamic-pituitary-adrenal axis) and neuronal (via the sympathetic and parasympathetic autonomous fibres) anti-inflammatory signals. The spleen has been identified as a pivotal organ in this response. Since none of the immune organs receives parasympathetic innervation, it has been proposed that the cholinergic activation is generated locally, with the help of sympathetic signalling to the organ.³⁹

A special role among cellular ACh receptors is occupied by the $\alpha 7$ nicotinic receptor subunit. Previously thought to exclusively form homopentameric ion-channel receptors, its functional characteristics have recently been extended. It has been found to form heteropentamers with $\beta 2$ subunits, akin to the prominent $\alpha 4\beta 2$ receptors in the brain, and an expression in immune cells also is likely. A hybrid duplication of the *CHRNA7* gene with *FAM7*, *CHRFAM7A*, is translated to a functional protein (dup $\alpha 7$). However, it seems to lack ACh binding ability, and thus is thought to act as a dominant negative regulator of $\alpha 7$ receptor function. In addition to its ionotropic function, mainly by means of Ca^{2+} transduction, the $\alpha 7$ receptor has been found to possess G-protein coupled metabotropic effects that can extend the duration of cholinergic activation. The $\alpha 7$ receptor can also, independently of Ca^{2+} , activate the JAK2/STAT3 pathway (see Section 1.2.7) in macrophages, leading to suppression of NF- κ B signalling.

On the other hand, cholinergic activation via M_1 and/or M_5 muscarinic receptors can lead to a positive immune response. The difference between muscarinic and nicotinic immune-signalling is elucidated by transgenic receptor knockout (KO) animals: Splenar cells from selective M_1/M_5 -KO mice secreted significantly lower amounts of the neuromodulators tumour necrosis factor (TNF)- α , interferon (IFN)- γ , and interleukin (IL)-6 than those from wild type (WT) mice. Conversely, antigen-stimulated splenar cells from $\alpha 7$ -KO mice produced significantly greater amounts of TNF- α , IFN- γ , and IL-6 than WT. In summary, the effects of cholinergic stimulation of the immune system is bidirectional and strongly context-dependent, and specific pharmacological intervention can shift homeostasis in both pro- as well as anti-inflammatory directions.

1.2.4 NEUROINFLAMMATION

Neurodegenerative as well as non-degenerative neurologic diseases are increasingly being associated with immunologic phenomena, prompting the need for integrative and translational approaches.⁴⁰ Transient and chronic inflammatory events can influence neuronal function and even survival in a dramatic fashion. Further, failure to resolve the acute inflammatory states might lead to maladaptive states, cases of »frustrated resolution«,⁴¹ in which the goal of adaptive immunity is not met. As was recently shown,⁴² resolution of inflammation is not just the »phasing-out« of inflammatory events, but rather bridges the gap between innate and adaptive immunity. While many acute-phase T_{H1} -type cytokines might have evolved to drive inflammation, their protracted influence may derail these post-inflammatory events and thus lead to maladaptive responses and chronic inflammation. T_{H1} -type cytokines include TNF, IFNs, IL-1 β and IL-6, and downstream mediators discussed in this

context are manifold: phosphoinositide 3-kinase (PI3K); cyclic adenosine monophosphate (cAMP); myeloid leukaemia cell differentiation protein 1 (MCL1); the complex of B cell lymphoma 2 (BCL-2), Serine/Threonine Kinase 1 (AKT), and BCL-2-associated agonist of cell death (BAD); all variants of mitogen-activated protein kinase (MAPK), i.e., extracellular-signal-regulated kinase (ERK) 1 and 2, JUN N-terminal kinase (JNK), and p38 MAPK; and the NF- κ B pathway. This list is not comprehensive, for a more detailed overview, see Fullerton & Gilroy.⁴¹

While none of the cardinal symptoms of inflammation are easily assessed in a CNS context, Virchow's fifth cardinal sign, *functio laesa*, is particularly difficult to tie to chronic neuroinflammation. Complex behavioural syndromes such as the functional deficits accompanying neurologic diseases might be influenced by protracted, maladaptive immunity, but the affected areas, brain structures, and timelines cannot be measured with current methods in neuroimaging. Only recently, it has become known that the brain is not immunologically pristine, but rather possesses a very specialised immune system, showing grave distinctions from, but also overlap with, peripheral immune systems.^{43,44} The mechanisms of immune privilege of the brain are constantly being refined; there is crosstalk between brain and periphery with blood-to-brain and brain-to-blood messaging,³⁹ and even migration of immune cells into the CNS, mainly as a response to sustained inflammation.

The first line of defence in CNS tissues are microglia, which in physiological state are resident ramified macrophages, and upon antigen sensing can produce an immediate native immune reaction. Further, the nascent immune system of the CNS comprises similar cells as the peripheral systems (T-cells, B-cells, NK-cells, dendritic cells), albeit with significant differences: antigen presenting cells express significantly fewer major histocompatibility complex (MHC) I and II molecules (which can however be induced by cytokine release upon inflammation); and the endocrine conditions (secretion of immune mediators from neurons) entail a more rapid response (seconds instead of days) with shorter duration of inflammation than in the periphery.⁴⁴

Under the surveillance of resident microglia, there is constitutive and inducible migration of immune cells between CNS tissues and periphery, in a loop of infiltration and drainage. Starting at antigen presentation inside the CNS, activated immune cells leave the CNS through one of two routes, both ending at the deep cervical lymph nodes: either through the cribiform plate into the nasal mucosa, or into the meningeal lymphatic vessels accompanying the sagittal and transversal sinuses in the *dura mater*. Following an immunological stimulus, activated immune cells in the deep cervical lymph nodes facilitate a secondary immune response and protect nervous tissue through secretion of cytokines⁴⁵ and re-migration of secondary immune cells to the CNS. Regulatory cytokines include IL-1 and IL-6, CCLs and CXCLs, leukaemia inhibitory factor (LIF), and epidermal and fibroblast growth factors. Re-migrating T cells can utilise various adhesion molecules expressed by endothelia along the blood-brain-barrier to cross into the CNS in a controlled fashion.⁴⁴

1.2.5 STROKE

Stroke is a medical emergency in which reduced blood flow leads to massive neuron death in the brain. There are two types of stroke: Haemorrhagic stroke, which is caused by bleeding due to a rupture of brain vessels, and ischemic stroke, misperfusion of a brain region caused by a clot in a cranial artery. Ischemic stroke makes up the vast majority of all strokes, about 90%. Stroke is currently the second most frequent cause of death in developed countries, only second to coronary artery disease. Those who survive the stroke are in many cases permanently and severely disabled. The prognosis of stroke patients is mainly dependent on clinical complications during initial care.

Infections are a leading cause of death in stroke patients, the CNS injury itself is an independent risk factor for the development of life-threatening infection. The most frequent complications accompanying stroke are fever and pneumonia, the fever in turn being most often caused by the infection. Affecting more than 20% of stroke patients, pneumonia is the most common serious post-stroke complication, featuring a mortality rate of more than 30%.⁴⁶ Stroke patients demonstrate a significant immunosuppression, resulting in lower count and functionality of immune cells. The ability of monocytes to synthesise cytokines is drastically reduced, a finding that has been reproduced in animal models of cerebral ischemia.

T cell activation and proliferation in the deep cervical lymph nodes is elevated following CNS injury, implicating that the drainage of immune cells from the site of injury plays an important role in immune system stimulation.⁴⁵ Depletion of those T cell leads to neuron death, suggesting that this naive response is favourable in stroke and similar conditions. However, the specifics of activation and migration, and the mediators (cytokines, antibodies) influencing the post-injurious response are still largely unclear.⁴³

The immunological situation after stroke is dominated by two opposing factors, the pro-inflammatory bodily response to injury and, often, infection, and the counter-regulatory immunodepressive response including the cholinergic anti-inflammatory reflex (see Section 1.2.3). The inflammatory response can become pathologic in the case of excess stimulation, which can result in »systemic inflammatory response syndrome« (SIRS), which in extreme cases can lead to shock and organ failure. As a counterbalancing measure, the body responds with a »compensatory anti-inflammatory response syndrome« (CARS), designed to allow fighting the infection while also protecting the body from excessive immune stimulation. However, in CNS injury, the anti-inflammatory component might overwhelm inflammatory processes, leading to a pathological »CNS injury-induced immunodepression syndrome«, CIDS.

In addition to the humoral and neurohumoral immunomodulatory pathways described above, the brain can directly steer immune processes by the release of cytokines. This might be particularly impactful in CNS injury, where blood-brain-barrier and homeostasis are disrupted. Contrary to the selective uptake of substances into the CNS, export from the CNS is mostly instantaneous and not

tightly controlled. Among the cytokines found in circulation after stroke are transforming growth factor (TGF)- β , IL-1 β , IL-6, and TNF- α . So far, it is unclear which of the immunomodulatory axes (humoral, neurohumoral, or direct release of cytokines from the brain) contributes to CIDS, and how they relate to each other. As a consequence, it is also unclear if directed intervention against CIDS would be beneficial, or even feasible.⁴⁶

1.2.6 CIRCADIAN ASPECTS OF CHOLINERGIC SYSTEMS

Cholinergic systems and psychiatric diseases share another common theme: regulation of circadian time and sleep patterns. Cholinergic nuclei have been associated with the resetting of the circadian clock in the suprachiasmatic nuclei (SCN). Retrograde tracing from the SCN⁴⁷ has identified basal forebrain nuclei (Ch1 & Ch4), as well as the pedunculo-ponitine nucleus (PPN, Ch5), laterodorsal tegmentum (LDT, Ch6), and the parabigeminal nucleus (parabigeminal nucleus (PBG), Ch8) as regulatory input to the SCN (compare also Figure 1.1). Basal forebrain cholinergic neurons are active in wakefulness and in the rapid eye movement (REM) phase of sleep,⁴⁸ and optogenetic activation of PPN and LDT cholinergic neurons (channelrhodopsin 2 under the ChAT promoter) during non-REM sleep was sufficient to induce REM sleep in mice.⁴⁹ Basal forebrain projections to other brain regions seem to functionally diverge from the projections to the SNC. In a study analysing prefrontal and hippocampal cholinergic activities, the increase in tonic ACh release during REM sleep was contingent on subsequent wakefulness,⁵⁰ and thus might convey a stronger »wake-up« signal than projections to the SCN alone.

Muscarinic receptors M1 and M3 are essential for REM sleep: REM-sleep is completely abolished in combined M1/M3 receptor KO mice.⁵¹ Arousal-induced phase shifts induced by activation of Ch4 cholinergic neurons projecting to the SCN were blocked in animals pretreated with (anti-muscarinic) atropine injections to the SCN, demonstrating that cholinergic activity at muscarinic receptors in the SCN is necessary for arousal-induced phase shifting.⁵² However, in their atropine perfusion experiment (locally via injection), the authors did not preclude cholinergic influences from the other nuclei.

In parallel, psychiatric diseases regularly exhibit symptoms of disturbed circadian rhythm. In the cholinergic-catecholaminergic imbalance hypothesis of BD,¹⁶ the imbalance follows variable transcriptionally regulated rhythms, and affected individuals exhibit decreased REM latency (the duration from onset of sleep to the first REM phase), which can be modulated by muscarinic agonists/antagonists.⁵³ Conversely, sleep deprivation exerts short-term antidepressant effects,⁵⁴ reduced cortical ACh levels,⁵⁵ and vast transcriptional changes in basal forebrain cholinergic neurons.⁵⁶

While the SCN regulates circadian timing of the organism, individual cellular timings are controlled by a group of transcriptional activators and de-activators, called clock genes. The autoregulatory feedback loop thus creates oscillates between day and night timing, under the influence of external factors. How exactly the individual cellular clocks are synchronised by the SCN is still unclear.⁵⁷

The first molecular circadian controller, circadian locomotor output cycles kaput (CLOCK), was identified by Joseph Takahashi and colleagues via mutagenesis screening in mice in 1997.⁵⁸ The transcription factors CLOCK and brain and muscle ARNT-like protein 1 (BMAL1) form heterodimers and bind to E-box elements in the promoters of period (PER) 1 and 2, and cryptochrome (CRY) 1 and 2, which lead to negative feedback regulation; or to E-box elements in the promoters of NR1D1 (giving rise to the Rev-Erb α protein) and NR1F1 (giving rise to ROR α), which compete for the ROR element in the BMAL1 promoter. ROR α induces BMAL1 expression, while Rev-Erb α represses it, thus leading to an oscillating expression pattern. In neurons, CLOCK can be substituted by its parologue NPAS2.

1.2.7 NEUROKINES

In comparison to the widely studied cholinergic projection neurons originating in the basal forebrain (Ch1-Ch4) that are known to depend on a retrograde survival signal by means of nerve growth factor (NGF), trophic influences on other cholinergic populations such as the cortical interneurons are unclear. NGF was described by Rita Levi-Montalcini in the 1950s as the first known instance of trophic peptides required for the survival of sympathetic ganglia.⁵⁹ The group of neurotrophic substances since discovered (most prominently, the brain-derived neurotrophic factor BDNF) are commonly referred to as »neurotrophins«. They convey their trophic effects through a family of transmembrane receptors; NGF binds to neurotrophic receptor tyrosine kinase 1 (NTRK1) with high affinity, BDNF binds to neurotrophic receptor tyrosine kinase 2 (NTRK2) with high affinity. However, both also bind to a third receptor, nerve growth factor receptor (NGFR), which is also known as p75, although with low affinity. NGFR function is complex, depending on the context it seems to be able to suppress as well as enhance the primary neurotrophic signal mediated by NTRK1/2(cite). The dependence of basal forebrain cholinergic neurons on retrograde NGF signalling was discovered in the 1980s.⁶⁰

A second group of trophic peptides with cholinergic implications are the so-called »neurokines«; the name results from the fact that this particular subgroup of cytokines has been associated with neuronal function in the central and peripheral nervous systems. Most prominently, they include the ciliary neurotrophic factor (CNTF), LIF, and IL-6, all of which coincidentally have been known by the acronym CDF. In the late 1980s, two groups of scientists (McManaman⁶¹ and Rao⁶²) independently identified proteins in extracts of muscle fibre that induced a differentiation of neurons towards a cholinergic type, and thus termed these proteins »choline acetyltransferase development factor« or »cholinergic differentiation factor« (both abbreviated CDF). Only later, through sequencing of the peptides, it became known that they had in fact discovered two distinct neurokines, LIF (Rao) and CNTF (McManaman, personal communication). IL-6, on the other hand, is abbreviated CDF for an entirely different reason: in this case it is short for »CTL (cytolytic T lymphocyte) differentiation factor«.

CNTF, LIF, and IL-6 convey their impact on neuronal activity through a partly redundant neurokine receptor pathway.²⁹ There are two basic types of neurokine receptors: soluble and transmembrane. The primary receptors for CNTF (CNTFR) and IL-6 (IL6R) are soluble proteins that are secreted into the extracellular space and, upon binding of a neurokine, bind to transmembrane receptor dimers on the cell surface. These transmembrane receptors are the LIF receptor (LIFR) and the »interleukin 6 signal transducer« (IL6ST), which is also known as gp130. Due to the latter's predominance, neurokines are also referred to as gp130 receptor family cytokines⁶³. Every neurokine has its preferred constellation of soluble and transmembrane receptors: CNTF binds to the soluble CNTF receptor and a dimer consisting of one gp130 and one LIFR protein; IL-6 binds to the soluble IL6R and a dimer of two units of gp130; LIF does not usually bind a soluble receptor but rather binds immediately to a dimer comprising one of each gp130 and LIFR; however, there are significant redundancy, pleiotropy, and crosstalk between those systems.^{64,63,65}

remove miRs
to later? All receptor constellations result in a main effect of activation of the JAK/STAT cascade (Fig. 1.2).

More specifically, neurokines can activate janus kinases (JAKs) 1 and 2 or the homologous tyrosine kinase (TYK) 2, and, successively, STAT (»signal transducer and activator of transcription«) isoforms 1, 3, 5A, and 5B, which then convey a multitude of cellular effects (e.g. in immunity or differentiation) through transcriptional activation. The STAT cascade is inherently self-limiting in that it usually leads to expression of transcription factors that serve as repressors of the STAT genes by SOCS (suppressors of cytokine signalling), PIAS (protein inhibitors of activated STATs), and PTPs (protein tyrosine phosphatases).⁶⁴

Neurokines, particularly IL-6, might serve as a link between the immunological and cholinergic aspects of physiological or disease processes. Since IL-6 is implicated in neurodegenerative, psychiatric, and injurious CNS diseases (Section 1.2), which all also possess a cholinergic facet, it makes sense to not only see it in the light of an immunomodulator, but also as a potential influence on neuronal function in cholinergic systems. A third of basal IL-6 levels are generated in adipose tissue in healthy humans,⁶⁶ and central (fatty) obesity increases risk for AD about 3-fold.^{67,68} SCZ and BD also are associated with obesity, although causality still is unclear. While obesity itself is more predominant in SCZ than in BD, obesity in BD patients is associated with decreased global cognitive ability as well as with poorer performance on individual tests of processing speed, reasoning/problem-solving, and sustained attention.⁶⁹ Low-grade chronic inflammation is recognised in obesity⁷⁰ as well as in neurodegenerative⁷¹ and non-degenerative psychiatric diseases.^{72,73} Sleep disturbance in animal models of mood disorders is accompanied by elevation in blood levels of IL-1, IL-6, and TNF- α .⁷⁴ Additionally, it has been shown that LIF can lead to a catecholaminergic-to-cholinergic neurotransmitter switch in peripheral neurons in a mouse model of protracted inflammation accompanying collagen-induced arthritis.⁷⁵ Though being a marginal phenomenon, it is not unthinkable that similar processes in central nervous cell might contribute to a disruption of homeostasis of cholinergic systems, and thus, to disease.

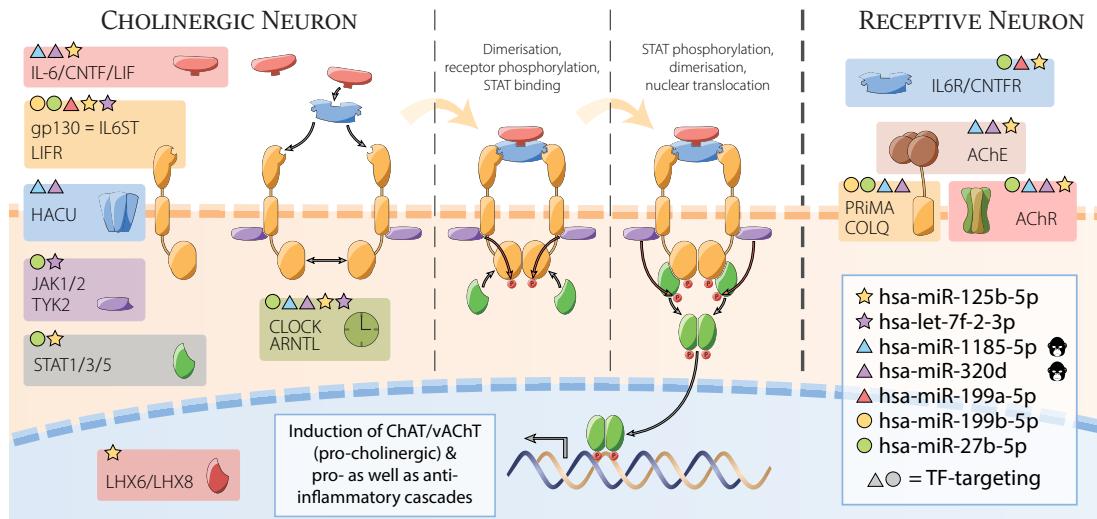


Figure 1.2: The Neurokine Pathway. The neurokines, such as CNTF, LIF, and IL-6, signal through a combination of soluble and membrane-bound receptors. Activation of a transmembrane neurokine receptor is usually followed by JAK recruitment and phosphorylation, and successively by STAT activation and translocation to the nucleus. Gp130-family neurokine, cholinergic, and circadian signalling pathways are controlled by primate-specific and evolutionarily conserved miRNAs. miRNA targeting of individual genes (indicated by coloured symbols) yields complex transcriptional interactions. Several miRNAs directly targeting the cholinergic pathway also target TFs controlling this pathway (circles and triangles).

1.3 TRANSCRIPTIONAL CONNECTOMICS

The term »connectomics« is not strictly limited to one scientific discipline; it is frequently used when the studied matter is defined by complex relationships between interaction partners. The most frequent use outside of transcriptional matters is neuronal connectomics, i.e., the relationships and projections between brain regions. In this dissertation, connectomics generally refers to epi-transcriptional interaction, the processes surrounding protein-coding gene expression. For the sake of simplicity, in this dissertation all descriptions of genomics and transcriptomics matters, of genes and their small RNA regulators, are to be seen in the context of *Homo sapiens*, unless explicitly stated otherwise.

NO MATTER THEIR LOCATION, CHOLINERGIC NEURONS ARE DEFINED BY THEIR ABILITY TO SYNTHESISE ACh AND RELEASE IT TO NEIGHBOURING CELLS TO A CERTAIN EFFECT. To fulfil this task, two particular proteins are essential: the choline acetyltransferase (ChAT) to synthesise ACh from choline and acetyl-Coenzyme A, and the vesicular acetylcholine transporter (vAChT, official gene symbol SLC18A3), which concentrates ACh in vesicles for later release. A notable genetic feature connects these two proteins beyond their functional association: the small *SLC18A3* gene - only 2420 nucleotides (nt) in size - sits inside the first intron of the *CHAT* gene and thus is already included in its primary transcript, and is subject to the *CHAT* promoter. However, oftentimes the (mature) transcript levels of *CHAT* and *SLC18A3* mRNA seem to be independently regulated; from the perspective of the organism, the possibility of differential regulation between these two genes makes sense. Since *SLC18A3* does not possess its own promoter, this differential regulation has to

be conveyed epigenetically.

This dissertation deals in large parts with approaches aiming to decipher these interactions; and while its primary topic revolves around cholinergic systems, the methods described in the following are designed to be applicable to the entirety of the genome/epigenome. Four particular types of cellular actors are subjects of these methods and therefore will be briefly introduced: genes in the classical sense as the conveyors of cellular function by encoding for proteins; TFs, a subclass of protein coding genes that are able to regulate the expression of other genes; miRNAs, a class of small non-coding RNA (smRNA) that has been known for approximately two decades and is reasonably well described functionally and mechanistically; and transfer RNA fragments (tRFs), a second class of regulatory smRNA that has only recently been rediscovered and is significantly less well described regarding its functionality.

1.3.1 TRANSCRIPTION FACTORS

Transcription factors (TFs) were among the first intracellular regulatory mechanisms to be discovered (the earliest article referencing the term »transcription factor« in its title on PubMed was published in 1972). TFs commonly translocate from the cytosol into the nucleus upon activation (often by phosphorylation), where they bind specific DNA sequences that usually range in size from 6 to 12 nt. The regions containing these binding sites (about 100 - 1000 nt in size) determine the effect upon binding, which can be one of two main modes: either a promoter, leading to an increased activity of transcription in the downstream vicinity of the binding site, or a repressor, having the opposite effect.

There exists a vast body of knowledge on TF-interactions with genes, mostly due to the long period of time since their discovery and the multitude of scientific publications, most often studying single TFs and their interactions with few genes, but cumulatively curated by several organisations. One of the currently largest curations of TF data, TRANSFAC, saw its original release in 1988. While these curation efforts can be extensive, they may present with serious bias towards particular TFs that might hold more scientific interest and thus are published far more frequently than others. Recently, comprehensive efforts have extended the available data significantly. Driven by the advent of RNA-seq, computational approaches have become able to not only comprehensively predict TF-gene interactions, but to do so in a highly tissue-specific manner (see Section 2.2.3). The human body is estimated to express up to 2600 distinct DNA-binding proteins, most of them presumed TFs,⁷⁶ although other studies give lower estimates.

1.3.2 MICRORNAs

THE FIRST ENDOGENOUS »SMALL RNA WITH ANTISENSE COMPLEMENTARITY« was described in 1993,⁷⁷ but microRNAs (miRNAs) were only recognised as a distinct regulatory class of molecules in the early 2000s. They are typically between 18 and 22 nt-long, single stranded RNA fragments, and their function is now largely undisputed: miRNAs serve as targeting molecules for a protein complex whose primary purpose is to repress translation of mRNA, and, in some cases, lead to mRNA degradation. The complex, therefore, is called RNA-induced silencing complex (RISC); central to its function is the family of argonaute (Ago) proteins, which can bind the mature miRNA and orient it for interaction with its targets. Guidance of RISC to the target mRNA is generally mediated via sequence complementarity between miRNA and the targeted mRNA. Specifically, a »seed« region, usually bases 2-8 on the miRNA, is mainly responsible for the interaction; in case of perfect complementarity of this seed to the mRNA sequence, the interaction is considered »canonical«.

In early miRNA research, the 3' untranslated region (UTR) of the mRNA was believed to contain most miRNA binding sites due to its greater accessibility (i.e., the lack of active ribosomes); however, cumulative recent reports suggest that binding inside the coding region of the mRNA is a regular occurrence(cite). The rules governing miRNA binding to target sequences show considerable flexibility; a recent study shows about 30% of analysed relationships to be of »non-canonical« nature(cite). In those cases, seed pairing with the mRNA is often imperfect. To ameliorate this loss of stability, compensation occurs typically by a secondary complementary structure after a small gap of non-complementary bases, leading to a »bridge«-type constellation. This flexibility has implications in applications involving targeting algorithms; those that consider only the seed region are more prone to false negatives than models that consider, for instance, the free energy of the whole molecule (see Section 2.2.4). figure?

BIOGENESIS

miRNAs, similar to coding genes, are transcribed from loci on the genome, many inside introns or even exons of coding genes.⁷⁸ The primary transcript (primary miRNA or pri-miRNA) typically contains a hairpin-like structure that usually results in a double-stranded molecule because of internal complementarity, and can contain up to six mature miRNAs. This hairpin structure is recognised by the DGCR8 protein (DiGeorge Syndrome Critical Region 8, in invertebrates called »Pasha«); the complex then associates with the RNA-cleaving protein »Drosha«, which removes bases on the opposite side of the hairpin, creating a miRNA precursor (or pre-miRNA), which is subsequently exported from the nucleus by the shuttle protein Exportin-5. In a final step in the cytosol, the ribonuclease »Dicer« removes the loop joining the 3' and 5' arms of the pre-miRNA, resulting in a duplex of mature miRNA, about 20 nt long. Initially, it was thought to contain only one active

miRNA, resulting in a designation of »miRNA*« for the complementary strand (commonly, the strand with lower expression). However, this notion has been disproven, and to reflect the possibility of both strands performing miRNA functions, nomenclature has changed to specify the arm of the pre-miRNA from which the mature form originates (suffix »-3p« for the 3' arm, and »-5p« for the 5' arm).

miRNA genes, in the same way as protein coding genes, can be subject to promoters and repressors, adding another layer of expression control by TFs. However, these TF-miRNA relationships are far less well described than common coding gene interactions, because miRNAs due to their shortness are not amenable to many standard gene expression assay forms. Estimation of the number of distinct gene targets of any one miRNA varies widely; however, it is generally accepted to not be less than several dozen targets per miRNA, and up to thousands of genes per miRNA (although that estimate might be overenthusiastic).

ORGANISATION AND CURATION

miRNAs are organised and curated by means of a periodically updated web-based platform, miRBase.⁷⁹ For *Homo sapiens*, miRBase v21 contains 2588 mature miRNAs from 1881 precursors. Evolutionarily, the miRNA repertoire has grown from rodents to primates, resulting in a number of primate-specific miRNAs that may convey additional function. miRNA nomenclature is organised⁸⁰ in a way that assigns evolutionarily conserved miRNAs the same designation (number) in all species in which they are expressed. In their full names, a prefix stating the organism of origin is added; for example, hsa-miR-125b-5p (for *Homo sapiens*) and mmu-miR-125b-5p (for *Mus musculus*) share the same sequence and most of their functionalities.

miRNAs are subcategorised in families (designated »mir« with lowercase »r«) by their genomic origin and phylogenetic homology aspects. As the annotation itself, family affiliations are in flux and change with each miRBase version. miRBase v21 lists 151 distinct miRNA families with 721 individual members in total. The remaining 1867 miRNAs do not (yet) belong to a larger family; the majority (80%) of those is newly discovered, as indicated by a 4-digit designation number.

DISEASE ASSOCIATION

miRNAs have been associated with a number of CNS diseases, including AD, Parkinson's Disease (PD), BD, and SCZ. However, the largest contribution since their discovery by far has been made by cancer research; of the approximately 90 000 publications found on PubMed with the term miRNA, about 42 000 involve cancer (search term »miRNA AND cancer«). In comparison, »miRNA AND Alzheimer's Disease« results in about 600 hits, while a search for »miRNA AND Schizophrenia« yields just 363 publications (as of October 2019).

In AD, several groups of miRNAs have been found to show characteristic perturbations before the

onset of symptoms, which makes them interesting biomarker candidates.⁸¹ Some miRNAs have been extensively studied in a variety of contexts, most prominently hsa-miR-132-3p. Among its targets are several key neuronal regulators (e.g. FOXP2, FOXO3, P300, MeCP2), and it is in turn controlled by many pivotal neuronal elements (e.g. REST, ERK1/2, CREB); this presents an explanation for the many physiological and pathological situations that miR-132-3p has been found to play a role in. Its functions include the control of neuronal survival/apoptosis, migration and neurite extension, neuronal differentiation, and synaptic plasticity.

miRNAs are able to fulfil their regulatory purpose in a context- and cell-type-dependent manner,⁸² such that the perturbation of one single miRNA might provide different functional outcomes in different tissues (e.g., glial cells and neurons), or different stages of disease. However, this »jack-of-all-trades« behaviour also poses significant problems in establishing miRNAs as pharmacological targets: In the case of antagonising of mimicking an existing miRNA, the amount of off-target effects would not only be enormous, the entire definition of an off-target effect would continuously change between tissues and during the course of the disease. For this reason, the design of custom oligonucleotides with limited capabilities might be preferable in the development of therapeutics based on RNA interference (See also Section 5.3).

1.3.3 TRANSFER RNA FRAGMENTS

TRANSFER RNA (tRNA) BREAKDOWN PRODUCTS have been known for decades, with first descriptions in the 1970s; back then, they were associated with a higher turnover of tRNA in cancer cells,⁸³ and proposed as urine-based biomarkers for certain malignancies.⁸⁴ However, their genesis was attributed to random processes, and due to lacking molecular biology characterisation techniques, interest in those fragments quickly faded. It was not until recently that studies have shown tRNA to be a major source of stable expression of small noncoding RNA^{85,86} in most mammalian tissues. Indeed, replicating the reports from the 1970s, tRNA breakdown products are the dominant form of small RNA in secreted fluids, such as urine and bile, and make up large parts of other bodily fluids as well.⁸⁷ They exist in two major forms: transfer RNA halves (tiRNAs), and the smaller transfer RNA fragments (tRFs). tiRNAs derive from either end of the tRNA, and are created by angiogenin cleavage at the anticodon loop.^{88,89} Smaller fragments are derived from the 3' and 5' ends of the tRNA (3'-tRF/5'-tRF) or internal tRNA parts (i-tRF), respectively, and may incorporate into Ago protein complexes and act like miRNAs to suppress their targets.^{90,91}

However, there is considerable controversy about the generalisation of tRF functions, as distinct publications discover very different and sometimes opposing mechanisms of action for their respective fragments. An obvious assumption is the miRNA-like functionality, at least for those tRFs that are in the length range of miRNAs. There have been several instances of tRFs proven to act as miRNA-like suppressors of translation in a RISC-associated manner,⁹¹ and of Dicer playing a large

part in their biogenesis.⁸⁵ There are even instances of small RNA molecules previously mislabeled miRNAs that have been discovered to actually be tRNA-derived, such as miR-1280.⁹²

On the other hand, multiple groups have identified tRFs to function not in an antisense-complementary manner, but by homology aspects. A valine-derived tRF was found to regulate translation by competing with mRNA directly at the binding site at the initiation complex and thereby displacing the original mRNA, leading to its translational repression.⁹³ Others have found multiple classes of tRFs derived from glutamine, aspartate, glycine, and tyrosine tRNAs, that displace multiple oncogenic transcripts from an RNA-binding protein (YBX1), conveying tumour-suppressive activity.⁹⁴ Most counterintuitive is the recent finding of a tRF proven to bind to several ribosomal protein mRNAs and enhancing their translation, and, when specifically inhibited, leading to apoptosis in rapidly dividing cells.⁹⁵

There is no consistent nomenclature yet to describe and organise tRFs, which are by nature more heterogeneous than miRNAs; while only 61 mature tRNAs are required in a cell to achieve a one-to-one »codon→amino acid« translation, one tRNA molecule can be the origin of several hundred distinct tRF molecules. Additionally, the amount of human tRNA genes is estimated at 500-600,⁹⁶ and there are many more pseudo-tRNA genes. To communicate the identity of individual tRFs, multiple approaches are common in current literature; most prominently, tRFs are tied to the parent tRNA and the amino acid carried by this tRNA. To illustrate: The 22-nt LeuCAG3' tRF (meaning: a fragment of 22 bases starting at the 3' end of the leucine-carrying tRNA with anticodon »CAG«) was shown to play an important role in regulating ribosome biogenesis.⁹⁵ Since there is no repository of the likes of miRBase yet, this approach can be cumbersome for replication purposes, and explicit statement of the exact sequence of each fragment is a must in publication. In fact, since the aforementioned paper does not mention the sequence explicitly, there exist 6 distinct possibilities of fragments fitting this description. While manageable on this small scale, this system prohibits efficient analysis of larger sets of tRFs that cannot be individually controlled. For this reason, the approach of Loher and colleagues⁹⁷ might be preferable: they propose the generation of a "license plate" based on the sequence of the fragment directly, composed of the prefix »tRF«, the length of the fragment, and a custom oligonucleotide string encoding (e.g., »B3« codes for »AAAGT«). This way, tRF names are unique and unmistakably linked to the sequence, nomenclature is species-independent, and tRNA

Disease?

origin can be quickly determined by sequence lookup.

1.4 NESTED MULTIMODAL TRANSCRIPTIONAL INTERACTIONS

- THE NEED FOR CONNECTOMICS

The ultimate aim of transcriptional connectomics is the combination of all interacting cellular components in a model that satisfactorily explains our real-life observations and is able to predict the functional outcome of a modification of one of these players. Even in the simplified case of only

studying the interactions between coding genes, TFs, miRNAs, and tRFs, the complexity of the required model exceeds our current capabilities by far. The more we know about the functioning of these intertwined systems, the more we understand how much there is still to learn.

For example, only recently has it become clear how complex transcriptional regulation by means of TFs really is, and, incidentally, the two systems studied foremost in this dissertation (nerve and immune cells) are the two most transcriptionally complex systems in any mammal. Through study of comprehensive genomic information of 394 tissue types in approximately 1000 human primary cell, tissue, and culture samples (from the FANTOM5 consortium) it was estimated that the mean number of active TFs towards any given gene is highest in immune (12 TFs per gene) and nervous cells (10 TFs per gene), and that any one TF in nervous and immune cells controls expression of a mean of 175 and 160 genes, respectively⁹⁸ (see also Section 2.2.3).

Similarly, it has been found that miRNAs, particularly in the nervous system, possess a much higher tissue specificity than coding genes, resulting in an expression landscape that varies widely between individual neuron types that are in close proximity in the brain. With the exception of single cell RNA-seq, no modern analysis method is capable of a resolution appropriate for accurate characterisation of these expression patterns, resulting in extinction of the signal of miRNAs that are not expressed consistently across cell types (similar to »housekeeping« genes) because of statistical interference. Very recent studies show that miRNA-gene co-expression networks are tightly linked to cell types in the nervous system, and that groups of miRs as functional modules associate with particular phenotypes in developmental and mature states.⁹⁹ This functional association with cell phenotype was found in quality comparable to the expression patterns of TFs, yet in quantity conveys smaller impact and thus is thought to be a fine-tuning mechanism, subtle and precise in purpose.

Another aspect of the tissue specificity of CNS-associated miRNAs is the high likelihood of under-representation or even non-discovery of those very specifically expressed miRNAs. Adding to the problem is the experimental bias towards rodent models when it comes to thorough studies of the CNS, where human or other primate samples are a rarity compared to rats or mice. Assessments of the numbers of yet unknown novel primate- and tissue specific miRNAs estimate their magnitude in the thousands,¹⁰⁰ resulting in an effective doubling of currently known miRNAs.

These high numbers of potentially interacting players present computational challenges: If estimating the number of expressed genes in a human cell at 20 000 (and the number of TFs at a low 1000), this makes for an estimated minimum of 200 000 »real« interactions in the possible $C = \frac{1000!}{10!(1000-10)!} \cdot 20\,000$, which practically equals infinity; this is without accounting for different tissue types or cell states (e.g., differentiation or disease). Similarly, the amount of mature miRNAs (2588 in miRBase v21) and their ability to target even more distinct transcripts than TFs with one single molecule present immense computational requirements for even listing all possible or actual relationships. An interaction table describing targeting of genes by miRNAs in one type of tissue has $2588 \cdot 20\,000 \approx 50\,000\,000$ individual fields.

Combining the different modes of transcriptional interaction presents additional challenges. A simple model system to visualise (in only one type of cell) the interaction of TFs targeting genes, and of miRNAs targeting genes as well as TFs, contains about 20 000 genes (a subset of which of the size of about 2000 are TFs), 2588 mature miRNAs, and a total of $2588 \cdot 20\,000 + 2000 \cdot 20\,000 \approx 90\,000\,000$ potential interactions. In standard application scenarios, such as the generation of an interaction network around a group of genes (e.g., the cholinergic genes), the processing requirements grow linearly with each added interaction partner, and exponentially with every regulatory layer that is added.

Practically, this information has to be provided, gathered, and integrated, which further multiplies the amount of storage and processing power required. miRWalk 2.0, a collection of miRNA interaction data, has collected 12 of the most popular miRNA-targeting prediction datasets, each of which has their strengths and weaknesses (see 2.2.4). Experimentally validated interactions (e.g. as collected in DIANA TarBase or miRTarBase) are gold standard, but far from comprehensive and strictly speaking only relevant for the cellular context in which the experiment was originally performed; there are also different evidence qualities to be accounted for, depending on the type of experiment performed. Ideally, all of these data are still accessible when performing the analysis, so a database created for this purpose should be able to incorporate all this information without any data loss while still remaining feasible in terms of computation time as well as space and working memory requirements.

This dissertation will first describe the creation of such a database and what has been learned during its various stages, and then go on to apply the database to different biological problems from real world experiments, such as the cholinergic differentiation of human male and female cultured neuronal cells, and the blood of stroke victims.

example of
standard inter-
action gene x
miR x TF?

more exten-
sive descrip-
tion of con-
tent?

»Wir sehen in der Natur nie etwas als Einzelheit, sondern wir sehen alles in Verbindung mit etwas anderem, das vor ihm, neben ihm, hinter ihm, unter ihm und über ihm sich befindet.«

Johann Wolfgang von Goethe

2

miRNetDB: Creation of a Comprehensive Connectomics Database

Natural philosophy, as represented by the thought of Johann Wolfgang von Goethe, is concerned with the holistic description of nature and the explanation and interpretation of its particular mechanisms. Although natural philosophy is the predecessor of modern, empirical science, its concepts and approaches are still valuable in today's data driven world; as the data we collect grows to dimensions that can only be interpreted computationally, functional reductionism becomes all the more important: By studying the facets of nature, we strive to understand it as a whole. Similarly, we regularly encounter Goethe's paraphrase of »all things are connected« in neuro-immunology, and in transcriptional connectomics.

BIOINFORMATIC SUPPORT IN CONNECTOMICS is indispensable, which can be seen by the sheer multitude of possible interactions between the participating factors. However, when I began working on this project (October 2015), there was no integrative database available for this purpose. Earlier that year, miRWalk 2.0 had been published, for the first time providing a relatively comprehensive source of predicted as well as experimentally validated miRNA targeting data¹⁰¹ (see 1.3.2). One year later, Marbach's »regulatory circuits« were published,⁹⁸ enabling analysis of comprehensive TF→gene relationships in 394 human tissues (see Section 1.3.1). These collections (as well as the data they were derived from) are the basis of the database further called *miRNetDB*, the development of which will be described in the following chapter.

Since a large part of the scientific progress of this dissertation deals with practical problems of mul-

timodal connectomics, I will begin by describing the infrastructure that makes effective computation of these problems possible. After this technical description of database structure and creation, I will explain the types and organisation of its content. The remainder of the chapter will then deal with the application of this infrastructure to real-world problems in transcriptional connectomics, and the statistical approaches suited to this special case.

2.1 IMPLEMENTATION

For any biological question to be asked in a bioinformatics setting, the effectiveness of the computational query determines the practicality of the approach. Because resources (i.e., processing power, storage, and working memory) are limited, the database that is queried should be organised in a way that facilitates retrieval of the desired information without excess processing of useless information. In the simplified case of only miRNAs interacting with genes in one direction (miRNA→gene), this means retrieval of only those interactions relevant for the queried genes or miRNAs.

Traditional table-based approaches (also known as relational databases) such as SQL (»Structured Query Language«) cannot provide such an implementation, since individual entries for genes and miRNAs (rows and columns) have to be accessed in their entirety, whether there is a connection between gene and miRNA (1) or not (0). Additionally, adding layers to these interactions (e.g., distinct prediction algorithms, tissues, or the interaction between TFs and genes) require the addition of entire tables the same size as the database, which is detrimental to effective use of space; and more complex queries also necessitate the transfer of information between those distinct tables (in SQL typically via a `JOIN` command), which claims additional working memory and processing time. Overall, the so-called »many-to-many« organisation of data does not lend itself to representation in a relational database.

The actual performance is determined by the processing power of the machine it is running on and several structural properties, such as organisation, indexing, monotony, and of course the size of the database; therefore, an estimation of processing time for queries is bound to be inaccurate. However, processing times typically do not vary on the scale of orders of magnitude, and thus general estimations can be made. Well optimised SQL databases with a size of 5 to 10 GB on disk usually require tens of minutes if not hours to complete one single complex query;¹⁰² *miRNetDB* in its current form takes up approximately 15 GB of storage. Since one analysis typically consists of several hundreds (and, in the case of permutation analyses, several hundreds of thousands) of these queries, processing times in SQL implementation are too long to be practically useful. (It seems important to note that, as of 2018, SQL also offers a graph-based organisation in addition to the traditional, relational layout. These two are separate systems, and not to be confused. The advantages of Neo4j as explained in the following should be seen from the perspective of 2015, when the database was established, and when there was no graph-based SQL implementation.)

Figure to explain tables?

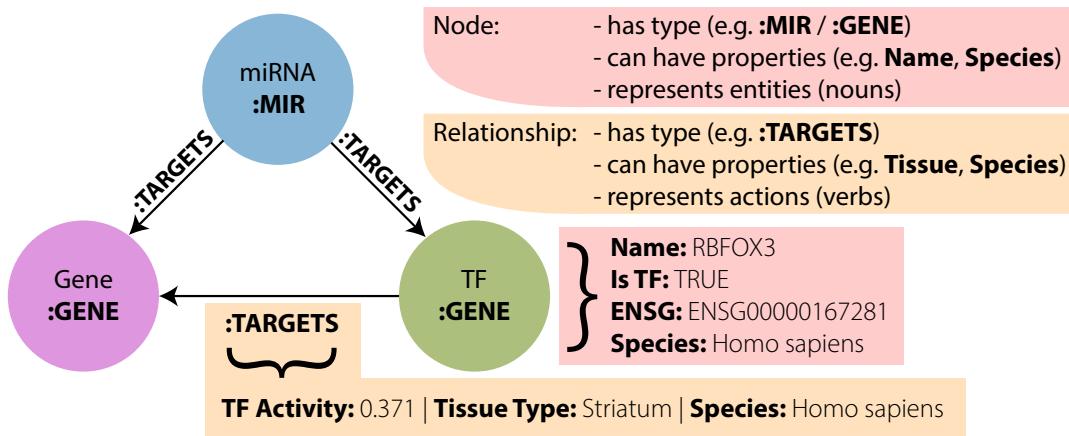


Figure 2.1: Organisation of a graph database. A graph consists of two basic building blocks: **Nodes**, representing entities, and **edges**, representing connections between entities. Each database entry (node or edge) is an instance of a particular *type* and can possess an arbitrary amount of *properties* detailing its specifics.

2.1.1 NEO4J: A GRAPH-BASED INFRASTRUCTURE

To query and display biological data that are organised in a network-like structure (many-to-many), a database that lends itself to the efficient processing and storage of network data is optimal. »Neo4j« utilises a database structure that is built on the save and recall of data points in *nodes* and *edges*, which represent entities (nodes) and relationships between those entities (edges); both nodes and edges can have any number of attributes and a unique property called »type«, usually describing the class of the entry (such as *gene* or *miRNA*). This database organisation replicates the network-like structure of the biological data studied (Fig. 2.1). Neo4j combines the network-like data structure with an efficient indexing system for quickly finding the entries queried for, and then »walks« along the edges of the nodes that have been found, thus only searching and returning the data that is relevant to the current query. Theoretically, this makes the database more likely to be efficient in the setting of transcriptional interactions, an estimation that turned out to be true.

Depending on the input, these queries can also be rather large; however, the main pitfall of tabular databases such as SQL is circumvented: there is no need to process entire rows or columns of the table to make sure that the query is satisfied in its entirety. This is particularly useful in a setting of sparse information. To illustrate: Only 30 of the 2588 miRNAs target a specific gene, which is common; a relational database, after finding the index of the queried gene, would have to search 2588 fields for 1/0. The graph database, on the other hand, has to execute only 30 searches (or, more accurately, 30 »walks« along the edges connected to the indexed node). In practice, even in the very first prototype implementations, this accelerated standard-case computations immensely, and was even able to accommodate advanced approaches in situations that had been inaccessible in the tabular implementation.

2.1.2 HIGH-THROUGHPUT DATABASE GENERATION

Neo4j provides several API (»application programming interface«) possibilities in implementation. For the purpose of entering large amounts of data into the database at once, the Java implementation is superior to the other forms in that it provides a batch processing mode via its `BatchInserter` class. I thus wrote a custom Java program for the purpose of creating an initial state of the database from the largest set of data, the complete miRWalk 2.0 content with 12 algorithms and validated interactions. The downloaded data was organised in a plain text based file format, with one text file for each miRNA, totalling in size about 6 GB (for *H. sapiens*). The database was set up in a way that allows only one node for each individual miRNA and gene entered to avoid duplications, using the commands

- `createDeferredConstraint()`
- `assertPropertyIsUnique()`
- `createDeferredSchemaIndex()`

of the Neo4j Java package. This approach made sure to create only one node for each miRNA (type: MIR) and gene (type: GENE) in the data, which is essential for proper functioning of the database. Each of these nodes received several properties to store individual data, such as the various gene/miRNA identifiers, miRNA sequence, and species.

Between those basic nodes, the batch insertion process created edges for each relationship that was found in the original data, assigning a type identifier to each edge detailing the origin of this interaction (type: name of the prediction algorithm or »VALIDATED« for experimental data). Thus, while the nodes for genes and miRNAs themselves are unique, an arbitrary number of relationships can exist between any two nodes, depending on how many interactions they share.

2.1.3 MAINTENANCE AND QUALITY CONTROL

All additional datasets, such as the TF regulatory circuits or tRF targeting predictions, were entered into *miRNetDB* using the regular operation mode. Testing was also performed in regular operation, with manual as well as automated tests to assert the correct transfer of information from raw data to the graph database, and to avoid unpredictable behaviour. At times, conflicts had to be resolved manually, for instance when miRNA names conflicted between old »miRNA*« and new »3p/5p« notation; all manual edits are documented in the code, which was published alongside Lobentanzer *et al.*⁶

Except for the rapid import of large amounts of data in creation of a database, the Java implementation of Neo4j does not offer many advantages over the native R implementation, »RNeo4j«. Thus, after creation and a short period of experimentation with graphical user interfaces, I abandoned the Java program in favour of the more flexible R programming. While Java is an object-based programming language, whose benefits lie in extreme flexibility in regards to platform and purpose, high

modularity, and speedy processing, R as a procedural language is the work horse of modern bioinformatics. Its procedural design (the division of data and functions that operate on that data) facilitates the transfer of approaches between distinct datasets, and the enormous vibrant community of data scientists using R provides a wealth of third party packages to tackle almost any bioinformatic task. In the remainder of this dissertation, all analyses are performed in R, unless specifically stated otherwise.

2.2 MATERIALS

All materials used in the creation of *miRNetDB* have been acquired from resources that are non-commercial, web-available, and open-source (in the case of code). All properties and relationships derived from this data were entered into *miRNetDB* as either nodes, properties of nodes, edges, or properties of edges.

2.2.1 GENE ANNOTATION

Even though »regular« protein coding genes have been known for a long time, there is no consensus yet about their nomenclature and organisation. Complicated by newly discovered functions and properties of phylogenetic nature, the scientific representation of the human genome is in constant flux. Several large organisations strive to provide a robust annotation of the human gene catalog, but also in many cases contradict one another. There are three nomenclature systems that are of high importance in modern genomics:

- The traditional naming system of acronyms (e.g. CHAT) and fantasy-names (such as »Sonic Hedgehog«), also occasionally called »gene symbol«, is still widely popular because of its accessibility to humans, but is also not particularly robust because of a high amount of synonyms with high confusion potential (see e.g. Section 1.2.7 on CDF) and instances of genes without names having to carry unwieldy systematic names.
- The American National Center for Biotechnology Information (NCBI), a branch of the National Institute of Health (NIH), curates and hosts a multitude of biological and medical data, and for the organisation of gene information uses its own systematic nomenclature termed »Entrez« ID. Entrez is a molecular biology database that integrates many aspects of biology and medicine in a gene-centred manner, and therefore Entrez IDs are useful to quickly connect a gene to its function, nucleotide sequence, or associated diseases. Entrez IDs are regular integers without additional characters.
- Akin to the NCBI effort, ENSEMBL is a project of the European Bioinformatics Institute (EBI) as part of the European Molecular Biology Laboratory (EMBL). Compared to the Entrez database, it is more focused on study and maintenance of the genome itself, and therefore has a more intricate nomenclature that allows for differentiation of, for example, genes and their various tran-

script isoforms (ENSEMBL IDs carry character prefixes for class identification, e.g., ENSG for genes, ENST for transcripts).

All of these are being used on a regular basis in many publications, and, often, they are used exclusively. As a result, the end user of the published data has to have access to all possible annotation forms, or, at least, a means to translate one into the other; often, this also introduces conflicts. For this reason, all ID types were entered into *miRNetDB* upon creation or during maintenance, for convenience and to minimise analysis prolongations due to conflict resolution.

2.2.2 MICRORNA ANNOTATION

miRBase provides a consistent annotation for miRNAs. Due to their relatively recent discovery, there still are major changes from version to version; the syntax, however, is stable. In addition to the miRNA »names« that are composed of species, the string »miR«, pre-miRNA designation number, and strand origin (not in all cases!), such as »hsa-miR-125b-5p«, miRBase provides IDs for pre-miRNA molecules (also called ancestors) termed »MIID«, and IDs for mature miRNA molecules termed »MIMAT«. However, in practice, these are rarely used. Similarly, miRNA families are annotated using the »MIPF« ID.

2.2.3 TRANSCRIPTION FACTOR TARGETING

The FANTOM5 project has applied 5' cap analysis of gene expression (CAGE) to a large number of human samples from diverse tissues to determine the accurate 5' ends of each transcript.¹⁰³ Knowledge of this fact enables accurate prediction of promoters likely to control a transcript's expression. Marbach and colleagues used this information in combination with detailed human gene expression data to derive a complex interaction network of TFs and genes (»regulatory circuits«), and in doing so aggregated samples with similar expression patterns and origins into 394 fictional tissues.⁹⁸ For every tissue, each TF was assigned transcriptional activities towards all genes that it supposedly targets (with the sum of all activities in any given tissue being 1). Marbach and colleagues have shown that the cumulative transcriptional activities towards any given gene correlate well with the actual gene expression in corresponding samples from an independent repository.

Even in its fifth iteration, FANTOM data is not entirely comprehensive, which came to my attention due to a cholinergic anomaly: The 5' CAGE peaks of the *CHAT* and *CHRNA7* (the nicotinic $\alpha 7$ receptor subunit) genes in raw FANTOM5 data do not pass the expression threshold, and therefore are not included in, e.g., Marbach's »regulatory circuits«. Both are critically important not only for neuronal cholinergic systems, but also for the non-neuronal aspect of immune processes. For instance, macrophages have been shown to produce ACh via ChAT, and the $\alpha 7$ homomeric ACh receptor conveys direct immune suppression by its expression on monocytes.³⁷ Paradoxically, the CAGE peak of *SLC18A3*, which lies in the first intron of *CHAT*, crosses the threshold and there-

fore is included in the data. Unfortunately, I was not able to remedy these circumstances even upon personal communication with Daniel Marbach (author of »regulatory circuits«) and Hideya Kawaiji of the FANTOM5 consortium, although the latter acknowledged the possibility of a gene annotation deficit leading to misattribution of the *CHAT* signal to *SLC18A3* due to the closeness of their 5' ends. Thus, it seems viable to substitute *SLC18A3* targeting data for the absent *CHAT* data in certain situations.

The entire collection of transcriptional activities in all tissues was downloaded from the project's web page,⁹⁸ and neuronal and immune tissues were manually curated and entered into *miRNetDB*. The collected data comprises 33 neuronal tissues and 26 immune cell tissues (Appendix A), and 1 130 196 TF→gene relationships in total (not all 394 tissues were entered).

2.2.4 MICRORNA INTERACTIONS

The content of miRWalk 2.0 is freely available online;¹⁰⁴ however, there is no option of downloading the complete set. The targeting data thus was downloaded per miRNA using a custom crawler, with standard options for all 12 prediction algorithms (miRWalk, miRDB, PITA, MicroT4, miRMap, RNA22, miRanda, miRNAMap, RNAhybrid, miRBridge, PICTAR2, and TargetScan) in plain text format. For experimentally validated interactions, the main sources were DIANA TarBase¹⁰⁵ and miRTarBase,¹⁰⁶ both of which offer complete download options. As of 2019, the 3.0 version of miRWalk allows complete species downloads; however, the developers have abandoned their third party algorithm plurality reducing the number of available alternatives from 12 to 4, which can be considered a significant disadvantage:

While sequence complementarity, particularly of the »seed«-region, is the primary paradigm of miRNA-mRNA interaction, prediction algorithms vary widely in their implementation, general purpose, and approach to interaction prediction (for a comprehensive review of approaches and rules, see Yue *et al.*¹⁰⁷). A large group of available options utilise sequence conservation aspects to increase candidate viability (such as miRanda, PicTar, TargetScan, and microT4). Others, such as RNA22 and PITA, utilise biophysical aspects such as free energy of binding or the accessibility of target sites due to secondary RNA structures as prediction arguments. All of these approaches have their up- and downsides, e.g. considering their general precision and sensitivity, or their adequate prediction of particular cases, such as multiple site targeting. Thus, it has been proposed to use a combination of complementary approaches instead of only one algorithm per analysis.¹⁰⁸ For this reason, I might have preferred the 2.0 version of miRWalk, even if 3.0 had been available at the time.

One advantage of the collection of all data in a quickly accessible database is the opportunity to compare the different approaches to target prediction. A statistical evaluation of the collected interaction data from miRWalk 2.0 showed vast differences in general prediction quantity (Table 2.1) as well as prediction accuracy and sensitivity when compared to the validated subset of data (Table 2.2). Since the ground truth is not known, this is an additional argument for the combination of multiple

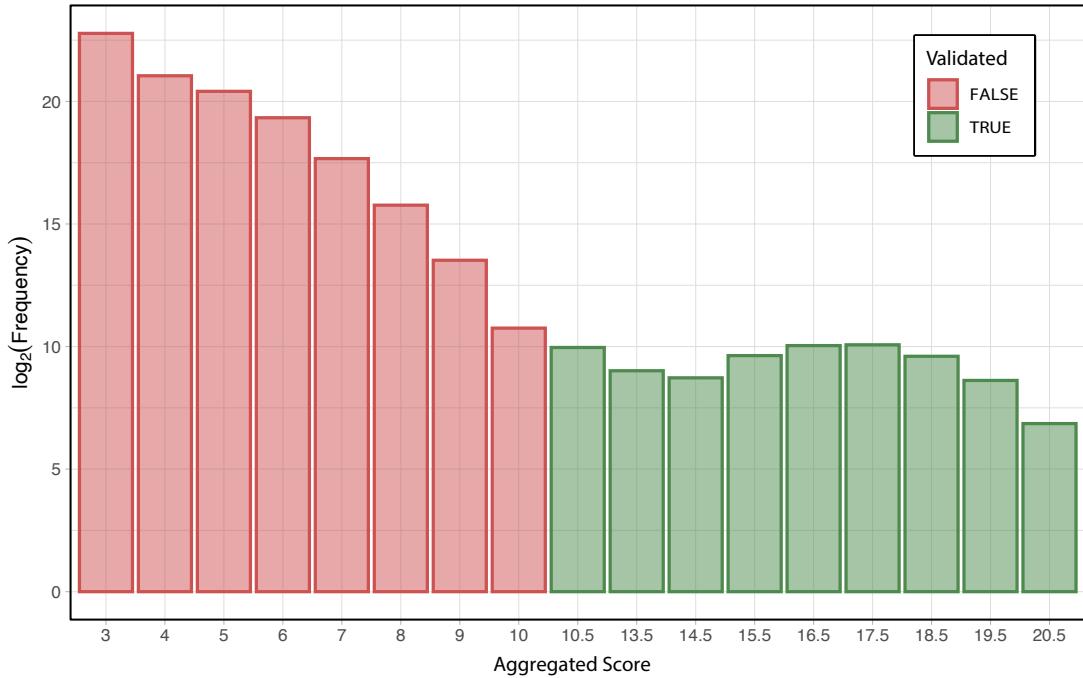


Figure 2.2: Histogram of miRNA → gene score distribution. Aggregation of individual algorithms yields a score range of 3 to 10 per individual miRNA → gene interaction. In case of additional existence of experimental validation (evidence level high) for any predicted interaction, score is increased by 10.5. The distribution shows a sharp decrease in predicted interactions towards higher scores, and a maximum of validated interactions at prediction scores 6 and 7.

algorithms instead of the use of a single set. Apart from RNAhybrid and miRBridge, all algorithms presented reasonable base hit frequencies and increases in the validated test set. While miRBridge already has the lowest positive frequency of all the algorithms, it is the only one to achieve a negative score in the validated test set. On the other hand, RNAhybrid has a vastly higher base hit frequency than the second highest scoring algorithm (by more than 300%), making it very likely to produce false positive results, and less valuable in the aggregation scoring system. The remaining 10 algorithms were included in *miRNetDB* targeting data. For ease of use, an additional relationship type was created from the aggregated single algorithm hits of any miRNA → gene relationship, with the sum of algorithms predicting the interaction as a score variable. This yields a theoretical score range from 3 to 10 (miRNA → gene relationships with only one or two hits were ignored for the sake of space). To account for experimentally validated interactions, each miRNA → gene relationship that was supported by strong evidence of interaction was modified by addition of 10.5 score points (a half point for quick identification of a validated relationship), extending the maximum score to 20.5 points. The resulting optimised graph contains 11 687 931 human miRNA → gene targeting relationships with a distinct score distribution (Fig. 2.2). In comparison, only 6146 miRNA → gene relationships are experimentally validated with »strong« evidence.

algorithm	hit frequency
RNAHYBRID	71.62%
MIRMAP	19.90%
MIRWALK	19.74%
TARGETSCAN	16.33%
RNA22	12.34%
MICROT4	11.81%
MIRANDA	10.65%
PITA	4.90%
MIRDB	1.17%
MIRNAMAP	0.75%
PICTAR2	0.62%
MIRBRIDGE	0.15%

Table 2.1: Prediction algorithms ordered by the fraction of all possible interactions they predict as being real (positive rate). Different algorithms display a wide variation of hit rates in the entirety of predicted interactions between any miRNA and gene. Red: excluded from analysis.

algorithm	validated hit frequency	hit rate increase
PICTAR2	6.98%	1129.40%
MIRDB	9.80%	838.43%
MIRANDA	51.73%	485.94%
TARGETSCAN	70.63%	432.51%
MIRNAMAP	3.10%	410.95%
PITA	15.57%	317.20%
MICROT4	32.60%	276.10%
MIRMAP	53.86%	270.65%
MIRWALK	50.95%	258.15%
RNA22	22.51%	182.38%
RNAHYBRID	90.47%	126.32%
MIRBRIDGE	0.01%	0.00%

Table 2.2: Prediction algorithms ordered by their increase in true positive rate when considering only validated interactions. The hit rate increase when comparing experimentally validated interactions with the entire predicted data (Table 2.1) is also subject to strong variation. Hit rate increase is the increase of hit rate if only considering validated data as opposed to all predicted interactions. None of the studied algorithms unite a good precision (hit rate increase) and coverage (validated hit frequency).

2.2.5 FILTERING OF AGGREGATED PREDICTION SCORES

For the estimation of the »true« miRNA→gene interactions in the predicted-only data in *miRNetDB*, two premises are relevant: First, the enormous amount of hits with a score of 3 in all likelihood is an over-estimation, and second, the amount of currently validated interactions can be but a small fraction of »true« interactions. Assuming the truth lies on the axis between these two extremes (i.e., at some score value inside the *miRNetDB* interactions), the true amount of human miRNA→gene interactions must approximately fall within the range of 2^{10} to 2^{20} . Looking at the score distribution of all *miRNetDB* interactions (Fig. 2.2), the maximum amount of validated interactions is predicted by a combination of 6 or 7 algorithms (i.e., a score of 16.5 or 17.5). Thus, to approximate the true state, I chose to apply a low-cut filter to *miRNetDB* queries at a minimum score of 6. This is the standard case referred to as »*miRNetDB* query« in the remainder of this dissertation. In some cases, such as the graphical analysis of whole-genome miRNA targeting (see e.g. Section 3.8), the score threshold was raised to 7 to circumvent computational limitations.

2.2.6 DE-NOVO PREDICTION OF TRF TARGETING

Due to the recency of their (re-)discovery, no comprehensive interaction sources exist for transfer RNA fragments. There have been documented cases of miRNA-like behaviours of distinct RNA fragments,^{85,91} justifying an attempt to predict interactions in a comprehensive manner. Of the available options for nucleotide interaction prediction algorithms, TargetScan¹⁰⁹ seems particularly suited for this task because it provides the option of evaluating the evolutionary conservation of target

sites in the putatively targeted genes, thereby providing an additional layer of security: The sequence of 3' UTRs is evolutionarily less stable than the coding part of genes; thus, high conservation of the binding site might indicate evolutionary pressure to keep up the interaction with the fragment, making an actual function of the interaction more likely. TargetScan also presents with reasonable sensitivity and specificity as confirmed by an independent group,¹¹⁰ and through an additional algorithm allows the attribution of a score based on the branch length (on the species tree) of conserved targeting.¹¹¹

miRNA-like behaviour implies the existence of a region on the tRF similar to a miRNA »seed«, and TargetScan also expects a seed as input to its targeting algorithm. Since there has been no definitive answer to the question as to where the seed region in tRFs might be, it is safest to assume nothing and explore all possibilities, i.e., simulate every possible seed position for interaction discovery. For this purpose, all discovered sequences of tRFs (exceeding a base mean expression of 10 counts) were chopped into 7-nt pieces (7mers), which is the length of miRNA seeds, and statistically improbable enough to appear in the genome at random; the average length of a human 3' UTR is 800 nt, so the probability of finding any 7mer randomly in any one 3' UTR is $p = \frac{800}{4^7} = 0.049$, which agrees with the 5% false discovery ratio (FDR) convention.

Describe TargetScan process

2.2.7 MICRORNA PRIMATE SPECIFICITY

During the course of evolution, higher organisms typically attained more complexity in a variety of functional categories. The CNS as the system of highest complexity underwent several drastic developments from invertebrates to lower mammals to higher mammals still. miRNAs are no exception. While many miRNAs are functionally as well as literally conserved in all mammals, primates in particular have gained a substantial amount of novel miRNAs whose function is in large parts elusive. Due to the restrictions on experimentation on higher mammals, particularly primates, many of those miRNAs can only be studied observationally, or by transgenic experiments in rodents. A cholinergic example of a gain-of-function in higher mammalian miRNA regulation is the vesicular acetylcholine transporter, SLC18A3. As described in Section 2.2.3, the SLC18A3 gene is situated in the first intron of CHAT, and thus is always primarily co-expressed with the latter. However, a primate-specific miRNA, miR-298, targets the 3' UTR of SLC18A3.¹¹² Thus, the primate neuron has gained a mechanism of independent SLC18A3/CHAT regulation that the mouse, for example, does not possess. It is easily imagined that such a gain of neuronal flexibility, in many instances, can aid the development of a more effective brain. However, the primate specificity of miRNAs is not yet consensus, and thus not found in annotation databases such as miRBase, even though they list all miRNAs discovered in any species. To get an impression of the amount of possible gain of function, I performed a review of miRNAs expressed in a representative variety of annotated species. From hereon out, largely method-related paragraphs will be set in sans-serif font face.

Species Selection

The tested species were selected from miRBase v21. Some of the available species are severely limited in the extent of miRNA annotation, likely because of a research bias. Therefore, only the most well-annotated species were selected. These are (number of annotated primary and mature miRNAs in brackets):

- *Homo sapiens* (human; 1881, 2588)
- *Gorilla gorilla* (gorilla; 352, 357)
- *Pan troglodytes* (chimp; 655, 587)
- *Pongo pygmaeus* (orangutan; 642, 660)
- *Macaca mulatta* (rhesus macaque; 619, 914)
- *Bos taurus* (cow; 808, 793)
- *Canis familiaris* (dog; 502, 435)
- *Mus musculus* (mouse; 1193, 1915)
- *Rattus norvegicus* (rat; 495, 765)

The first four species belong to the hominid group; the first five are primates. It is likely that these collections are not complete, with the degree of completeness depending on the amount of research performed on the species (as demonstrated, e.g., by the difference between mouse and the other non-primates). This places considerable difficulty on asserting primate specificity of miRNA, and in turn on assertion of the effects of evolution on the miRNA regulatory system.

Single miRNA Inter-Species Homology Computation

To determine the homology of miRNAs between the studied species, reference genomes were downloaded from the respective sources and analysed phylogenetically, using the genomic coordinates provided by miRBase. Sequence homology was determined via dynamic programming using the Smith-Waterman algorithm.¹¹³ Briefly, this algorithm can be used to determine the similarity of two genomic sequences, based on a scoring system rewarding matches and penalising mismatches. Smith and Waterman extended the original approach by Needleman and Wunsch,¹¹⁴ which is used to compare two complete sequences. Both algorithms rate an alignment by dynamic scoring inside a 2D-matrix, with the sequences to be compared as the x- and y-axes (one letter per cell). By a change in the scoring system, the Smith-Waterman algorithm finds the best local alignments, instead of comparing the two sequences in their entirety. In the case of miRNAs, this behaviour is useful because, between species, there are frequent additions or deletions of single nt on both ends of the homologous miRNA. Genomes were procured from the following sources:

- *Homo sapiens*: GRCh38 (NCBI)
- *Gorilla gorilla*: gorGor3 (UCSC)
- *Pan troglodytes*: panTro4 (UCSC)
- *Pongo pygmaeus*: PPyG2 (Ensembl)
- *Macaca mulatta*: rheMac3 (UCSC)
- *Bos taurus*: bosTau6 (UCSC)
- *Canis familiaris*: canFam3 (UCSC)
- *Mus musculus*: mm10 (UCSC)
- *Rattus norvegicus*: rn5 (UCSC)

Using the genome coordinates provided by miRBase, the genomic sequences of miRNAs and pre-miRNAs of each species were determined. Using the Smith-Waterman algorithm, all identified homologs of human miRNAs were subjected to homology scoring, and score results were visualised as a heatmap.

INTER-SPECIES DISTRIBUTION OF miRNAs

The inter-species relationships of annotated miRNAs do not follow a simple evolutionary distribution from less complex to more complex organisms, but rather seem to partially result from parallel development (Fig. 2.3). Taking into account the high probability of missing annotations in several species (particularly hominids), it seems prudent to define primate specificity of miRNAs not by presence in primates, but rather by absence of the miRNAs in non-primate species (also excluding miRNAs *only* annotated in human). Thus, primate specificity of a human miRNA is assumed if the miRNA is expressed in at least one primate species, and absent from all non-primate species in this roster. This definition yields a list of 377 primary and 350 mature putative “primate specific” miRNAs in miRBase v21 (Appendix B). Judging from recent analyses,¹⁰⁰ there probably exist many more. The primate-specificity attribute was entered into *miRNetDB* as miRNA node property.

2.3 MIRNETDB USAGE

Neo4j uses a language (called »Cypher«) akin to SQL, which utilises keyphrases to issue commands, but combines it with a semi-graphical syntax to account for the graph-based layout of the data. In the following, I will describe its basic usage and the advantages it provides in the matter of transcriptional connectomics. The basic »finder« function (similar to **SELECT** in SQL) is called **MATCH** in Cypher, and, when combined with the semi-graphical syntax, can be used to identify nodes or more complex patterns in the database. The graphical syntax consists of two main building blocks that represent the basic types of data inside the database: nodes as regular brackets »()« and edges between nodes as a construct of hyphens and box brackets, that can also have a direction indicated by the greater sign »()-[]->()«. To specify the elements to be found, attributes of nodes and/or edges can be filtered by using curly brackets in the node definition, or the **WHERE** clause. To be returned, elements need to be assigned arbitrary variable names:

Listing 2.1: MATCH

```
1 MATCH (gene:GENE {species: 'HSA'})  
2 WHERE gene.name = 'CHAT'  
3 RETURN gene
```

Query 2.1 identifies a node (arbitrarily designated »gene«) with type GENE (indicated by the colon), with attributes »species« (HSA, i.e. *H. sapiens*) and »name« (CHAT), and returns the node with all its attributes. Since the nodes of type GENE are restrained, there can only be one gene of

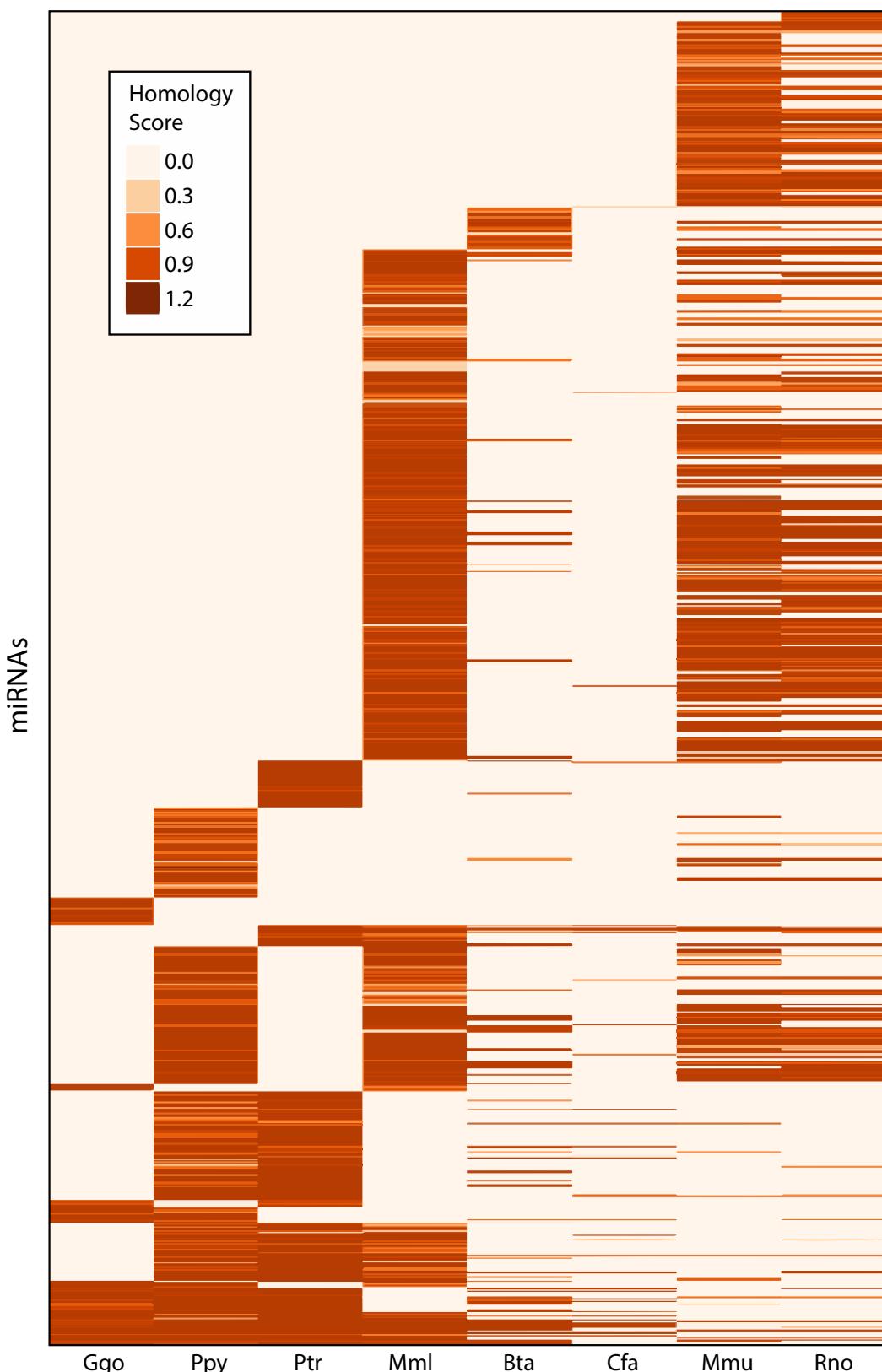


Figure 2.3: Homologues of Human microRNAs in Primate- and Non-Primate-Species. Homology to human miRNAs was determined by Smith-Waterman local alignment for each homologous miRNA of 8 species. Homology scores were visualised on a heatmap, each column represents the homology to human of the miRNAs of the respective species. The heatmap is ordered from bottom to top by the amount of miRNA homologues in primates. The miRNAs at the very bottom are shared by human as well as all four primate species, followed by the miRNAs shared by three primate species, and so on. Ggo: *Gorilla gorilla*, Ppy: *Pongo pygmaeus* (Orangutan), Ptr: *Pan troglodytes* (Chimp), Mmt: *Macaca mulatta* (Rhesus macaque), Bta: *Bos taurus* (Cow), Cfa: *Canis familiaris* (Dog), Mmu: *Mus musculus* (Mouse), Rno: *Rattus norvegicus* (Rat).

species *H. sapiens* with this name in the database, and thus, only one data point will be returned. The graphical syntax further allows for pattern matching of, for instance, miRNA→gene relationships:

Listing 2.2: Patterns

```
1 MATCH (mir:MIR)-[rel:TARGETS]->(gene:GENE {species: 'HSA'})  
2 WHERE gene.name = 'CHAT'  
3 RETURN mir, rel, gene
```

Query 2.2, similar to query 2.1, starts by identifying the node of species HSA with the name CHAT, and proceeds to look for miRNA→gene relationship edges arriving at this node; the relationships have to be of the type TARGETS (the pre-aggregated score-based accumulation of targeting). As soon as no further edges are found, the process terminates and returns all found miRNAs (»mir«), relationships (»rel«), and genes (»gene«) in discrete form, including all their attributes, such as the ENSG and Entrez IDs, the MIMAT IDs for all found miRNAs, or the score value of their targeting relationship. In this query, since there is a constraint on genes, the only gene returned is *CHAT*. However, Cypher is not limited to filtering on unique attributes; it allows for query and return of as many data points as are needed. For example, if one is interested in all miRNA→gene interactions in the cholinergic system, the query might look as follows:

Listing 2.3: Filtering

```
1 MATCH (mir:MIR)-[rel:TARGETS]->(gene:GENE {species: 'HSA'})  
2 WHERE gene.name IN {cholinergic_genes}  
3 RETURN mir, rel, gene
```

The effectiveness of graph-based databases becomes clear in this approach: Query 2.3 is processed starting at a user-defined filter, the list of cholinergic genes as an input (containing *CHAT*, *SLC18A3*, cholinergic receptor genes, acetylcholinesterase, etc.). In a first step, all nodes are found that fulfil the criteria: type GENE, from species *H. sapiens*, that are in the list of names given. Since the gene nodes are indexed, this only requires milliseconds. Then, through the connection of edges to these nodes, it finds all miRNA nodes that have a miRNA→gene relationship towards any of the cholinergic genes. By using the gene nodes as starting point, the query can end as soon as no other edges fulfilling these criteria are found on any of the nodes. In comparison, to satisfy this query in a relational database, the rows representing these cholinergic genes would have to be assessed in their entirety, not only in those columns that represent an extant relationship, thus prolonging execution.

The database then returns all miRNA→gene relationships in this set, representing the network of cholinergic miRNA regulators, including all of their attributes. The advantages of graph-based data do not end there; say one wants to return only »master« regulators of cholinergic systems, defined as miRNAs that target at least 4 of the genes in the cholinergic set. In a relational database, this would

have to be done post-hoc, by aggregation of relationships and removal of any results that do not exceed this threshold. This requires storage of the entire result in memory, and additional computational steps that can be very taxing depending on the size of the result table. In Cypher, this can be done during the query (code comments indicated by `//` explain single steps):

Listing 2.4: Two-stage Filtering

```

1 MATCH (gene:GENE {species: 'HSA'})
2 WHERE gene.name IN {cholinergic_genes}
3 WITH gene //the found genes are used as input for the second query
4 MATCH (mir:MIR)-[rel:TARGETS]->(gene)
5 WHERE count(rel) >= 4
6 RETURN mir, rel, gene

```

Query 2.4 essentially proceeds in the same way as query 2.3 in that it identifies the gene nodes filtered for and looks for the miRNAs connected to those nodes by TARGETS-type relationships; however, in the second step (which is performed per gene node as returned by the `WITH` clause), it returns only those patterns that have at least 4 incoming miRNA→gene relationships. Query 2.4 only requires little additional processing compared to query 2.3, and thus does not require nearly as much time as the post-hoc filtering required in a relational database query. This filtering can be applied in many stages, and in many forms, such as sums, averages, maximum and minimum, or other combinations of arithmetic and logical classifiers. Additionally, the patterns can be extended to represent complex relationships inside the graph. For instance, the following query 2.5 was used to find miRNAs that regulate any given gene in the database, and, simultaneously, affect TFs that are involved in regulation of this same gene (this type of interaction is called feedforward loop, see also Section ??).

Listing 2.5: Feedforward Loop Identification

```

1 MATCH (gene:GENE) //find gene
2 WHERE gene.id = ID //by identifier (Entrez)
3 WITH gene //use as input for next step
4 MATCH (tf:GENE {species: 'HSA', tf:TRUE})-[rel]->(gene)
5 //find TFs targeting that gene
6 WHERE type(rel) IN {tissue_types} //TFs only from specific tissues
7 //for instance, CNS cell types (Appendix A)
8 WITH gene, rel, tf //use as input for next step
9 MATCH (gene)<-[rel_m1:TARGETS]-(mir:MIR {species:
  'HSA'})-[rel_m2:TARGETS]->(tf)
10 //find miRNAs that target both gene and TF

```

```
11 WHERE rel_m1.score > 5 AND rel_m2.score > 5
12 //low-cut filter at a minimum cumulative score of 6
13 RETURN gene, tf, rel, type(rel) AS tissue, mir, rel_m1, rel_m2
```

This analysis can be performed in real time, on the whole genome and miRNome, and merely takes seconds for one iteration, a performance unimaginable in a relational database approach; advanced statistical approaches such as permutation only become viable at this timescale.

2.4 STATISTICAL APPROACH TO TRANSCRIPTIONAL CONNECTOMICS

The enormous amounts of data generated by modern molecular biology methods, such as RNA-seq and bioinformatics, present new challenges to statistical methodology. A major objective in the analysis of large datasets is a robust statistical representation of the distribution of this data. Traditionally used approaches such as Student's t-test are not automatically applicable to the intermediary results of these modern methods, because the premise of a normal distribution often does not hold, or has to be proven first. This section will describe the statistical problems encountered in the analysis of intermediary data produced by *miRNetDB*; the statistical properties of large count data directly generated by RNA-seq will be discussed in Section 3.6.3.

2.4.1 Permutation

The evaluation of comprehensive prediction datasets regarding miRNA→gene interactions on a genome scale is statistically challenging. Molecular interaction studies have explored only a minority of all possible targeting relationships, and as such, the ground truth of miRNA→gene interaction is unknown (see Section 2.2.4). Since there is no negative interaction data, validated interactions can only be defined in the positive space. Additionally, the various prediction algorithms also heavily diverge in their predictions, which leads to the question of how to approach the estimation of false discovery ratio (FDR) while simultaneously avoiding high false negative rates.

One possible approach that can aid in identification of the most pertinent effects in this case is random permutation. In this approach, the result of an analysis (e.g., a numeric targeting score of a miRNA→gene interaction, or a Spearman correlation between two gene sets) is compared to a null distribution that was generated from an iterative analysis similar to the initial one, but with randomised input (e.g., a group of miRNAs of the same size as the original set, randomly selected from all miRNAs, or the gene sets from the original analysis with randomly scrambled group affiliations). This permutation of the analysis is performed many times (usually between 10 000 and 1 000 000 iterations, depending on the context), and results in a distribution of possible outcomes that can be arranged from lowest to highest, often resulting in a normal (or »normal-like«) distribution, thus facilitating the estimation of confidence intervals, and, similarly, p-values for the »real« result.

A positive side-effect of performing a permutation analysis on a base collection of data, such as *miRNetDB*, is the automatic correction of inherent biases. For instance, should a particular gene by its genetic structure invite a large amount of false positive predictions as to the miRNA→gene interactions towards it, these will be

present in the test as well as in the permutation comparison, and thus cancel out and yield a high p-value for this interaction, effectively transforming the false positive into a true negative.

2.4.2 Gene Set Enrichment Analysis

The objective of gene set enrichment is the identification of statistically over-represented entities in a dataset. The standard use case in biomedicine is the Gene Set Enrichment Analysis (GSEA), that is used to identify the most important classes of genes in large datasets, such as the ones produced by RNA-seq. Briefly, the analysis follows these steps: the studied genes are scored by a certain method, such as p-values from differential expression analysis, which enables the identification of a relevant subgroup, the test set (e.g., the 100 genes with lowest p-values). This test set is then compared to a background of genes (usually, all detected genes, or a large amount of genes from the entire dataset) by a statistical method fit to determine their enrichment in pre-defined categories. Often, ontological categories are used, such as the »biological process« type of Gene Ontology (GO), or KEGG pathways.

For each of these categories, the method tests for a representation of genes in the test set exceeding the frequency statistically expected by random sampling from the background of genes; thus enabling an estimation of the functionality these test set genes might inhabit in the process that is studied. Statistical approaches often employed in gene set enrichment are Kolmogorov-Smirnov statistics, permutations, or, more generally, hypergeometric tests such as Fisher's exact test. There are a wide variety of software solutions available for the implementation of gene set enrichment testing.

Gene Ontology curates an enormous catalogue of coding gene products and their functions. At the current time, GO hosts 7 330 378 annotations (2 836 377 for »biological process«, 2 289 165 for »molecular function«, and 2 204 836 for »cellular component«), subdividing 1 405 197 individual gene products from 4493 species (205 with more than 1000 annotations) into 44 733 ontological terms (29 457 »biological process«, 11 093 »molecular function«, and 4183 »cellular component« terms). The individual GO categories are organised in a hierarchical manner, more specifically, a directed acyclic graph (DAG). Each branch of the DAG tree contains related terms, progressing from the most general terms (top) to the most specific ones (at the bottom).

Whenever a GO analysis is described in this dissertation, it means a gene set enrichment analysis performed on a particular subset of genes (that might e.g. be the targets of a group of miRNAs) towards the elucidation of their biological function, i.e., the »biological process« category of GO annotation.

*There is no scientific study more vital to man
than the study of his own brain. Our entire view
of the universe depends on it.*

Francis Crick

3

microRNA Dynamics in Cholinergic Differentiation of Human Neuronal Cells

This chapter will discuss the current state of knowledge on brain transcriptomics, generally and in the specific case of cholinergic neurons in the CNS, and then go on to explain the steps we undertook to elucidate small RNA processes in central cholinergic systems. First, our aim was to clarify co-expression patterns of central cholinergic neurons, which required analysis of transcriptome data in single-cell resolution. Based on this information, we selected two human models of cholinergic neuronal differentiation and established a differentiation protocol amenable to RNA extraction and successive molecular biology assays, most importantly, RNA-seq. The expression patterns so obtained were then used to perform bioinformatics analyses using the database introduced in Chapter 2, *miR-NetDB*.

3.1 NEURONAL TRANSCRIPTOMES - BACKGROUND

The mammalian brain requires a constant supply of oxygen and nutrients, because it does not provide storage for either. Though it only makes up approximately 2% of the entire human body mass, its energy expenditure is around 20% of the whole.¹¹⁵ For this reason, each square millimetre of brain tissue (except for the ventricles) is infiltrated by hundreds of capillaries.¹¹⁶ Since the blood-brain-barrier is essentially provided by supporting glia cells surrounding all capillaries from the »inside« (see Fig. 3.1, modified from Lobentanzer & Klein¹¹⁷), neurons numerically constitute only a minority of brain tissues (but burn two thirds of its energy).

Until very recently, studies aiming to clarify the transcriptional profiles of neurons applied either

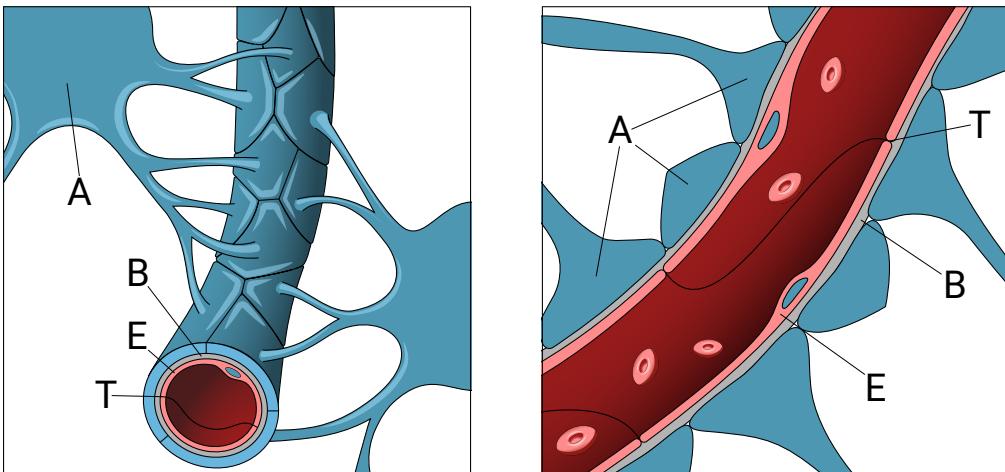


Figure 3.1: Schematic display of the blood-brain-barrier. The blood-brain-barrier surrounds virtually every capillary in the CNS. A: Astrocyte, B: Basal Membrane, E: Endothelial Cell, T: Tight Junction. Modified from Lobentanzer & Klein, 2019.¹¹⁷

microarray technology or RNA-seq (also known as deep sequencing or next generation sequencing). For these methods, several cubic millimetres of brain tissue are required at the least; often, cubic centimetres are used. In contrast, the diameter of neuronal somata is usually in the micrometre range. Thus, the resolution of the method and the actual cellular resolution differ by a factor of approximately 1000. Additionally, even among the neuronal population, there is considerable heterogeneity and transcriptomic plurality; single brain regions rarely consist of less than 30 different neuron types, tightly packed next to each other, each with their own transcriptional identity.^{118,119,120,121} Newest studies, deciphering the murine nervous system by sequencing of 500 000 individual cells, show that neuron diversity is very similar regardless of brain region.¹²² These circumstances hold true for any mammal, and most of our knowledge stems from the analysis of our favourite research animal, the mouse. In humans, the diversity is only exacerbated; in fact, the elevation in CNS complexity, which is only made possible by enhanced transcriptional control, may be the reason for our superior cognitive abilities(cite).

Cholinergic neurons always constitute a minority in any neuronal population, sometimes to extremes. Most tissues are dominated by few neuron types, such as pyramidal cells in the cortex. The most common neurotransmitter types are GABAergic (inhibitory) and glutamatergic (excitatory), each with several subtypes. It is estimated that more than 80% of cortical neurons are excitatory, and more than 90% of synapses release glutamate.¹¹⁵ There are two major cholinergic regions in the mammalian brain: The striatum is fairly well-populated with rather large cholinergic interneurons, and the basal forebrain holds a large amount of (smaller) cholinergic projection neurons (compare Fig. 1.1). However, in transcriptomic analyses, these tissues are seldom used, maybe due to lack of scientific interest, or because they are notoriously hard to access (the basal forebrain is small and deeply imbedded in the midbrain). The cortex, particularly the neocortex, is most often the tissue of choice

in these studies, due to its scientific interest and accessibility. Though it contains only a minuscule amount of cholinergic interneurons whose transcriptional identity is still a matter of debate, several of the recent single-cell RNA-seq approaches have independently identified cholinergic interneurons in cortical regions (see Fig. 3.2).

3.2 CORTICAL SINGLE-CELL RNA SEQUENCING

THE IMPACT OF TRANSCRIPTIONAL DYNAMICS on any disease depends on co-expression of the relevant genes in the affected cell. Selection of a model therefore has to take co-expression into account. In particular, if neurokines are to possess any relevance for cholinergic properties of central nervous cells, the cells in question would have to express molecular machinery required to receive neurokinin signals. The advent of single-cell RNA-seq for the first time enables the resolution of gene expression on a cellular basis, and thus the disentangling of spatially close individual neuron types (and other, non-neuronal CNS cells); most of this information is lost in RNA-seq performed on brain homogenate, even of a small biopsy. Differences in genes are reduced to the universally expressed »housekeeping« genes, save the most extreme perturbations. In miRNAs, this circumstance is only exacerbated, in parallel to their even more tissue-specific expression.

To provide a detailed tally of transcriptional subtypes in the CNS, publicly available single-cell RNA-seq datasets of suitable tissues were analysed towards their cholinergic properties. All studies that were available at the time focused on some subsection of the cortex (visual or somatosensory) or the hippocampus. Additionally, the data provided by those studies was in some cases pre-aggregated to represent classes of single neurons with similar transcriptomes (Fig. 3.2 A&B^{119,120}); in other cases, every single neuron was represented (Fig. 3.2 C&D^{118,121}).

An important quality-related parameter of a single-cell RNA-seq experiment is the sequencing depth achieved per single sequenced cell. Some of the screened datasets do not provide sufficient depth to resolve genes with medium expression, which includes our primary cholinergic markers *CHAT* and *SLC18A3*. The datasets which did provide adequate sequencing depth were filtered for their expression of these markers, and additionally characterised by their expression of common markers for cell types to be expected in the CNS. Raw data were downloaded from their respective sources and imported into the R environment, where they were converted into similar format. The numeric expression values of each dataset were normalised to transcripts per million (TPM) to allow comparison (with counts n and transcript length ℓ of gene A and all genes i per sample):

$$\frac{\frac{n_A}{\ell_A}}{\sum_i \frac{n_i}{\ell_i}} \times 10^6$$

For graphical display, TPM were further normalised to a range of 0-1. The cholinergic genes were filtered from each dataset and plotted as heatmaps. Plotted were only samples that expressed *CHAT*, *SLC18A3* (also known as vAChT), and/or *SLC5A7* (also known as HACU).

The identified samples provide an overview of potentially cholinergic cells in the sampled brain regions, and allow an assessment of the functional type and gene co-expression patterns in central cholinergic cells (Fig. 3.2). Most cells identified as cholinergic by this definition expressed the gen-

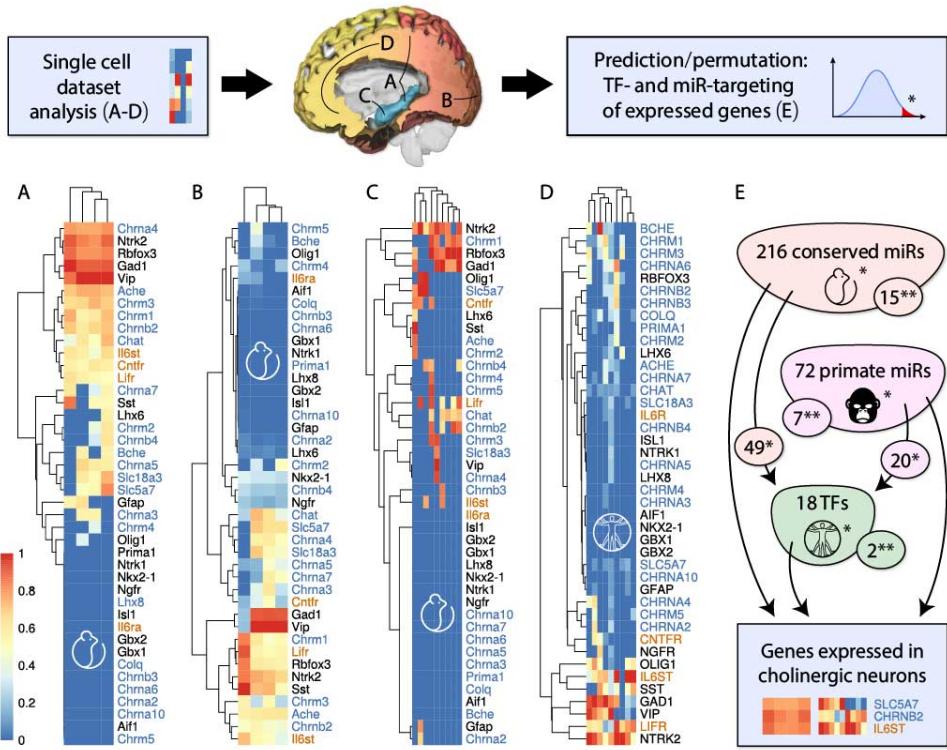


Figure 3.2: Single-Cell Sequencing of CNS Tissues. Expression patterns of cholinergic and cholinergic-related genes were analysed using web-available single-cell sequencing datasets. Expression was normalised to reflect a span between 0 and 1. **A)** Clustered single-cell sequences from transgenic mouse somatosensory cortex and hippocampus.¹¹⁹ **B)** Clustered single-cell sequences from transgenic mouse visual cortex.¹²⁰ **C)** Single-nucleus sequencing of adult mouse hippocampus.¹²¹ **D)** Single-cell sequencing of the human developing neocortex.¹¹⁸

Several neuronal marker *RBFOX3*, also known by its trivial name NeuN, but not the microglial marker *AIF1*. Few cells (or clusters of cells) expressed non-neuronal markers such as *GFAP* (astrocytes) or *OLIG1* (oligodendrocytes), hinting at sparse non-neuronal cholinergic functions. In agreement with my findings, cells or clusters identified as cholinergic by the authors of the respective studies^{119,120} (also by personal communication with Peter Lönnérberg) had been classified as interneurons and co-expressed a number of known phenotypic neuronal markers, such as *somatostatin (SST)* and *vasoactive intestinal peptide (VIP)*.

The identified cholinergic cells also revealed a constant co-expression with neurokinin-related genes, particularly the transmembrane neurokinin receptors *LIFR* and *IL6ST*, demonstrating a capacity to receive and process neurokinin signals. In contrast, the high affinity receptor for NGF, *NTRK1*, is not co-expressed in mature (NeuN-positive) cholinergic neurons in the analysed regions, fundamentally distinguishing these cells from the basal forebrain cholinergic projection neurons.

3.3 MICRORNA AND TRANSCRIPTION FACTOR TARGETING PREDICTIONS

Making use of the information aggregated in *miRNetDB*, the genes identified as being expressed in cholinergic neurons were subjected to permutation targeting analyses of miRNAs and TFs. Genes were assumed to be expressed in cholinergic neurons if they were expressed in more than one individual sample in all single-cell

RNA-seq datasets (Fig. 3.2 A-D). The TFs identified as active towards cholinergic genes in cholinergic neurons were additionally subjected to another round of miRNA targeting permutation analysis. Targeting of genes with random selections of miRNAs and TFs were permuted XX times to estimate FDR. Statistical significance of the miRNA→gene or TF→gene interactions was assumed at FDR < 0.05.

Permutation targeting analyses revealed a nested regulatory interaction between 72 primate-specific miRNAs, 216 conserved miRNAs, and 18 TFs towards cholinergic genes expressed in cholinergic neurons (Fig. 3.2 E). TFs targeting cholinergic genes were in turn targeted by 49 conserved and 20 primate-specific miRNAs that also targeted cholinergic genes directly.

3.4 Gene Clustering Based On Expression

Hierarchic clustering was applied to expression data to identify functional grouping of genes and cells based on co-expression. Initially, samples (i.e., single cells, pre-aggregated clusters of cells, or brain regions) are compared using a similarity- or distance-matrix (where similarity = 1 - distance). The similarity measure is based on a computation according to the method used. For instance, Euclidean distance between two gene expression vectors (i.e., samples) of length n is the distance between points p and q in n -dimensional space, defined by:

$$d_E(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

Applying this measure to all pairwise combinations of samples results in a dissimilarity matrix that can be converted to a hierarchy using one of several clustering algorithms. Generally, samples are grouped by their similarity. Initially, each sample is assigned to its own cluster, and then, cluster number is iteratively reduced by joining the closest clusters. This results in a hierarchic tree of samples, that can be »cut« at any height to yield an arbitrary number of clusters. In biological analyses, the method after Ward (in R, »Ward.D2«) is often used.¹²³

Due to the structure of the data (small number of entities compared to whole genome analysis, repetition of zeroes in individual samples), the Bray-Curtis dissimilarity¹²⁴ is superior to Euclidean distance. Bray-Curtis dissimilarity is defined as:

$$d_{BC}(p, q) = \frac{2C_{p,q}}{S_p + S_q}$$

Where C is the sum of the lesser expression values common to both vectors p and q , and S is the total number of genes expressed in each sample (i.e., values greater than zero in each vector). Based on this measure, the samples were clustered according to their cholinergic gene expression levels using Ward's method to yield five separate clusters. Intermediary clustering results (not shown) revealed a uniform distribution of ATP citrate lyase (ACLY), yielding no additional information; thus, it was removed. Also removed for the purpose of clustering were the non-neuronal nicotinic receptor subunits $\alpha 1$, $\beta 1$, γ , δ , and ϵ .

3.4.1 Co-EXPRESSION OF FUNCTIONAL GROUPS OF CHOLINERGIC GENES

Hierarchic clustering of cholinergic genes in each of the datasets revealed a grouping of cholinergic genes according to their biological function. Table 3.1 shows considerable uniformity in two single-cell mouse datasets, which diverge substantially from the brain-region- and TF-based human set. Gen-

erally, clustering shows separation of at least 3 groups of cells, one of which is the classic *cholinergic* neuron with genes for synthesis and transport of acetylcholine. Due to the frequent co-expression of *CHAT* and *SLC18A3*, it is safe to assume the *SLC18A3* as a viable substitute for *chat* expression and clustering in the FANTOM5 data of Marbach *et al.*⁹⁸ (for more details, see Section 2.2.3). In the single-cell datasets, the *CHAT* gene is expressed in parallel with the two cholinergic transporters, without exemption. The other groups could be described as *receptive* neuron (not cholinergic as the aforementioned, but different types of cholinergic receptors and esterase) and other, rather specialised groups, probably comprising various glial cells. These last, specialised groups are not very visible in the human dataset, which lacks the single cell resolution of the mouse datasets and therefore includes glial cells in every sample of any region. Therefore, differences in cholinergic gene expression patterns derived from Marbach et al are likely the result of the numbers and dominant types of cholinergic neurons in the respective region.

Functional stratification of cholinergic genes is also visible in a dendrogram of gene clusters from all four analysed single-cell sequencing datasets (Fig. 3.3). While there is variability in the composition of receptor subunits (which is to be expected regarding the different sampled brain regions), the core cholinergic genes (such as *CHAT*, *SLC18A3*, *SLC5A7*, and *ACHE*) associate similarly in all datasets. Notably, the distinction between a *cholinergic* and a *cholinoreceptive* neuron is always visible by a grouping of, on one hand the synthesis, vesicular packaging, and reuptake of ACh, and on the other hand, cholinergic receptors and signal termination by AChE.

3.5 THE CELLULAR MODEL

We selected two mono-cultures of human neuronal cells for subsequent experiments: LA-N-2 and LA-N-5. During the selection process, multiple options were considered. Multicellular models would, in principle, allow disentanglement of the functions of distinct cell types, for instance glia and neurons. This could be achieved by *in vivo* or *ex vivo* approaches in rodents. However, our diseases of interest (Section 1.2) display a noticeable lack of transferability from lower mammals to human(cite). Alternatively, co-cultures of human cells in mono-layer or as 3D-culture have been proposed, but these still lack experimental stability.

3.5.1 THE SH-SY5Y NEUROBLASTOMA CELL LINE

A prominent example of human neuronal cell culture used in the identification and elucidation of cholinergic processes is the immortalised neuroblastoma cell line SH-SY5Y.¹²⁵ Derived from its parent line SK-N-SH, an adrenergic neuroblastoma,¹²⁶ it expresses ample amounts of *ACHE*, and thus had become a work horse in many cholinergic fields, such as Alzheimer’s Disease (which is treated with AChE inhibitors), pesticide development, and warfare(cite). However, in spite of its usefulness for processes involving *ACHE*, it turned out a less than optimal choice for the study of molecu-

cluster	Zeisel et al	Tasic et al	Marbach et al
I	Ache, Chrm1, Chrm2, Chrm3, Chrm4, Chrna4, Chrna5, Chrna7, Chrnb2	Ache, Chrm1, Chrm2, Chrm3, Chrm4, Chrna2, Chrna4, Chrna5, Chrna7, Chrnb2, Chrnb4	CHRM1, CHRM2, CHRM5, CHRNA2, CHRNA4, CHRNA6, CHRN2, CHRN3
Ib			ACHE, BCHE, CHRNA3, CHRN2, PRIMA1
Ic			CHRM3
Id			CHRNA5, CHRNA9
II	Chat, Chrnb4, Slc18a3, Slc5a7	Chat, Chrm5, Chrna3, Slc18a3, Slc5a7	SLC18A3, SLC5A7
III	Chrm5, Chrna10, Chrna3	Chrna10	
IV	Bche, Prima1	Bche, Prima1	
V	Chrna2, Chrna6, Chrnb3	Chrna6, Chrnb3	

Table 3.1: Cholinergic gene clusters according to cell type vs. brain region. The two transgenic mouse datasets from Zeisel et al and Tasic et al show high similarity in gene distribution. With high likeliness, cluster I is a group of postsynaptically cholinergic, "receptive" cells. Cluster II represents the classic "cholinergic" neuron, with synthesis, vesicular packaging and ACh-reuptake genes. The transcription factor-based dataset of Marbach et al depends on whole brain regions instead of single cells to determine similarity, and thus yields distinctly different classification. However, it also distinguishes between cholinergic synthesis (with SLC18A3 as a substitute for CHAT expression) and cholinoreceptive functions.

lar events surrounding *CHAT* and *SLC18A3*, as it barely expresses both genes(cite), and cannot be coerced to elevate *CHAT* expression by the usual differentiation techniques (own experimentation, data not shown). Thus, for the questions asked in this chapter of the dissertation, SH-SY5Y does not qualify as adequate representation of a »cholinergic neuron«.

3.5.2 THE LA-N NEUROBLASTOMA CELL LINES

Following the elimination of SH-SY5Y as a suitable subject, a literature search for candidates representing a cholinergic neuronal transcriptome revealed, among others, representatives of the LA-N neuroblastoma cell lines developed by R.C. Seeger around 1980.^{127,128} Neuroblastoma is a form of neuronal cancer often affecting small children, and, consequently, the two cell lines used in my experiments are immortalised biopsies of a 3 year old girl (LA-N-2¹²⁷) and of a 4 month old boy (LA-N-5¹²⁸). The decision to use LA-N-2 as initial cellular model was influenced by three factors: it is well described in literature, although most studies had been published in the 1980s and 90s; it expresses substantial amounts of *CHAT* and *SLC18A3*(cite); and it responds to neurokine-mediated differentiation by assuming a neuronal morphology accompanied by further elevation of *CHAT* and *SLC18A3* expression. LA-N-5 was not nearly as well described as LA-N-2, but later added to the experimental roster because of the complementary sex and hints towards cholinergic differentiation

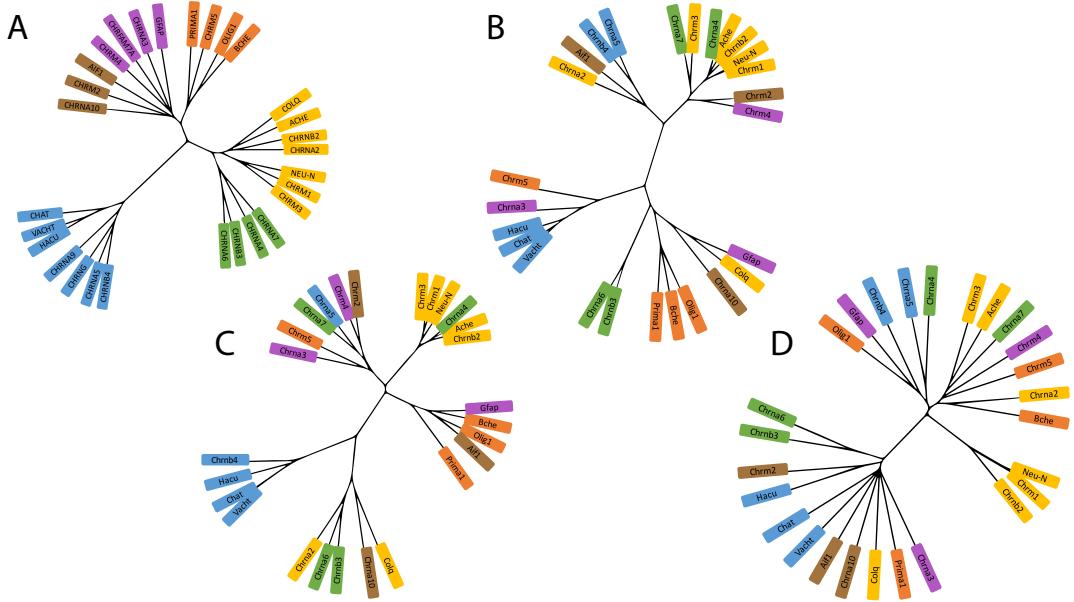


Figure 3.3: Clusters of Cholinergic Genes in Single-Cell Sequencing. Cholinergic genes were clustered using Bray-Curtis dissimilarity in four public data sets of single-cell sequencing. The displayed dendograms visualise the distance between the genes across all samples. Gene clusters were coloured by grouping in Darmanis et al.¹¹⁸ (A). Notably, genes clustered according to their biological function, for instance, CHAT, vAChT and HACU always are closely associated (blue), as are the genes comprising the putative »cholinceptive« neuron (yellow). A) Single-cell sequencing of the human developing neocortex.¹¹⁸ B) Clustered single-cell sequences from transgenic mouse visual cortex.¹²⁰ C) Clustered single-cell sequences from transgenic mouse somatosensory cortex and hippocampus.¹¹⁹ D) Single-nucleus sequencing of adult mouse hippocampus.¹²¹ Marker genes are: NeuN - neurons; OLIG1 - oligodendrocytes; AIF1 - microglia; GFAP - astrocytes.

under retinoic acid.¹²⁹

3.5.3 Culture

LA-N-2 and LA-N-5 are very similar in their culture requirements. They have comparatively high duplication times, which can be lowered by using certain conditions that affect medium composition, nutrition, and CO₂ content. The cells were acquired at DSMZ (Braunschweig, Germany), which recommends keeping them in a 50:50 mixture of Dulbecco's modified eagle medium (DMEM) and Roswell Park Memorial Institute medium (RPMI1640), with 20% fetal calf serum (FCS) added. Sometimes, recommendations also suggest Leibovitz's L-15 medium, which is specifically designed for low CO₂ conditions, and others have suggested increased CO₂ levels inside the incubator. A combination of the DSMZ-recommended medium with 8% CO₂ atmosphere inside a 37°C incubator to accelerate growth to a degree that the cells could be split 1:3 to 1:4 in a weekly cycle. This protocol was used for all further experiments, which were performed between splits 2 to 8 after thawing of a batch from -80°C. All handling during maintenance and experimentation was performed under a laminar flow hood.

3.5.4 Differentiation

Neuronal differentiation of neuroblastoma cell lines has been performed in many instances, utilising a wide variety of differentiation agents such as the very general retinoic acid or 5-bromo-uracil, or very specific reagents, such as the neurokines IL-6 and CNTF(cite). LA-N cells have also been described to react to a selection of these substances; however, due to our elevated interest in neurokine mechanisms, we opted for a neurokine-based

differentiation protocol. In personal communication, James McManaman revealed that the »CHAT development factor« that he had discovered⁶¹ was, in fact, CNTF, which had never been published. Additionally, of the neurokines used for differentiation purposes, CNTF is best described in literature and easily acquired in dried form from Merck (formerly SigmaAldrich, Darmstadt, Germany). CNTF was resuspended in pure water to a concentration of $25 \mu\text{g ml}^{-1}$ and stored for experimentation in aliquots at -20°C.

LA-N cells are very sensitive to repeated temperature changes (or other handling-related disturbances), which resulted in increased amounts of apoptotic cells following repeated removal from the incubator after seeding or medium changes during the experiment (Lobentanzer, not published). For this reason, the differentiation reagent was only added once, 24h after initial seeding of the cells, and further disturbances avoided until the time of lysis. For the maximum duration of the experiments, 120h from seeding until lysis, the initially supplied medium was sufficient for survival.

Differentiation was performed in regular growth medium without changes in FCS content, and CNTF was added to the medium after an initial growth period of 24h. Cells were seeded into 12-well plates at approximately 200 000 cells/well, with 1 ml of growth medium. To determine the optimal amount of CNTF for differentiation, time-dose experiments were performed for both cell lines in a range from 1 ng ml^{-1} to 100 ng ml^{-1} for several time points during four days. Here, we discovered the first pharmacological difference between LA-N-2 and LA-N-5: the maximum of their cholinergic response to neurokin stimulation (i.e., an elevation in CHAT and SLC18A3 transcription) occurs at different concentrations of CNTF. While LA-N-2 cells respond most strongly to 100 ng ml^{-1} , LA-N-5 cells show an »inverted u«-type dose response with a maximum around 10 ng ml^{-1} CNTF (Fig. 3.4). James McManaman, who studied LA-N differentiation thoroughly in the 1990s,¹³⁰ believes both lines to respond in an »inverted u«-type manner (personal communication); thus, in all likelihood the LA-N-2 response would also diminish at CNTF concentrations significantly higher than 100 ng ml^{-1} . Also, CNTF concentrations could likely be significantly lowered by removal of the high amount of FCS in the medium, however, that would require the use of a special serum-free medium, which would have to be established up front, and may have other, unforeseen consequences. Regardless, CNTF concentrations around 100 ng ml^{-1} (i.e., pico- to nano-molar) still are well within the physiological range of concentrations that the mammalian brain is able to reach by paracrine secretion via, e.g., astrocytes.¹³¹

To study the small RNA dynamics following CNTF exposure of LA-N-2 and LA-N-5, the experiment was stopped at 4 time points and the cells were quickly lysed *in situ* to preserve total RNA in that state: for the quick, immediate-early-like phase, at 30 and 60 minutes after the addition of CNTF, and, for the long-term effects of differentiation, at 48 and 96 hours after the addition of CNTF (Fig. 3.5, from Lobentanzer *et al.*⁶). Each time point was controlled by a pseudo-treated (using pure water) culture from the same batch that had been seeded at the same time as the experimental group. In the final series used for the parallel sequencing of LA-N-2 and LA-N-5, all experiments were carried out in quadruplicates.

combine fig
3.4, 3.5

vacht qpcr
figure?

3.5.5 RNA Isolation

Total RNA was isolated using TRIzol (ThermoFisher Scientific), essentially as suggested by the manufacturer, with slight changes to the protocol to enrich small RNA species. The cells, growing in a monolayer in 12-well-plates, were cleared of medium, washed two times with 500 µl of cell culture grade phosphate buffered saline (PBS) (Gibco), and immediately suspended in 1 ml of TRIzol, pipetting up and down until visibly dissolved. After incubation for 5 minutes at room temperature, the samples were stored in -20°C for short periods of time until

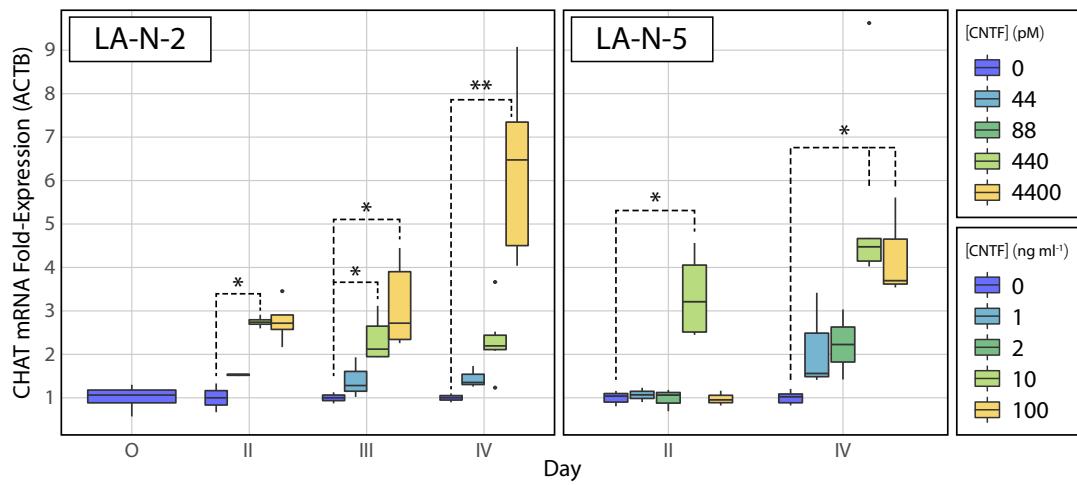


Figure 3.4: Time-dose curve of CNTF-mediated differentiation of LA-N-2 and LA-N-5. Cells were stimulated with varying doses of CNTF, and lysed at various time points to determine ChAT mRNA levels via qPCR. Expression ($\Delta\Delta C_t$) was normalised to housekeeping genes (ACTB, GAPDH, RPLP0) and to control sample without CNTF to determine fold-changes. LA-N-5 reacts strongest to a concentration of 10 ng ml^{-1} , while LA-N-2 reacts strongest to 100 ng ml^{-1} . *: $p < 0.05$, **: $p < 0.001$

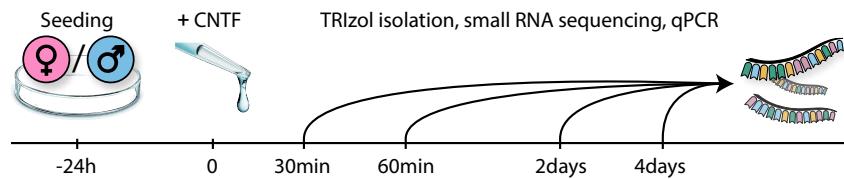


Figure 3.5: LA-N-2 / LA-N-5 Differentiation Timeline. Cells were seeded at $\sim 2E05$ cells/well in a 12-well-plate. After 24h, CNTF was added to the existing medium as quickly as possible to avoid disturbance. Cells were lysed *in situ* at time points 30 minutes, 60 minutes, 48 hours, and 96 hours using TRIzol for downstream RNA processing.

RNA isolation.

TRIzol-suspended lysates (1 ml) were added to RNA-separation centrifuge tubes (PhaseMaker Tubes, ThermoFisher Scientific), adding 200 µl of pure chloroform and mixing vigorously for 15 seconds. After two minutes, the mixture was centrifuged at 12 000 g and 4°C for 15 minutes, and the upper, watery phase containing the RNA was extracted. This was mixed with approximately 2 parts of pure ethanol and incubated for 10 minutes at room temperature to precipitate the RNA. The precipitate was spun at 12 000 g and 4°C for another 10 minutes, and the supernatant discarded. The pellet was washed with 85% ethanol (vortexed briefly) and centrifuged again for 5 minutes at 7500 g and 4°C.

After the final centrifugation step, the samples were transferred to the laminar flow hood, and air dried after removal of most of the supernatant via micropipettors. The pellet was allowed to dry almost until completion and resuspended in 30 µl to 50 µl pure RNase-free water. RNA concentration was measured at a Nanodrop 2000 instrument (ThermoFisher Scientific) and samples were diluted to a uniform concentration of 100 ng µl⁻¹. Finally, RNA samples were aliquoted according to later purpose and stored at -80°C.

RNA quality was determined by analysis on a 2100 Bioanalyzer instrument (Agilent) using a nano chip and 1 µl of sample; RNA integrity number (RIN) was near optimal for all samples (>9).

3.6 Small RNA Sequencing and Differential Expression Analysis

For the detection and analysis of small RNA species, RNA-seq is the current gold standard method. It allows the mapping of a comprehensive transcriptome and thus is vastly superior to small scale and consecutive methods such as real-time quantitative polymerase chain reaction (RT-qPCR), and even the larger scale microarrays. Microarrays, while also potentially allowing a »snapshot« of entire transcriptomes, are limited by the predetermined sequences on the chip. RNA-seq, on the other hand, is not biased towards any structural property of the sample; this is particularly important in the analysis of small RNA species, since their sequences are very variable (tRFs) and still not completely catalogued (miRNAs). Assuming an adequate sequencing depth ($\geq 1E06$ reads/sample), RNA-seq allows a comparison of all expressed small RNA species at once, which is immensely helpful when dealing with processes on the combinatorial scale of miRNA regulation.

3.6.1 Sequencing

For small RNA sequencing, the aliquoted samples were shipped on dry ice to the cooperating institute at the Hebrew University of Jerusalem, the Silberman Institute of Molecular Biology, the laboratory of Prof. Hermona Soreq. 600 ng of total RNA per sample were prepared for sequencing using the NEBNext Small RNA Library Prep Set for Illumina (New England BioLabs). The libraries were multiplexed with coloured barcodes, allowing for sequencing of all 48 samples on one chip. Briefly, this includes ligation of sequencing adapters to both 3' and 5' ends of all (single-stranded) RNA fragments in the sample, followed by 12-15 cycles of reverse transcription to form the RNA library. Ligated and amplified libraries were then size selected via gel electrophoresis on a 6% Polyacrylamide gel. The band representing small RNA species on the gel was excised and prepared for loading onto the sequencing chip. After loading, the chip was sequenced in a NextSeq 550 series instrument (Illumina) with a read length of 80 nucleotides (nt), single-end.

The quantity of reads per sample was determined by analysis of the raw fastq files. The read count across all samples before filtering was $7.8E06 \pm SD 2.5E06$, read count after quality filter and adapter removal was

$6.8E06 \pm SD 2.2E06$ ($n = 48$); a mean of 87% of reads remained after filtering, exceeding the recommended minimum amount (~1 million) by 4- to 12-fold (Fig. 3.6 A). Overall, ~326 million reads remained to be passed down to subsequent analyses. Sequencing quality was determined by analysis of the raw reads using the FastQC software(cite). Even before adapter removal and quality filtering, FastQC detected no »reads of poor quality« in any sample. Fig. 3.6 B gives a representative example of read quality per base (Sample 1).

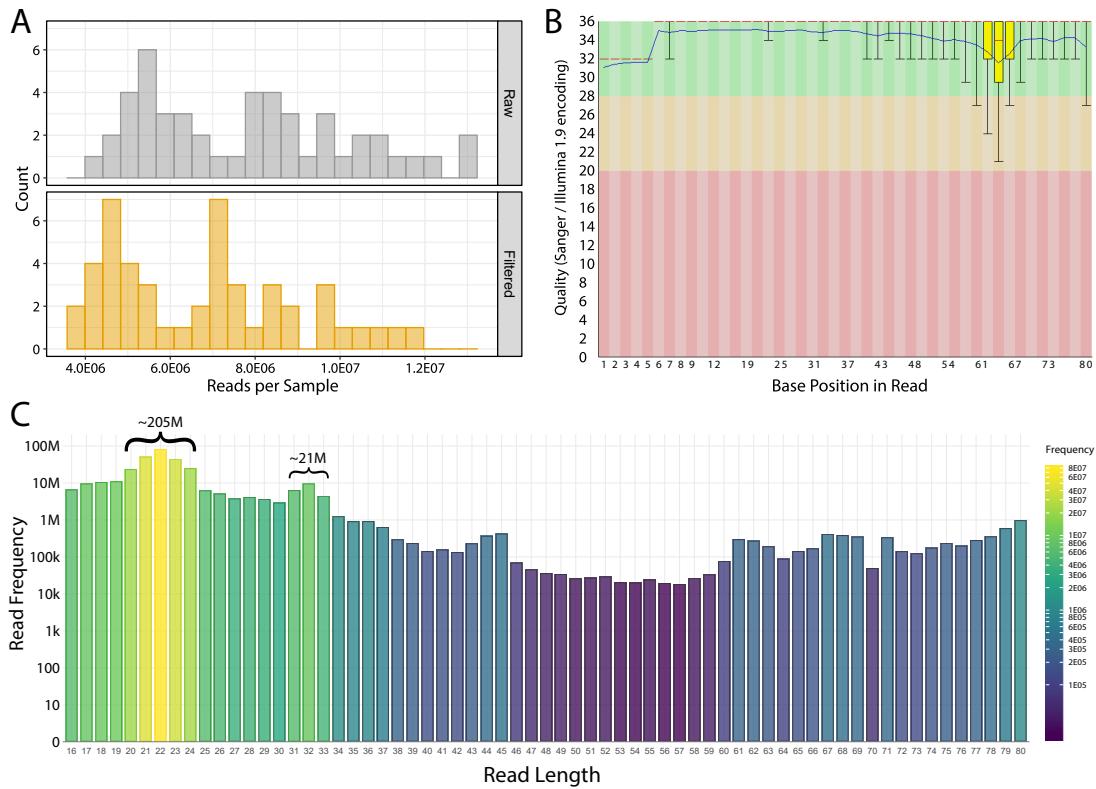


Figure 3.6: Small RNA Sequencing - Read Count, Quality, and Length. All samples provided near optimal quality. **A)** Per sample read count had a mean of $7.8E06 \pm SD 2.5E06$ in raw samples (top) and $6.8E06 \pm SD 2.2E06$ after quality filtering and adapter removal (bottom). 87% of reads were retained after filtering, with samples spanning read count values between 4 and 12 million. **B)** Representative example of quality score per base position in the sequencing (FastQC output of sample 1). Quality scores are always near the optimum, with a characteristic slight dip around nt 65. This occurs in all samples and is likely a technical result of the sequencing process. Possibly, it reflects the most common adapter ligation position after size selection of the RNA pool. **C)** Read length was determined for every one of the ~326 million reads. Nearly 80 million reads have a length of 22 nt, and the peak from 21 to 24 nt comprises ~205 million reads. This represents the bulk of miRNAs, and probably a significant amount of tRFs. The second peak, from 31 to 33 nt, still comprises ~21 million reads; these in all likelihood represent the longer tRNAs. The reads above a length of 33 nt only sum up to an amount of ~6 million, and may contain RNA of viral origin, or even mature tRNAs.

Raw reads were adapter-trimmed and quality filtered using the flexbar software ¹³² with parameters

```
-a adapters.fa -q TAIL -qf sanger -qw 4
-min-read-length 16 -n 1 --zip-output GZ
```

The sequence used in the *adapters.fa* file, as recommended by the manufacturer, was

AGATCGGAAGAGCACACGTCTGAACTCAGTCAC

Paired-end sequencing still is superfluous in small RNA-seq, because none of the common alignment pipelines can use the second (reverse) read, and manual paired alignment does not yield nearly as much benefit as the

depth increase in single-end sequencing (the read count per sample effectively doubles). 80 nt is the maximum read length possible in our small RNA workflow, and is excessive for the analysis of miRNAs. For transfer RNA fragments, however, a longer read can yield a more complete picture of expression, since the longer tRNAs can easily reach 40 nt in length. Indeed, the read length distribution after adapter removal shows a significant amount of small RNA species exceeding the length possible for miRNAs (Fig. 3.6 C).

align longer reads to genome? viral?

3.6.2 Sequence Alignment

For the alignment of miRNA sequences, parts of the miRExpress 2.0¹³³ pipeline were used according to the documentation. First, a lookup table for the current miRBase version 21 was created as per the instructions of the authors. The alignment was then performed using the commands *Raw_data_parse*, *statistics_reads*, *alignmentSIMD*, and *analysis*; *Trim_adaptor* was skipped because the adapters had already been trimmed in the quality filtering step. Additionally, since miRExpress is not accepting of sequences of any length, the raw data was length filtered to include only reads up to a length of 25 nt before input into miRExpress. Thus, raw reads were aligned to the miRNome provided by miRBase v21, yielding count tables of mature miRNAs and miRNA precursors for each sample. In total, 1913 mature miRNAs from miRBase v21 were discovered in the data.

discuss in text?

describe what miRExpress does?

3.6.3 Differential Expression Analysis - R/DESeq2

To determine the effect and dynamics of CNTF-mediated differentiation of LA-N-2 and LA-N-5, the expression state of each measured time point was compared to the respective control using the established R package *DESeq2*.¹³⁴ *DESeq2* determines differential expression (for gene i and sample j) in count-based data by application of a linear regression model to a negative binomial distribution based on a fitted mean μ_{ij} and a gene-specific dispersion value α_i . The mean is derived using a sample-specific »size factor«, s_j , and a parameter q_{ij} proportional to the expected true concentration of RNA fragments in the sample. The *DESeq2* differential expression pipeline is composed of the following commands:

- `estimateSizeFactors()` to estimate s_j
- `estimateDispersion()` to estimate α_i
- `nbinomWaldTest()` application of a generalised linear model to determine log-fold changes and statistics via the Wald test, using $\mu_{ij} = s_j q_{ij}$ and $\log_2(q_{ij}) = x_j \beta_i$

The Wald test, named after Abraham Wald,¹³⁵ is an approach to hypothesis testing that measures the distance between the tested unrestricted estimate and the null hypothesis, using the precision as a weighting factor. The larger the distance between tested values and the null, the more likely the measured values are »true«. RNA-seq data can be modelled using binomial distributions,¹³⁶ such as the Poisson distribution, and the difference between two Poisson means (e.g., »treated« vs »control«) can be tested by generalised linear models based on the distributions directly (Poisson regression), Fisher's exact test, or the likelihood ratio test. However, comparative analysis has shown that the Wald test on log-transformed data provides statistical power superior to these other methods,¹³⁷ particularly in lowly expressed fragments. The design formula for the linear regression was applied to LA-N-2 and LA-N-5 separately as a simple factor combination of condition and time point:

$$y \sim \text{condition_time}$$

To reduce the noise introduced by the high variance in low-count genes while preserving large, »real« differences, the authors propose the »shrinkage« of log-fold changes to avoid arbitrary low-cut filtering at a predefined expression (count) value. Multiple variants are available; for miRNA data, the adaptive algorithm »apeglm«¹³⁸ (adaptive t prior shrinkage estimator) yielded sensible results (see Fig. 3.7).

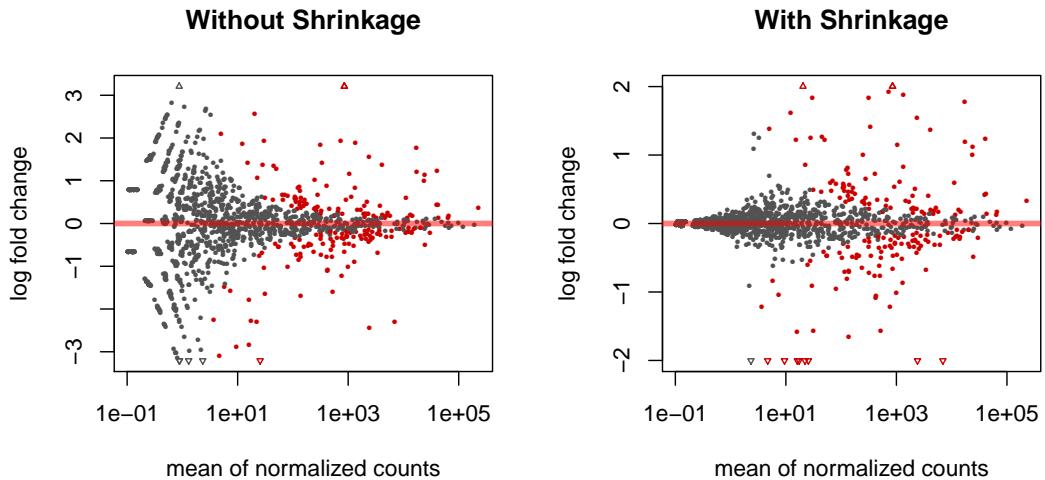


Figure 3.7: MD Plot Shrinkage Comparison. A mean-difference plot (MD Plot) is a plot of log-intensity ratios (differences, »M-values«) versus log-intensity averages (means, »A-values«); it is synonymous with »MA Plot«. The DESeq2 function `plotMD` shows the fold changes attributable to a given variable over the mean of normalised counts for all the samples in the data set. Points will be coloured red if the adjusted p value is less than 0.1. Points which fall out of the window are plotted as open triangles pointing either up or down. The left plot is generated from the standard linear model, the plot on the right is corrected by the »apeglm« algorithm¹³⁸ to reduce noise in the low-count fragments (data from LA-N-2 CNTF vs control on day 4).

3.6.4 MICRORNA DYNAMICS IN CNTF-MEDIATED CHOLINERGIC DIFFERENTIATION OF LA-N-2 AND LA-N-5

Differential expression analysis performed in this manner yielded 490 differentially expressed (DE) miRNAs across all groups, with characteristic distributions between cell lines and time points. The raw data and processed counts were deposited to NCBI Gene Expression Omnibus (GEO), accession GSE132951. An earlier sequencing experiment (deposited as GSE120520), which was similar in principle, but only comprised three biological replicates and only LA-N-2, reproduced 80% of DE miRNAs in the newer LA-N-2 samples. Considering the general reproducibility of RNA-seq and the lower replicate number, 80% is an excellent substantiation of the result. About 25% of miRNAs predicted in single-cell permutation targeting analysis (see Fig. 3.2 E) were found DE in LA-N-2 and LA-N-5 (Fig. 3.8 A) in all three groups, i.e., conserved, primate-specific, and TF-targeting miRNAs.

DIFFERENTIAL EXPRESSION IN BOTH CELL LINES

114 mature miRNAs were detected as DE in both cell lines, with some changes similar in both, while others were inverted (Fig. 3.8 B). In both cases, however, count-change values (see Box 2) correlated

highly between the two cell lines (similar: 76 miRNAs, Spearman's $\rho = 0.9066$, $p < 2.2\text{E-}16$; inverted: 38 miRNAs, $\rho = 0.9294$, $p < 2.2\text{E-}16$).

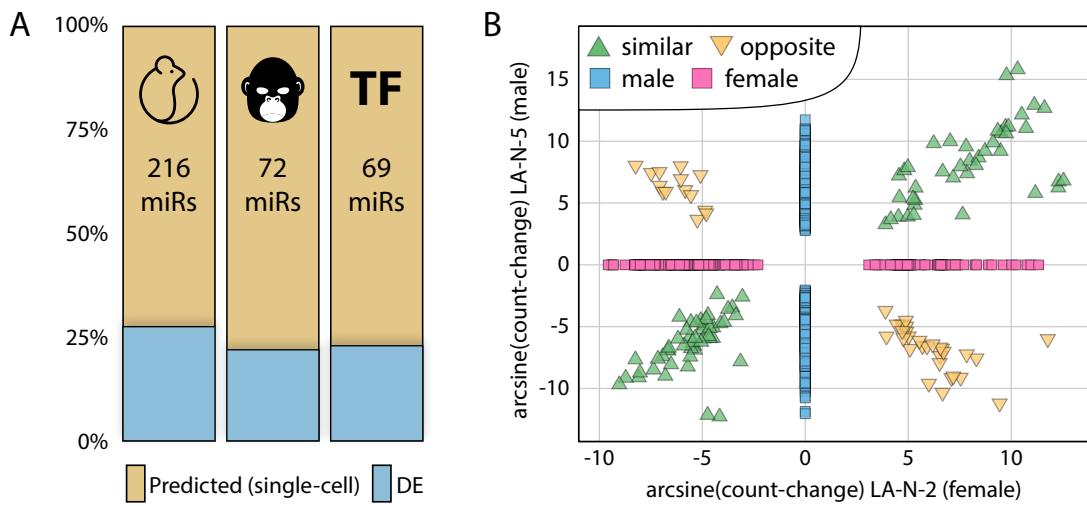


Figure 3.8: Differentially Expressed microRNAs in LA-N-2 and LA-N-5.

DIFFERENTIAL EXPRESSION ALONG THE TIMELINE

For consistency, from hereon out, time points 30 minutes and 60 minutes will be termed »early«, while 2 days and 4 days will be referred to as »late«. Differential expression was detected in all groups, lending credibility to the rapid changes in expression needed for a miRNA response of the »immediate-early« type. However, the response to long-term CNTF stimulation was larger in miRNA numbers as well as effect sizes (Fig. 3.9 A&B). Of all early perturbed miRNAs, only 3 and 13 miRNAs were exclusively perturbed immediate-early-like in LA-N-2 and LA-N-5, respectively; all others were still DE after 2 and/or 4 days. In LA-N-2, the late time points at 2 and, particularly, 4 days showed the

Box 2: The count-change metric

The frequently used log-fold change metric is not ideally suited for assessing the potential effect of expression changes for individual miRNAs because it does not reflect mean expression levels. To determine the absolute change in expression, the count-change metric was introduced, a combination of base mean expression and log-fold change, to weigh DE miRNAs against one another. The count-change is defined as follows:

$$CC = (BM \cdot 2^{LFC}) - BM$$

CC: count-change, BM: base mean expression, LFC: log-2-fold-change.

Importantly, by using the base mean expression, count-change correlates directly with sequencing depth. Generalisation, e.g. comparison between two individual experiments, is therefore not straightforward. A normalisation to raw reads would enhance comparability, however, other effects such as fragment distribution and quality aspects may also play a significant role.

greatest perturbation; in LA-N-5, the picture was more complex (Fig. 3.9 C&D). However, generally, there were large similarities as well as exclusivities between the time points 2 and 4 days in both cell lines. When comparing early and late time points between LA-N-2 and LA-N-5 directly, similarly complex patterns emerged (Fig. 3.9 E&F). Particularly at late time points (Fig. 3.9 F), every possible combination of overlap exists. 24 miRNAs were DE in all late conditions; 107 miRNAs were DE only in LA-N-2, and 269 miRNAs were DE only in LA-N-5.

DIFFERENTIAL EXPRESSION BETWEEN LA-N-2 AND LA-N-5

While there was considerable intersection in DE miRNAs between the cell lines, a substantial amount of miRNAs was only DE in one of the two lines. Generally, response to CNTF was higher in the male-originated LA-N-5 cells; however, there were also miRNAs found DE only in the female LA-N-2 (compare Fig. 3.9). Thus, not all of the differences in miRNA expression can be attributed to a higher sensitivity in LA-N-5. Similarly, LA-N-5 shows a »non-significant trend« toward higher count-change values (mean of absolute count-change across all DE time points, 20 907 versus 3066, Welch two-sample t test, $p = 0.08$).

The influence of genotype on the differentiating effect of CNTF was determined via a statistical interaction design in the *DESeq2* Wald test. Briefly, by including an interaction term in the linear regression formula, the effect of the condition (CNTF or control at each time point) between the two genotypes can be isolated:

$$y \sim \text{condition} + \text{genotype} + \text{condition : genotype}$$

Using the interaction term *condition : genotype*, miRNAs that reacted significantly different to CNTF stimulation in LA-N-5 compared to LA-N-2 were determined. Of note, the sexual dimorphism becomes more pronounced over the course of differentiation. While there is no significant difference between LA-N-2 and LA-N-5 at 30 minutes and only one miRNA DE at 60 minutes, numbers increase at 2 days and reach a maximum at 4 days, with significant overlap (Fig. 3.10 A). Although not all miRNAs found in this manner belong to the group of miRNAs with inverted expression between LA-N-2 and LA-N-5, several show significant differential regulation between the male and female cellular models (e.g., hsa-miR-615-3p, Fig. 3.10 B). To further examine the effect of genotype on the small RNA response to CNTF, the regular differential expression results (Section 3.6.4) were intersected with the interaction term for the late time points. This resulted in a complex pattern of intersecting miRNAs, in both cell lines (Fig. 3.10 C&D). Again, all possible overlaps between any two groups exist; 37 and 36 miRNAs are found in all four groups of LA-N-2 and LA-N-5, respectively. Among those, 16 mature miRNAs belong to all sets. All pertinent sets of miRNAs can be found in Appendix C.

something
special?

compare most
important tar-
gets early/late

target, GO for
which?

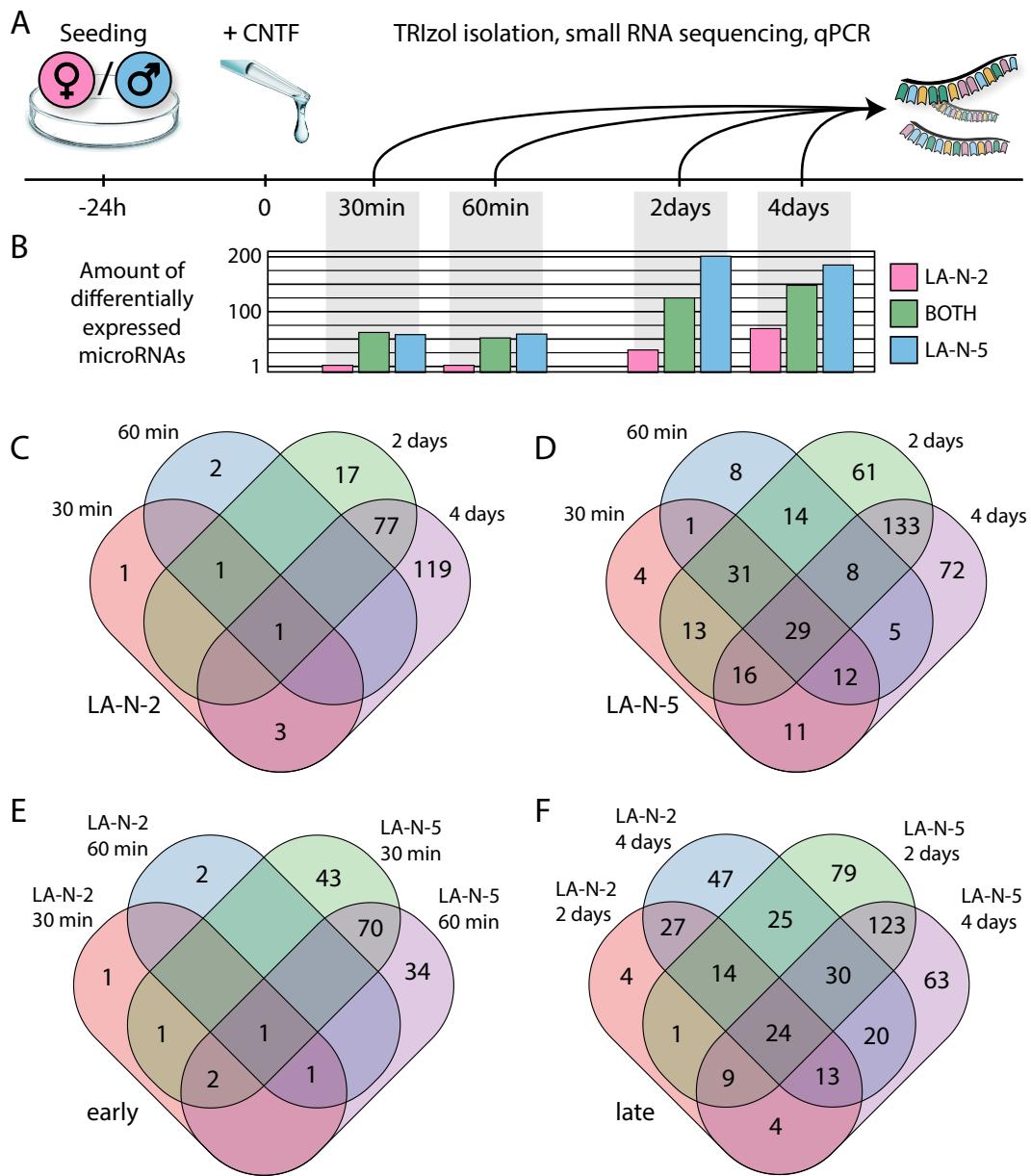


Figure 3.9: LA-N-2 / LA-N-5 Timeline and Differential Expression. A) Experimental timeline of CNTF differentiation. B) Bar plot of differentially expressed (DE) miRNAs per time point, divided by cell line where differential expression was measured (LA-N-2 only, LA-N-5 only, or both). C) Venn diagram of DE miRNAs in LA-N-2, divided by time point. Few early DE miRNAs, and continually more the longer differentiation lasts. D) Venn diagram of DE miRNAs in LA-N-5, divided by time point. Similar in pattern to C, but more pronounced in number. E) Intersection of early time points in LA-N-2 and LA-N-5. Despite the low differential expression in LA-N-2, there is overlap. F) Intersection of late time points in LA-N-2 and LA-N-5. Overlap is pronounced and complex, however, there are also cell line exclusive miRNAs.

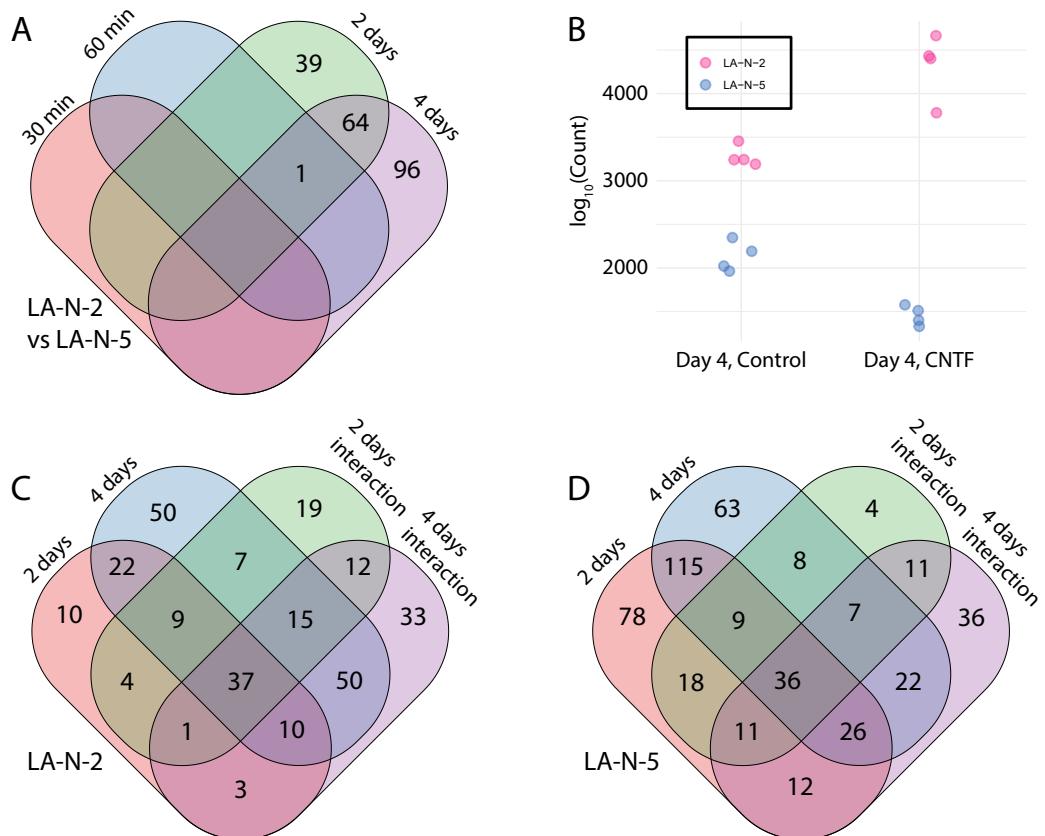


Figure 3.10: miRNAs DE between LA-N-2 and LA-N-5. Application of a design model formula which includes an interaction term enables display of the influence of the male or female genotype on differential miRNA expression. **A)** Venn diagram of miRNAs differentially expressed between LA-N-2 and LA-N-5 at all four time points. **B)** Counts plot of normalised raw expression values of hsa-miR-615-3p. Exemplary of a high influence of genotype on the differential expression caused by CNTF differentiation, hsa-miR-615-3p is more highly expressed in the female LA-N-2 and elevated after four days of CNTF-induced differentiation, while in the male LA-N-5, it is expressed slightly lower and suppressed upon differentiation. **C)** Venn diagram comparing late differential expression in LA-N-2 with late time points of differential expression between LA-N-2 and LA-N-5. All possible combinations exist, however, there are miRNAs affected by genotype that are not differentially expressed in the simple model. **D)** Venn diagram comparing late differential expression in LA-N-5 with late time points of differential expression between LA-N-2 and LA-N-5. Essentially similar to C), but with partly higher quantities of DE miRNAs.

3.6.5 MICRORNA FAMILY ENRICHMENT

To categorise and systematise the sexual dimorphism of CNTF differentiation of LA-N cells, statistically over-represented miRNA families in the differential expression datasets were determined. Of the 151 miRNA families listed in miRBase v21, members of 71 families are DE in LA-N-2 and LA-N-5. Enrichment of male, female, and ubiquitously DE miRNAs in these families was determined via hypergeometric gene set enrichment based on Fisher's exact test for each of the families. Five families were enriched in both male and female cells, and 12 families in only one of the two cell lines (Fig. 3.11 A, left side). The size range of enriched families was substantial, from small families with only 4 mature members to extensive families with dozens of mature miRNAs.

how many
DE miRs in
families?

GENE TARGETING OF ENRICHED FAMILIES

The targets of all individual miRNAs in the enriched families were determined via *miRNetDB* query. Of note, the amount of family members in any miRNA family did not correlate with the absolute amount of targets predicted (Fig. 3.11 A). Rather, the influence of individual miRNAs was the main factor determining the size of the gene target network. However, those families that were enriched in only one cell line presented with significantly smaller target sets than those that were found DE in both (mean targeted genes per miRNA 217 versus 378, Welch two-sample t test, $p = 0.001$). Relative to family size, 4 of the enriched families targeted less genes than all others: mir-10 ($p = 0.016$), mir-192 ($p = 0.042$), mir-379 ($p = 0.011$), and mir-515 ($p < 0.001$). Hypothetically, the spectrum of target amounts may correlate with the degree of functional specification of distinct miRNA families: on one end, broadly acting families such as let-7 with sex-independent function, on the other, families with a narrow target profile, such as mir-10, whose restricted function can associate with sex-specific effects.

which/how
many of the
pertinent mirs
above are in
families?

3.7 MICRORNA FAMILY GENE ONTOLOGY ENRICHMENT

A significant drawback of the recency of the discovery of regulatory small RNAs is the lack of comprehensive functional annotation. While protein coding genes are well annotated and neatly organised into an enormous amount of ontological categories (see Section 2.4.2), miRNAs have only been anecdotally associated with specific functions in the cell. Additionally, the functional roles of protein coding genes are much more limited than those of miRNAs; the number of potential functions of any miRNA correlates with the number of mRNA targets this miRNA has, and is also highly context-dependent (e.g. regarding cell type, cell state, disease). Thus, to systematically screen a large amount of miRNAs and families, I had to turn to an indirect approach: the GO analysis of targeted genes.

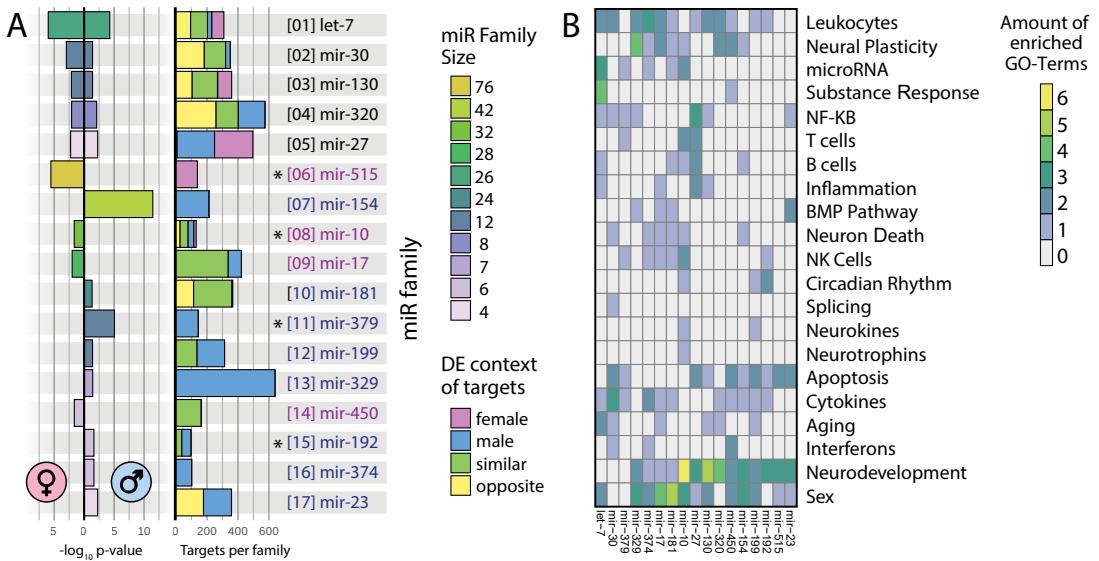


Figure 3.11: miRNA Families Enriched in Differential Expression and their Ontological Associations. 17 miRNA families were enriched significantly in the DE miRNAs following CNTF-mediated differentiation of LA-N-2 and LA-N-5 (Fisher's exact test, $p < 0.05$). **A, left side**) Bar plot of p-values of enriched families, ordered by family size; family size encoded by colour. **A, right side**) Stacked bar plot of the number of gene targets per family. Bars are divided by the DE pattern between LA-N-2 and LA-N-5 of each individual family member. DE context (encoded by colour) varies from detection in all categories (such as let-7 or mir-10) to detection only in one cell line (such as mir-515 or mir-154). Four families show significantly less target genes than all other families in relation to their size (denoted by asterisks). **B**) Gene Ontology enrichment analysis of gene targets of all enriched families, 737 distinct terms curated into CNS- or immunity-related categories. Families mir-10 and mir-199 show association with neurokines and circadian rhythm.

3.7.1 Creation of miRNA Family Gene Target Sets

GO analysis of the targets of a single miRNA is challenging, because the analysis requires a weighted scoring system of input genes. For single miRNAs, the options for scoring are limited to the aggregated targeting score or permutation p-values. Using families enables the introduction of a further scoring method: the aggregation of individual family members targeting the same gene. The reasoning behind this approach is to determine a general functional »area« of biological process that the miRNA family in question operates in. To account for the possibility of multiple areas being affected by a family, the test set of genes in any GO enrichment analysis should not be too small (i.e., rather the top 100 genes than the top 10).

Following this reasoning, the targets of all miRNAs in each family were determined via *miRNetDB* query. For each family, genes were ranked by their cumulative targeting score ρ from all family members. For gene i and number of miRNAs in family x , gene score ρ is calculated from individual miRNA→gene scores s :

$$\rho_i = \sum_{n=1}^x s_{ni}$$

3.7.2 GO Analysis of Target Sets

The gene target sets of individual miRNA families were ordered decreasingly by their cumulative score ρ and subjected to GO analysis via the R package *topGO*.¹³⁹ Briefly, *topGO* analysis extends the basic hypergeometric approach of GO enrichment analysis by de-correlating the DAG structure of GO annotation (see Section 2.4.2), allowing a weighted correction for the interdependency of neighbouring GO nodes. If a gene is found in both the parent node (more general) and the child node (more specific), the less specific parent node gene is weighted

less; in this way, the most specific node of each hierarchical branch can be found without confounding the result with less specific terms. While GO analysis always is subject to interpretation by the researcher, this weighted algorithm has been shown to reduce false positives while retaining a high true positive ratio.

topGO analysis was performed using the classic (i.e., Fisher's exact test) as well as weighted methods for comparison, however, to determine significance, the p-values calculated by the enhanced weighted algorithm were used. FDR was controlled at 5%. As recommended by the authors,¹³⁹ the ordered list of gene targets up to the 3000th position was used as a background for the analysis; the test set in each case was the top 10% of targeted genes.

3.7.3 LARGE SCALE GO TERM CURATION

The GO analysis performed in this manner for all 17 enriched families resulted in a list of 737 distinct GO terms related to any of the families. To generate an overview of functional implications of the individual families, the GO terms were filtered and aggregated manually. Terms not relating to CNS- or immune-function were removed, and the remaining terms were sorted into one of 21 categories (Fig. 3.11 B). Generally, the families associated with neurodevelopment and neural plasticity, diverse immune functions, cell cycle control, and sex. More general categories were found in most families, while more specific functions showed a sparser distribution.

Only two families associate significantly with neurokine-related function, mir-10 and mir-199. Both are involved in neurodevelopment- and sex-related function, and show the very specific association with circadian rhythm. Family mir-10 additionally is implicated in control of neurotrophin-related mechanisms, and in several blood-borne immune cells, such as T-, B-, and NK-cells.

3.8 WHOLE GENOME miRNA→GENE NETWORK GENERATION

A common approach to complex network relationships is physical modelling. A complex graph (with directed and weighted edges) can be coerced to self-organise by application of a force-directed layout. In this process (also known as spatialisation), the network, defined only by its nodes and edges, is transformed into a map, usually in two dimensions. An important prerequisite is the scale-free topology of the network, a structure that transcriptional connectomes usually present with(cite). A force-directed layout transforms a network by simulating a gravitational system, or a system of magnetic nodes connected by springs, in which the nodes repel each other, but edges between two nodes pull them towards each other. By manipulation of multiple physical attributes of the model, a mapped representation of the network's organisation can be produced. As a result, nodes (i.e., genes, TFs, and miRNAs) with close interaction are mapped in close proximity, while nodes with low interaction are far apart. Similarly, nodes with pivotal function in the network (»hubs«) gravitate towards the centre of the map, while »less important« nodes are shifted towards the fringes.

The network comprising all DE members of the 17 enriched miRNA families and 12 495 targeted genes as determined via *miRNetDB* query was subjected to force-directed mapping using the Java-based software Gephi

0.9 and its primary force-directed algorithm, ForceAtlas2.¹⁴⁰ Gephi, and ForceAtlas2, are designed to generally handle graphs with up to 10 000 unique relationships; however, the standard *miRNetDB* query resulted in a network with ~160 000 edges. To reach a computationally manageable number of relationships, the score threshold was raised to a minimum of 7, which resulted in a network of 46 937 unique edges. The resulting network was exported as a vector graph and manually edited in Adobe Illustrator to further enhance its readability (Fig. 3.12).

The resulting transcriptional connectome map illustrates the functional compartmentalisation of miRNA→gene interactions. miRNAs of distinct families are frequently found in close proximity to one another, most often forming one or two clusters. In the case of two clusters forming, the clusters are usually representative of the two complementary strands of the pre-miRNA(s), since 3' and 5' variants of any pre-miRNA usually possess fundamentally different seed sequences, and thus, targets. The let-7 family is distinguished by its removal from the bulk of other interactions, possibly representing a particularly specialised set of functions, at least for the 5' variants of the bottom cluster. Families with predominant differential expression in one of the two cell lines (sexes) inhabit different sides of the main graph and show little intermingling, pointing towards sexually dimorphic gene target distribution. The two neurokine-associated families, mir-10 and mir-199, are located near the centre of the graph, in two strand specific clusters (»[08a]&[12a]« and »[08b]&[12b]«.

To gather more detailed information than grouping of miRNAs with similar function, such as direct miRNA→gene interaction, the size of the studied networks must be reduced. For each family affected by CNTF-differentiation, a single graph was created, laid out by application of ForceAtlas2, and analysed for critical nodes. The distinct families and their gene targets yield immensely diverse graph layouts, that here cannot be described in their entirety. However, the entire collection of graphs in interactive visual form is accessible at <https://slobentanzer.github.io/cholinergic-neurokine>. Due to an elevated interest, the cholinergic/neurokine miRNA interface and the families mir-10 and mir-199 will be described in more detail, and in conjunction with sex-specific perturbations in neurologic diseases.

3.9 APPLICATION TO SCHIZOPHRENIA AND BIPOLAR DISORDER

A comprehensive structural analysis of perturbations on a genome scale is hardly possible without heavy truncation of results or dimensionality reduction methods. Truncation is commonly performed by ranking perturbations by their p-values in ascending order and only regarding the highest ranked entries, which often amounts to less than ten individual transcripts. On the other hand, commonly used dimensionality reduction techniques include principal component analysis (PCA), t-distributed stochastic neighbour embedding (t-SNE), and clustering/stratification approaches. While truncation enables human-readable presentation of results, in principle it does not lend itself to complex polygenic events such as neurologic disease. Common dimensionality reduction techniques are useful in providing structural overview of a high-dimensional dataset, but give little insight into causal

evolutionary?
discussion?

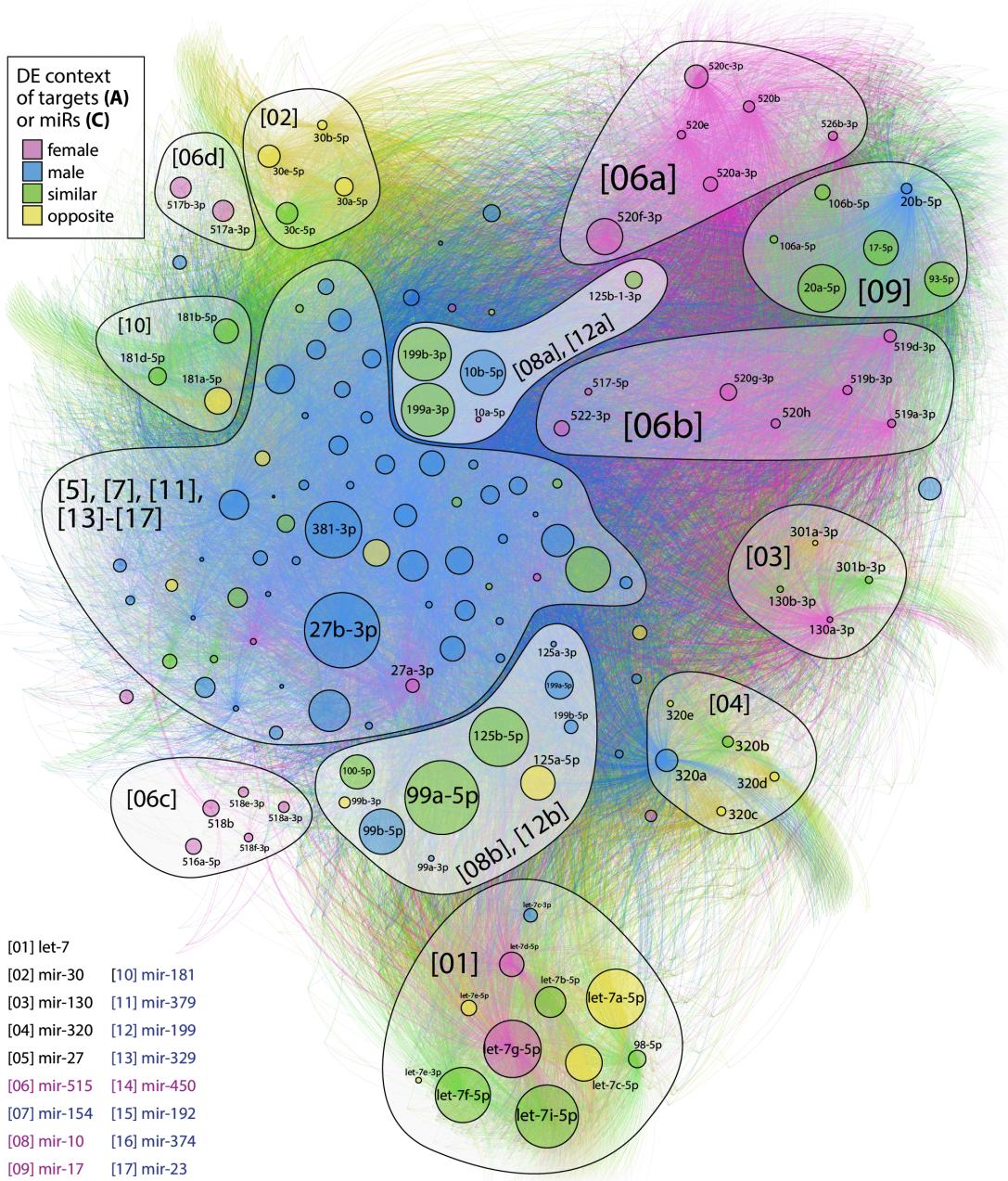


Figure 3.12: Full Connectome of LA-N-2 and LA-N-5 Differentially Expressed miRNA Families. The network of miRNA families and their 12 495 targeted genes self-organises into a connectome map with 46 937 unique edges. miRNA node size scaled by absolute count-change, nodes coloured by DE context. Numbers in brackets denote miRNA families, gene nodes have minimal size. By application of a force-directed layout, the miRNA families visibly self-segregate into clusters. The let-7 family, male-biased and female-biased clusters take up major parts of the network. Families mir-10 and mir-199, with neurokine association, form two mixed, sexually dimorphic clusters near the centre of the map (lighter shade).

source	accession	publication	technology	subjects/samples	disease
NCBI GEO	GSE35978	Chen <i>et al.</i> ¹⁴¹	microarray	150/312	SCZ & BD
NCBI GEO	GSE12649	Iwamoto <i>et al.</i> ¹⁴²	microarray	102/102	SCZ & BD
NCBI GEO	GSE53987	Lanz <i>et al.</i> ¹⁴³	microarray	76/205	SCZ & BD
NCBI GEO	GSE17612	Maycox <i>et al.</i> ¹⁴⁴	microarray	51/51	SCZ
NCBI GEO	GSE21138	Narayan <i>et al.</i> ¹⁴⁵	microarray	30/59	SCZ
NCBI GEO	GSE5392	Ryan <i>et al.</i> ¹⁴⁶	microarray	82/82	BD
NCBI GEO	GSE80655	Ramaker <i>et al.</i> ¹⁴⁷	RNA-seq	96/281	SCZ & BD
NCBI GEO	GSE106589	Hoffman <i>et al.</i> ¹⁴⁸	RNA-seq	94/94	SCZ
NCBI GEO	GSE68559	Webb <i>et al.</i> ¹⁴⁹	RNA-seq	10/98	NA
NCBI GEO	GSE96659	Fontenot <i>et al.</i> ¹⁵⁰	RNA-seq	5/209	NA
NCBI GEO	GSE45642	Li <i>et al.</i> ¹⁵¹	RNA-seq	86/670	NA
Synapse	CMC	Gulyás-Kovács <i>et al.</i> ¹⁵²	RNA-seq	579/579	SCZ & BD

Table 3.2: Data Sources for Microarray and RNA-seq Analyses.

relationships of single entities. We thus aimed to find an alternative approach to dimensionality reduction which conserves the internal relationships inherent to the data, and which profits from the network organisation of our input data.

For the application of *miRNetDB* data to real-world problems, suitable psychiatric and neurologic disease datasets were sought in the common repositories ArrayExpress, NCBI GEO, and Synapse. Among the datasets with agreeable quality, SCZ and BD were the only diseases with sample amounts that allowed a statistically valid analysis of sexual dimorphisms. While many neurologic disease studies are simply limited in their number of subjects, autism presented a different issue: the majority of donors were males (more than 90%). Direct analysis of miRNA expression patterns was not possible, because very few studies study miRNAs directly, yet. Thus, studies on mRNA were substituted to infer on miRNA dynamics.

3.9.1 Analysed Datasets

Twelve datasets including 1361 subjects were downloaded from their repositories (Table 3.2). Data of DLPFC RNA-seq of 579 SCZ patients and controls was obtained from the Common Mind Consortium (<http://www.synapse.org/CMC>). To address the diverse origins and technological aspects of data, care was taken to appropriately unify and normalise the data. The data preparation and meta analyses were performed essentially as described by Gandal and colleagues.³¹ Samples of brain regions not consistent with the research question (e.g., cerebellum), or from patients with diseases other than SCZ or BD, were removed from datasets on a case-by-case basis. RNA-seq datasets were used to individually confirm the perturbations found in the meta-analysis of microarray studies.

3.9.2 Microarray Quality Control and Data Preparation

Read-In and Normalisation

Illumina datasets were read, \log_2 -transformed, and quantile-normalised using R/lumi.¹⁵³ Affymetrix datasets were read and RMA-normalised (\log_2 -transformed, background corrected, quantile-normalised) using R/affy.¹⁵⁴

Affymetrix data were additionally corrected for 3'/5' bias using the *AffyRNAddeg()* function (not available for other chip manufacturers). All available biological (e.g., sex, age) and technical (e.g., batch, RIN, post-mortem interval) covariates were collected and used for the analysis. Individual correlations of case-control status S with any covariate C were assessed using a linear model (R/lm) with formula $C \sim S$; statistical significance was determined via ANOVA (R/anova). If necessary, case-control samples were balanced to eliminate significant covariate correlations with case-control status (all $p > 0.05$).

Outliers

Outlier removal was performed using the method proposed by Oldham, Langfelder & Horvath.¹⁵⁵ Briefly, the (dis-)similarity matrix of samples is transformed into a signed, weighted correlation network. Network adjacency (α) of samples (nodes) S_i and S_j is defined as:

$$\alpha_{ij} = \left(\frac{\text{cor}(S_i, S_j) + 1}{2} \right)^2$$

As such, the connectivity between samples can be measured by the standardised connectivity (Z.K), which describes the strength of correlation between any given node and all other nodes in the network. As proposed by Oldham *et al.*, outliers were removed if their Z-score was below the threshold of Z.K = -2.

Annotation

To enable comparison between datasets of diverse technical origin, probes were annotated using ENSEMBL gene identifiers using R/biomaRt.¹⁵⁶ To maintain comparability with the analysis by Gandal *et al.*,³¹ the same version of ENSEMBL DB (v75, Feb 2014) was used. Probes were collapsed onto single genes using the *collapseRows()* function of R/WGCNA,¹⁵⁷ using the maximum mean signal across all probes per gene. Of note, information loss occurred by multiple collapsing of probes and integration of datasets, which can only be performed using the genes common to all datasets (i.e., represented by microarray probes). The final gene set encompassed 12 391 individual genes, with several notable cholinergic/neurokine exceptions (CHRNA7, CHRM1, LHX8, CHKB, PRIMA1, CNTF). Missing genes result from annotation deficits between different probe sets, cannot be comprehensively manually controlled on a genome scale, and cannot be re-introduced at this stage.

3.9.3 Differential Expression Meta-Analysis

The individual experimental datasets were each corrected for covariate influences by multiple regression based on all available biological and technical covariates. Briefly, the linear regression model was solved using matrix algebra operations. In matrix form, a linear regression model of observations Y (i.e., gene expression levels), independent variables X (i.e., covariates), coefficients β , and error terms ε can be described as:

$$Y = X\beta + \varepsilon$$

As a consequence, the residual sums of squares can be expressed as the cross product:

$$\text{RSS} = (Y - X\beta)^T(Y - X\beta)$$

Then, the coefficients $\hat{\beta}$ can be estimated by solving the derivative:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Coefficients were estimated for all relevant technical and biological covariates (e.g. post-mortem interval, RIN, sex, age) and used to regress covariate influence on gene expression levels:

$$Y_{new} = Y - (X\hat{\beta})^T$$

After covariate regression, differential expression was calculated across all datasets for each disease group using a linear mixed model with a fixed effect for each study and case-control status (»group«), and a random effect for each individual subject. Computation was performed in R, using R/nlme,¹⁵⁸ with parameters

$$fixed = \sim group + study \text{ and } random = \sim 1 | subject$$

This yielded an array of log-fold changes between cases and controls for each gene and disease. To determine statistical significance, 10 000 permutations of the mixed-model regression were performed for each use case, randomly assigning case-control status. The resulting null distributions were used to determine FDR, with threshold for significance at 0.05.

Sex-Specific Meta-Analysis

Samples of all datasets were split between males and females (cases as well as controls), and individually subjected to the same procedure as the sex-independent data: covariate regression, differential expression via a linear mixed model, and estimation of statistical significance via permutation testing.

Transcriptome Correlation

Correlation of disease transcriptomes was performed by using Spearman's rank correlation coefficient. Spearman's ρ was determined between SCZ and BD sex-independently as well as separately in males and females.

Most Diverging Genes

Genes were ranked by their divergence between any two compared datasets, sex-independent data of SCZ and BD, and any meaningful combination of sex-dependent data in SCZ, BD, males, and females. The divergence δ of any gene G between datasets i and j was defined as (logFC: log-fold change):

$$\delta = \log FC(G)_i - \log FC(G)_j$$

Where positive values of δ indicate a positive bias of G towards dataset i .

3.9.4 SEXUAL DIMORPHISM IN SCHIZOPHRENIA AND BIPOLAR DISORDER

Sex-independent correlation replicated the finding of Gandal *et al.*,³¹ with Spearman's $\rho = 0.7100$ ($p < 0.001$). However, diverging from the established annotation of ENSEMBL v75 to later versions of the database significantly altered the correlation coefficient, leading to lower correlation in all tested

cases. Comparing the sex-independent data with only male or female subjects, those also show lower general correlation between SCZ and BD: in females, correlation was $\rho = 0.6150$ ($p < 0.001$), in males, $\rho = 0.5783$ ($p < 0.001$). While it is possible that these variations are caused by structural properties of the data unrelated to sexual dimorphism, such as the loss of power due to the reduction in size, the consistently lower correlation in sex-specific subsets also might indicate an averaging effect between male and female patients, leading to a higher correlation in spite of significant sexual dimorphism.

To address the potential differences between male and female brain transcriptomes, which might reflect the observed clinical dimorphism, we subjected the 100 most-diverging genes between any two datasets to GO enrichment analysis (Fig. 3.13, from Lobentanzer *et al.*⁶) in hopes of identifying the most discriminating molecular pathways between SCZ and BD, and afflicted males and females. Sex-independently, the most-diverging pathways between SCZ and BD principally involved mechanisms of inflammation and immunity (e.g., “acute inflammatory response,” $p = 0.003$; “cellular response to cytokine stimulus,” $p = 0.01$).

DIFFERENCES IN SEXUAL DIMORPHISM BETWEEN SCZ AND BD

Computation of diverging pathways between males and females in each disease indicated a larger divergence between sexes in SCZ than in BD. SCZ-biased genes of males and females showed no overlapping GO terms (Fig. 3.13 A), but BD-biased genes of males and females showed large GO term overlap, particularly in inflammatory components (Fig. 3.13 B). Notably, specific components of neurokine signalling were elevated in both males (IL-6, $p = 0.007$) and females (JAK/STAT, $p = 0.01$) with BD.

OVERLAP OF MALE-BIASED GENES BETWEEN SCZ AND BD

Shared transcriptional properties of SCZ and BD were identifiable only in male diverging genes. While female-biased SCZ genes showed no implications in CNS processes, male-biased SCZ and BD genes overlapped in functions concerning inflammation and immunity (Fig. 3.13 C). Female-biased BD genes were associated with CNS function and development (Fig. 3.13 D).

SPECIFICITY OF ONTOLOGICAL TERMS

When comparing different areas of biological function, such as neurotransmission, immunity, and inflammation, the different areas notably diverged in the specificity of identified terms. Considering the functionality of the applied method (*topGO*, see Section 3.7.2) to find the most specific node in any branch of the DAG tree while disregarding its (less specific) parent nodes, this might indicate a difference in the gravity of perturbation in the different systems. For instance, the GO terms indicating neurotransmission as affected system were much less specific than those indicating immunity-

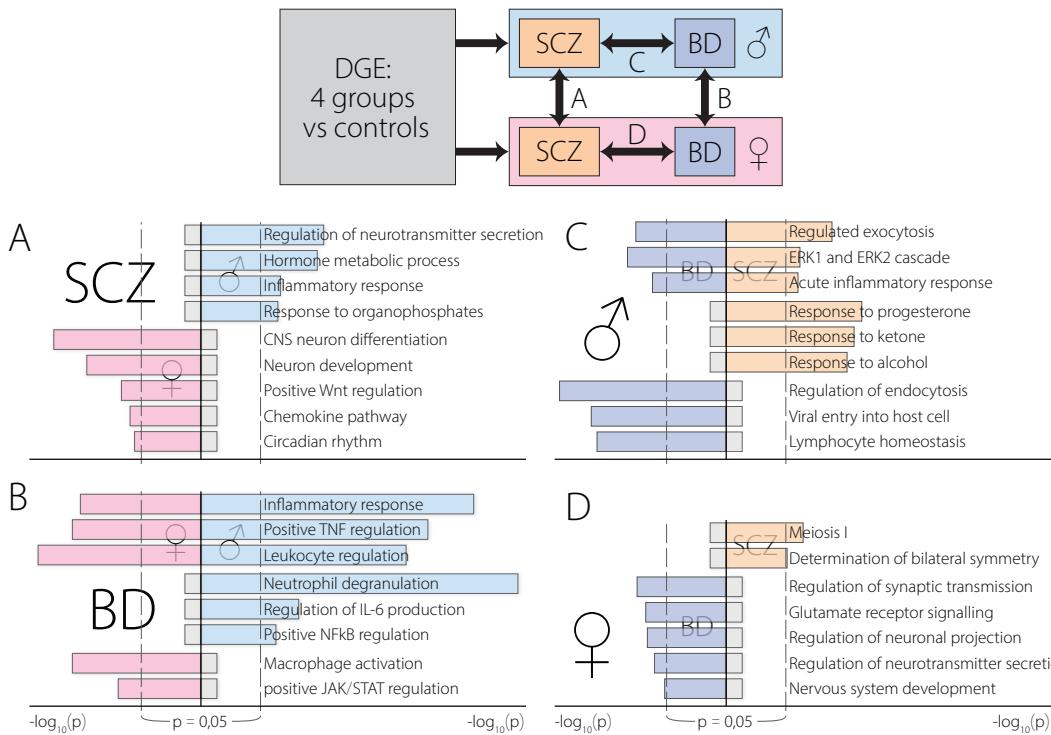


Figure 3.13: GO Enrichment of Diverging Genes. Results of differential gene expression were dually compared: SCZ versus BD and male versus female (indicated by colours). GO enrichment of the top 100 distinguishing genes in one dimension was compared with the other for each pair of combinations. **A)** SCZ-biased genes diverge between males and females. **B)** BD-biased genes share immunological ontology in both males and females. **C)** Male-biased genes share immunological ontology in BD and SCZ. **D)** Female-biased genes diverge between SCZ and BD.

related processes. While significant neurotransmission-related terms failed to implicate specific neuron types or neurotransmitters (e.g., GO:0021953, »CNS neuron differentiation«; GO:0046928, »Regulation of neurotransmitter secretion«), immunity-related terms were very specific towards regulatory subsystems, and regularly implicated neurokinin mechanisms (e.g., GO:0032675, »Regulation of interleukin-6 production«; GO:0046427, »Positive regulation of JAK/STAT cascade«).

3.9.5 COMBINATION OF DISEASE DATA AND CELL CULTURE

To implement the proposed complexity reduction technique, we applied a reductionist approach to the comprehensive network generated from perturbed miRNA families and their targeted genes (Fig. 3.12), based on the unbiased analysis of sexual dimorphism in SCZ and BD, which implicated processes of neuronal, immunological, and circadian origin (Figure 3.13). To merge these results with the implications of cholinergic cell culture, we added genes implicated in neurokinin signaling and circadian rhythm to the list of cholinergic genes (see Box 1). Returning to the collection of web-available patient data, we subjected this limited set of 76 genes and their 18 neuronal TFs to differential expression analysis.

The comprehensive network was then filtered by multiple consecutive steps. (I) Permutation

analysis of comprehensive miRNA targeting data specific for genes expressed in cholinergic neurons (Fig. 3.2) yielded a list of miRNA candidates that shows overlap with (II) miRNAs DE in our two models of neurokine-induced cholinergic differentiation (Fig. 3.8 A). (III) We included only families of miRNAs we found to be enriched in differential expression (Fig. 3.11). Sixty-nine miRNAs from 12 families passed this filtering process and were consecutively assembled in a force-directed network with the 94 genes of the previously compiled list. As a »spike-in«, we added miR-132-3p (DE in LA-N-5 cells), a miRNA which controls cholinergic processes^{159,160} and is known for its function in neurons¹⁶¹ and immunity¹⁶² and its perturbation in disease.¹⁶³ The resulting network (Fig. 3.14 A, from Lobentanzer *et al.*⁶) shows high structural homology to the comprehensive network shown in Figure 3.12. The miRNA families in this reduced network show spatial organisation similar to the comprehensive network.

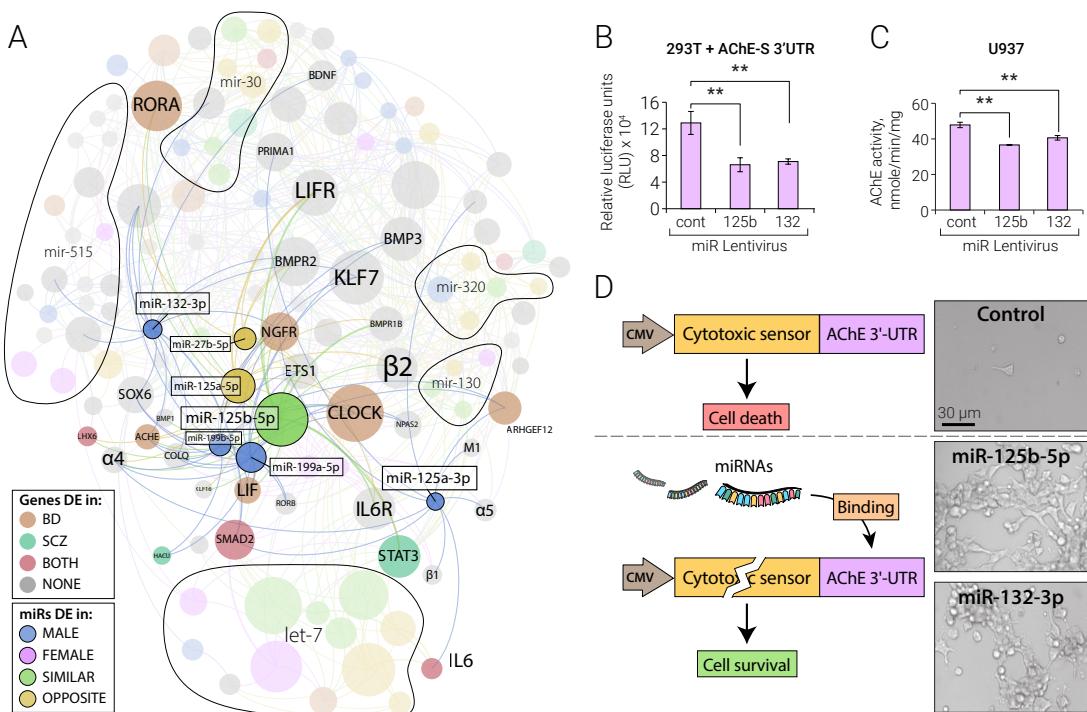


Figure 3.14: The cholinergic/neurokine interface. **A)** The miRNA families mir-10 and mir-199 pose a sexually dimorphic interface of cholinergic, neurokine, and circadian regulation by targeting nicotinic/muscarinic (e.g., a4b2 and M1) and neurokine receptors, transcriptional regulators of cholinergic differentiation (LHX and STAT) and circadian rhythm (CLOCK and RORA), the AChE and the AChE linker proteins PRIMA1/COLQ, and high-affinity choline uptake (HACU). Members of mir-10/199 families, spike-in miR-132-3p, and their targeted genes are shown in colour, and other miRNA families that passed the multiple filtering are indicated as areas. miRNA node size corresponds to count-change and gene node size to connectivity; colour and thicker edges indicate the DE context and experimentally validated connections. **B-D** Validation experiments of AChE targeting by miR-125b-5p, with miR-132-3p as a positive control. **B)** Lentiviral expression of miR-132 and miR-125b suppresses luciferase fused to the 3' UTR of AChE in HEK293T cells. Error bars indicate SE. **C)** Lentiviral expression of miR-132 and miR-125b suppresses the endogenous AChE hydrolytic activity of U937 cells with similar efficacy. Error bars indicate SE. **D)** Life/death assay of stably transfected HEK293T cells carrying the AChE 3' UTR fused to a cytotoxic sensor and co-transfected with miR-125b-5p, miR-132-3p, or control plasmids. Cells survive in case of binding of miR-132-3p and miR-125-5p to the 3' UTR.

In agreement with their localisation in the comprehensive network, miRNA families mir-10 and mir-199 inhabit a central role in the resulting interactome. Most-targeted genes in this network (as

indicated by their size) are the circadian regulators CLOCK and RORA. While CLOCK is located centrally, next to mir-10/199 miRNAs, RORA shows closeness to the mir-30/515 families. Generally, genes with larger cellular influence, such as transcription factors (STAT3, CLOCK) or TGF- β ligands (BMP family genes) are frequently targeted by miRNAs, while more specific transcripts, such

as the cholinergic receptor genes or neurokines, are targeted more selectively.

More so than the spiked-in miR-132-3p, mir-10/199 miRNAs target cholinergic genes, for instance, the neuronal nicotinic $\alpha 4\beta 2$ and muscarinic M1 receptors (mir-125) and HACU (miR-199). In addition, they target neurokine genes, such as the transmembrane neurokine receptor LIFR or STAT3, and circadian regulators (e.g., CLOCK and RORA). The two families react highly sexually dimorphic to CNTF-mediated differentiation; some are detected as DE only in one cell line, others exhibit inverted changes between cell lines. The 3p-variant of miR-125a distinguishes itself from the bulk of mir-10/199 miRNAs by exclusively targeting M1, $\alpha 5$ and $\beta 1$ receptors, and IL-6, and thus is slightly removed from the centre of the network.

The miRNA with most targets in this reduced interactome is miR-125b-5p, also displaying most experimentally validated interactions with neurokine genes (miRTarBase accessions: IL-6, MIRT-022105; IL-6R, MIRT006844; JAK2, MIRT734987; LIF, MIRT001037; LIFR, MIRT732494; STAT3, MIRT005006). miR-125b-5p also is the most perturbed miRNA (in this interactome) upon CNTF-mediated differentiation (highest count-change), and the only member of mir-10/mir-199 to be changed in similar direction in both cell lines (up-regulated). miR-125b-5p also targets multiple other inflammation-related genes (e.g., TNF, MIRT733472; IRF4, MIRT004534) and 5-lipoxygenase, which can influence inflammatory processes via production of eicosanoids.¹⁶⁴ miR-125b-5p has been directly associated with cytokine-mediated inflammation, as its over-expression increased the expression of TNF- α , IL-1 β , and IL-6, and markedly decreased I κ B- α .¹⁶⁵

A notable intersection of spike-in miR-132-3p and miR-125b-5p is the ACHE, an interaction which had not been validated for miR-125b-5p, but is known for miR-132-3p.^{162,159} Using miR-132-3p as a positive control, we performed ACHE-mRNA binding assays in validation of the predicted targeting by miR-125b-5p.

3.9.6 miR-125b-5p Acetylcholinesterase Targeting Assays

We performed three independent cell culture assays to confirm ACHE mRNA targeting by hsa-miR-125b-5p: luciferase suppression, AChE protein activity, and a cell death assay with a cytotoxic sensor. The 3' UTR of human ACHE mRNA¹⁶⁶ was cloned into the microRNA Target Selection System plasmid (System Biosciences, CA, USA) multiple cloning site, using EcoRI and NotI restriction enzymes (New England Biolabs). All plasmids were verified by DNA sequencing. For luciferase assays, HEK293T cells were transfected with miRNA Target Selection-AChE 3' UTR, and selected in the presence of Puromycin for 3 weeks. Stably transfected HEK293T (293T-AChE 3' UTR) cells were grown on 12-well plates and infected with lentiviruses expressing miR-125b-5p, miR-132-3p or a negative control sequence. After 48 hours incubation, cells were analysed using the Dual Luciferase Assay kit (Promega, WI USA) and Luciferase activity was measured using an Envision luminescent plate reader (Perkin-

more details?

Describe
other mem-
bers?

Elmer, Waltham, MA), essentially as previously described by Hanin *et al.*¹⁶⁷ For each reporter construct, renilla luciferase activity was normalized according to that of the firefly. Normalised activity after infection with miR-132-3p or miR-125b-5p was expressed as relative to that obtained after infection with the same plasmid with miRNA negative control. To show effects of changes in this miRNA's levels on real-life protein activities, we performed an AChE hydrolytic activity assay following infection of human monocyte-like U937 cells with hsa-miR-125b-5p, miR-132-3p or a negative control lentiviral vector. AChE hydrolytic activity levels were assessed by kinetic measurements of the hydrolysis rates of 1 mM acetylthiocholine (ATCh, Sigma) at room temperature, following 20 min incubation with and without 50 µM tetraisopropyl pyrophosphoramido (iso-OMPA, Sigma), a specific inhibitor of butyrylcholinesterase, to selectively assay for AChE-specific or total cholinesterase activity. For the life/death assay, stably transfected HEK293T cells were infected with lentiviruses expressing miR-125b-5p, miR-132-3p or a negative control sequence. 72 hours post-infection, a cytotoxic reporter fused to AChE 3' -UTR was added to the media and cells were kept for an additional 5 days to assess their viability. For all cell culture assays, statistical significance was determined using ANOVA with correction for multiple testing. Each sample was assayed in at least 3 biological replicates, and in all cases, hsa-miR-132-3p served as a positive control.

3.9.7 HSA-MIR-125B-5P TARGETS ACETYLCHOLINESTERASE

In all tested conditions, miR-125b-5p suppressed *ACHE* mRNA with equal potency as the positive control miR-132-3p (Fig. 3.14 B-D). Towards mRNA expression (luciferase) and functionality (cytotoxic sensor) as well as on protein level (AChE activity), miR-125b-5p demonstrated its interaction with *ACHE* mRNA 3' UTR. Luciferase units after miR-125b-5p transfection were approximately halved, indicating significant transcript degradation of the *ACHE* 3'UTR.

3.9.8 CHOLINERGIC/NEUROKINE MECHANISMS IN WEB-AVAILABLE RNA SEQUENCING EXPERIMENTS

To include recent developments in methodology, we analysed several recent RNA-seq studies addressing related questions. In a study of post-mortem brain transcriptome profiling of psychiatric disorders,¹⁴⁷ we found a down-regulation of IL-6, LIF, and several cholinergic receptors (M2, M4, α4, β2, α7), with sex-specific differences (males had significantly higher levels of neurokines than females). These changes were visible only in SCZ patients, not in BD or major depressive disorder. In a study of induced pluripotent stem cells (iPSCs) of SCZ patients and controls that were induced to show a neuronal phenotype,¹⁴⁸ we found an up-regulation of *CHAT* in SCZ-derived iPSCs, and a down-regulation of IL6R and the nicotinic α6 subunit. In this study, SCZ males showed a higher expression of the *SLC18A3* and lower expression of nicotinic subunits α 2, 7, and 9, and β3. In a study of differentiated human neuronal progenitor cells,¹⁵⁰ a knockdown of the circadian transcriptional controller CLOCK resulted in up-regulation of LIF and simultaneous down-regulation of neurokine transmembrane receptors LIFR and IL6ST, accompanied by slight bi-directional changes in several cholinergic receptors.

I know words. I have the best words.

Donald Trump

4

Dynamics Between Small and Large RNA in the Blood of Stroke Victims

Stroke is a dramatic incision into bodily homeostasis and affects a multitude of organ functions, first and foremost the brain. The immediate actions upon stroke are focused in preserving as much functional tissue as possible, so as to alleviate the cognitive damages resulting from neuron death. After this initial period of few hours, longer-lasting reactions determine the health and recovery of the patient. Many of these later events are related to immunity. The greatest danger to the patient after survival of the initial period are infections, such as pneumonia, usually between one and two weeks after the infarction. Pneumonia is often facilitated by aspiration of liquids or solids when the swallowing mechanism is impaired as a consequence of the cerebral damage. However, as introduced in Section 1.2.5, stroke-related immunodepression can play a role in post-stroke survival, and has been shown to have an impact on the transcriptome of blood-borne immune cells, at least for protein coding genes. The role of short RNA transcripts, and particularly of transfer RNA fragments, is much less clear.

4.1 RNA SEQUENCING, DIFFERENTIAL EXPRESSION, AND DESCRIPTIVE METHODS

We thus opted to analyse the blood of stroke victims taken upon hospitalisation, and screen it for changes in small and large RNA expression.

page break?

4.1.1 The PREDICT Cohort

The patient collective for the present study was recruited from a prospective, international, multi-center study with 11 study sites in Germany and Spain, led and approved by the neurologic department of Charité Berlin (www.clinicaltrials.gov, NCT01079728).¹⁶⁸ The study, called PREDICT, screened 484 stroke patients for clinical attributes and conventional biomarkers, with daily measurements in the first 5 days after stroke, and a three months follow-up. From these patients, a representative cohort of 49 patients were selected for blood small RNA sequencing.

4.1.2 Clinical Parameters Collected in the PREDICT Study

Stroke patients were assessed daily for the duration of hospitalisation, at least until four days after admission. Blood-based biomarkers that were measured at least once during this period include: monocyte human leukocyte antigen isotype DR (HLA-DR), interleukins IL-6, IL-8 and IL-10, IL-10 levels after 24h *in vitro* stimulation with lipopolysaccharide, lipopolysaccharide binding protein (LBP), mannan-binding lectin (MBL), and TNF- α . Also recorded were the time between admission and the collection of the blood sample, and the modified Rankin Scale (mRS). This scale is a rough categorisation of the severity of stroke, with 0 referring to no symptoms, and 6 signifying death. Scores 1-2 describe slight neurological deficits, 3 requires frequent help because of medium level deficits, 4 requires constant assistance with daily tasks, and 5 requires stationary care.

4.1.3 Sample Collection, RNA Isolation, and Sequencing

Blood was collected into RNA stabilising tubes (Tempus Blood RNA tubes, Applied Biosystems) on each day of hospitalisation, and we subjected blood samples collected on the second day to small and large RNA-sequencing. While choosing samples for sequencing we only considered samples from patients with modified Rankin Scale (mRS) values of 3 and below at discharge from the hospital, to exclude very severe cases of stroke, leaving n=240 relevant cases. The time from stroke occurrence to blood withdrawal varied between 0.94 to 2.63 days, with an average of 1.98 days. Blood samples from age- and ethnicity-matched healthy controls were obtained at matched circadian time from donors with ethical approvals from institutional review boards (ZenBio, North Carolina, USA).

RNA was extracted from 3 ml of whole blood of all 484 PREDICT patients using the Tempus Spin RNA isolation kit (Invitrogen, Thermo Fisher Scientific, Waltham MA, USA). RNA quality was determined by RNA gel for all samples and by Bioanalyzer 6000 (Agilent, Santa Clara CA, USA) for samples selected for RNA-sequencing, which showed high RNA quality with a median RIN of 8.8 (lowest RIN 7.9, highest RIN 9.9). We used 600 ng total RNA of 49 samples for small RNA library construction (NEBNext Multiplex Small RNA library prep set for Illumina, New England Biolabs, Ipswich MA, USA) and selected 24 out of the 49 short RNA-sequenced samples for PolyA-selected mRNA sequencing. These libraries were prepared from 1000 ng total RNA using the TruSeq RNA library preparation kit (Illumina, San Diego CA, USA) and were sequenced on the Illumina NextSeq 500 platform at the Hebrew University's Center for Genomic Technologies.

4.1.4 RNA Sequencing Alignment

Small RNA species were aligned after quality filtering using flexbar and miRExpress 2.0, as described in Section 3.6.2. Additionally, to assess tRF expression, small RNA reads were aligned to the exclusive tRNA space using the MINTmap pipeline.⁹⁷ Briefly, this pipeline compares short RNA sequencing reads with a collection of

sequences determined to only be contained inside mature tRNAs, without confounding from the many tRNA lookalikes in the human genome, e.g., in pseudogenes. The two RNA species were united into one expression matrix containing both miRNA and tRF expression.

Large RNA species were aligned to the human transcriptome (ENSEMBL transcriptome *Homo sapiens* GRCh38 release 79) using the fast dual-phase parallel inference algorithm *Salmon*.¹⁶⁹ Briefly, the method combines an »online« fragment mapping utilising continuous updating of a Bayesian prior with an »offline« phase that determines fragment quantities by application of the Bayesian model determined before via a standard expectation maximisation (EM) algorithm or a variable Bayesian EM. Additionally, the pipeline corrects for multiple typical biases in sequencing, such as position-specific biases, sequence-specific 3' and 5' end biases, fragment GC content bias, and fragment length distribution. The resulting quantified fragments were imported into R using the rsubreads package.[?]

4.1.5 Quality Control and Filtering

Raw and processed reads were quality-controlled using FastQC, as described in Section 3.6.1, with no samples falling below acceptable thresholds. Small and large RNA alignments were batch-corrected followed by analysis of inter-sample relationships via the method proposed by Oldham *et al.* (as described in Section 3.9.2). We excluded no large RNA samples and one small RNA sample (»11_40044_S12»).

4.1.6 RNA Sequencing Differential Expression Analysis

Quantified reads were subjected to differential expression analysis using DESeq2, essentially as described in Section 3.6.3. Small RNA species were analysed together by combining count tables for miRNAs and tRFs, large RNA were analysed separately. Both datasets were corrected for covariates *subject age* and *batch*. Correction for patient sex was not necessary because all patients in the final analyses were male. Log₂-fold changes were shrunk using *apeglm* as described in Section 3.6.3, at an alpha level of 0.1.

4.1.7 Gene Ontology Analyses

We performed GO analyses on the set of DE transcripts, using different ranking methods. GO analyses were performed using R/topGO as described in Section 3.7.2 using the weighted method.

Ranking by P-Value

Transcripts were ranked by p-value, and different test sets were tested against the background of the topmost two thousand transcripts. We tested the set of all DE transcripts (adjusted p-value < 0.05) and the separate sets of positively and negatively regulated transcripts. Additionally, for each test group, the criterion of log₂-fold changes > 1.4 was applied and re-tested.

Ranking By Count-Change

Alternatively, transcripts were ranked by count-change, and the top 100 significantly DE transcripts were tested against the background of the first two thousand transcripts. Similarly to the p-value ranking, test sets comprised all transcripts as well as only negatively or positively regulated transcripts.

4.1.8 Homology Computation Among tRNA Fragments

Transfer RNA fragment origin can be ambiguous, even in fragments derived from tRNA-exclusive space. To assess sequence-based relationships between tRFs, all detected fragments were subjected to pairwise homology analysis using local Smith-Waterman alignment (*pairwiseAlignment* function of the R/Biostrings package), and scores were transformed into a distance matrix to enable clustering and visualisation of relationships. t-SNE (t-Distributed Stochastic Neighbour Embedding) was employed to visualise tRF homologies in a 2D space.

4.1.9 t-Distributed Stochastic Neighbour Embedding

SNE (Stochastic Neighbour Embedding) replaces Euclidian distances between data points with conditional probabilities that represent similarities. The Gaussian distribution used in SNE to represent the probability density for any given data point in the low-dimensional space is replaced by a Student's t-Distribution in the updated t-SNE algorithm. In combination with the use of a symmetrised function with simpler gradients, this alleviates problems with optimisation of the cost function that is used to create forces between points on the low-dimensional map.¹⁷⁰

t-SNE was used in a variety of applications to reduce the dimensionality of high-dimensional data, for instance, the amino acid origin of tRFs, or the association of tRFs with distinct cell types in the blood. t-SNE analyses were performed in R, using the Rtsne package.(cite) t-SNE requires, apart from the input data, a parameter called *perplexity*, which determines the weighting of local as opposed to global effects in the data. So far, there are no strict rules governing the selection of a perplexity value, other than that the perplexity cannot exceed the number of individual data points. Since different perplexities can give widely varying results, which can sometimes be misleading, the resulting maps have to be screened with a range of perplexities to assess their robustness.

4.1.10 Cholinergic Association of Small RNA Species

To determine association of distinct smRNAs with cholinergic transcripts, we analysed the multiple-targeting relationships of each distinct smRNA towards our curated list of cholinergic-associated (CA) transcripts. We first created complete targeting data of all DE smRNAs towards all cholinergic transcripts, which we then successively filtered for multiple targeting behaviours. To assess the base level of multiple targeting of cholinergic transcripts, we utilised empirical cumulative density functions of the number of individual cholinergic targets of each miRNA and tRF (Figure 4.1). We assumed 80% to be a robust threshold of cholinergic targeting, and for diverging numbers between miRNAs and tRFs chose to use the higher (more stringent) threshold. smRNAs above this threshold (i.e., smRNAs targeting at least as many cholinergic transcripts as the threshold value) were considered CA.

4.2 DESCRIPTIVE ANALYSIS OF RNA DYNAMICS IN BLOOD AFTER STROKE

4.2.1 DIFFERENTIAL EXPRESSION OF LARGE RNA

At an alpha level of 0.05, we detected 694 differentially expressed (DE) transcripts, 204 of them up-, and 490 down-regulated (Figure 4.2 A). 18 of the up-regulated, and 109 of the down-regulated

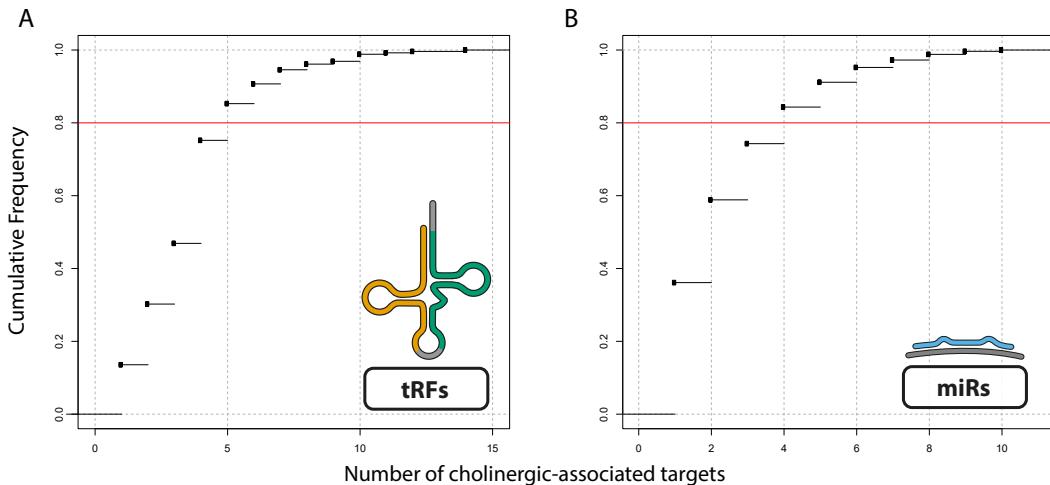


Figure 4.1: Cholinergic-associated Small RNA ECDF Curves. Cholinergic association was tested using *miRNeo* targeting data of miRNAs and tRFs. To assess the best-suited threshold for defining cholinergic association, empirical cumulative density functions were calculated for the number of cholinergic-associated (CA) genes targeted by each unique smRNA. **A)** Cumulative frequency of number of CA genes targeted by tRFs. Threshold of 80% (red line) is passed at five CA genes targeted. **B)** Cumulative frequency of number of CA genes targeted by miRNAs. Threshold of 80% (red line) is passed at four CA genes targeted.

transcripts exceeded the common \log_2 -fold change threshold of 1.4. To determine the most-impacted pathways, we performed GO analyses.

4.2.2 GENE ONTOLOGY ANALYSES OF DIFFERENTIALLY EXPRESSED GENES

Ranking of all transcripts (regardless of direction of regulation) by their differential expression p-value resulted in GO terms mainly related to circulatory system processes ($p = 0.018$) and immunity. Most notable immune-related terms included cytokine-mediated pathways ($p = 2.4E-04$), response to IFNs α ($p = 0.013$) and β ($p = 1.2E-03$), regulation of JAK/STAT cascade ($p = 0.013$), response to LPS ($p = 0.025$), and macrophage activation ($p = 0.026$). Limiting the test set to transcripts with \log_2 -fold change above 1.4 increased sensitivity towards immune processes, yielding lower p-values for the enrichment of positive ($1.7E-04$) and negative regulation of cytokine production ($5.7E-04$), type I interferon production ($3.9E-04$), response to bacterium ($5.9E-04$), innate immune response ($2.0E-03$), response to organophosphorous ($2.3E-03$), cytoplasmatic pattern recognition receptor signalling pathway ($2.8E-03$), and response to LPS ($9.1E-03$).

Positively regulated transcript pertained to circulatory system processes, such as platelet degranulation ($1.2E-03$) and aggregation (0.02), and sprouting angiogenesis ($4.8E-03$), but also antigen processing and presentation ($4.5E-03$). Test set limitation to \log_2 -fold change above 1.4 did not increase sensitivity of those terms, but presented essentially similar results. Negatively regulated transcripts were enriched in terms involving response to IFN α ($1.3E-03$) and β ($3.1E-04$), response to LPS ($1.5E-03$), rhythmic process ($2.5E-03$), positive regulation of T cell proliferation ($4.3E-03$), positive regulation of JAK-STAT cascade (0.015), and cellular response to IL-1 (0.019). Test set lim-

itation to \log_2 -fold change above 1.4 again increased sensitivity towards immune-related terms, but without changing the general pattern.

As a cross-check, DE transcripts were ranked by count-change, and re-analysed. The top 100 changed transcripts, without regard to direction (absolute count-change) yielded terms implying response to IFN α (3.8E-04), β (1.1E-04), and γ (1.4E-04), mitochondrial organisation (5.6E-03) and ATP synthesis (6.9E-03), response to IL-4 (6.9E-03), positive regulation of JAK-STAT cascade (8.8E-03), response to antibiotic (0.044), and platelet degranulation (0.045). The top 100 positively regulated transcripts yielded terms involving platelet degranulation (3.8E-03), mitochondrial ATP synthesis (4.2E-03), response to xenobiotic stimulus (0.017), platelet aggregation (0.013), and response to antibiotic (0.016), while the top 100 negatively regulated transcripts were associated with inflammatory response (1.3E-04), regulation of apoptosis (1.8E-04), cytokine secretion (6.8E-04), antigen processing and presentation (1.2E-03), regulation of lymphocyte apoptosis (2.3E-03) and proliferation (2.7E-03), response to antibiotic (5.2E-03), leukocyte homeostasis (7.6E-03), response to IL-1 (7.6E-03), and many more immune-specific processes. For a full list of all terms from these

DE genes?

[analyses](#), see Appendix D.

4.2.3 DIFFERENTIAL EXPRESSION OF SMALL RNA

In the simultaneous co-analysis of miRNAs and tRFs, we detected 420 DE miRNAs and 143 DE tRFs (adjusted p-value < 0.05, Figure 4.2 B&C). 63% of miRNAs (265) were down-regulated, while 87% of tRFs (124) were up-regulated. tRFs were mainly derived from the 3' end (3'-tRFs, 87) or from internal tRNA regions (i-tRFs, 48), while the tRFs from 5' ends (5'-tRFs) were in the minority (6). The amino acid distribution was shifted in favour of alanine- (35), glycine- (28), and proline-carrying (12) tRNAs (Figure 4.2 D). 30 of the 35 alanine-associated tRFs were 3'-tRFs, and all of those were up-regulated, indicating non-random generation of these fragments.

4.2.4 HOMOLOGY AMONG tRNA FRAGMENTS

Using pairwise homology among all DE tRFs, visualised via t-SNE (see Section 4.1.9), we identified clusters of highly similar fragments, that correlate with their amino-acid origin, i.e., the amino acid which is carried by the respective parent tRNA (Figure 4.2 E). This relationship persisted across distinct individual tRNAs coding for the same amino acid, and was particularly pronounced in tRNAs associated with alanine, glycine, leucine, proline, and methionine.

4.2.5 CHOLINERGIC ASSOCIATION OF SMALL RNA SPECIES

The association of smRNA species to distinct systems or pathways is not trivial because of the multiple-targeting nature of these RNAs. For the purpose of the following analyses, we define a small RNA as being associated with cholinergic processes by the positive association of the smRNA with a number

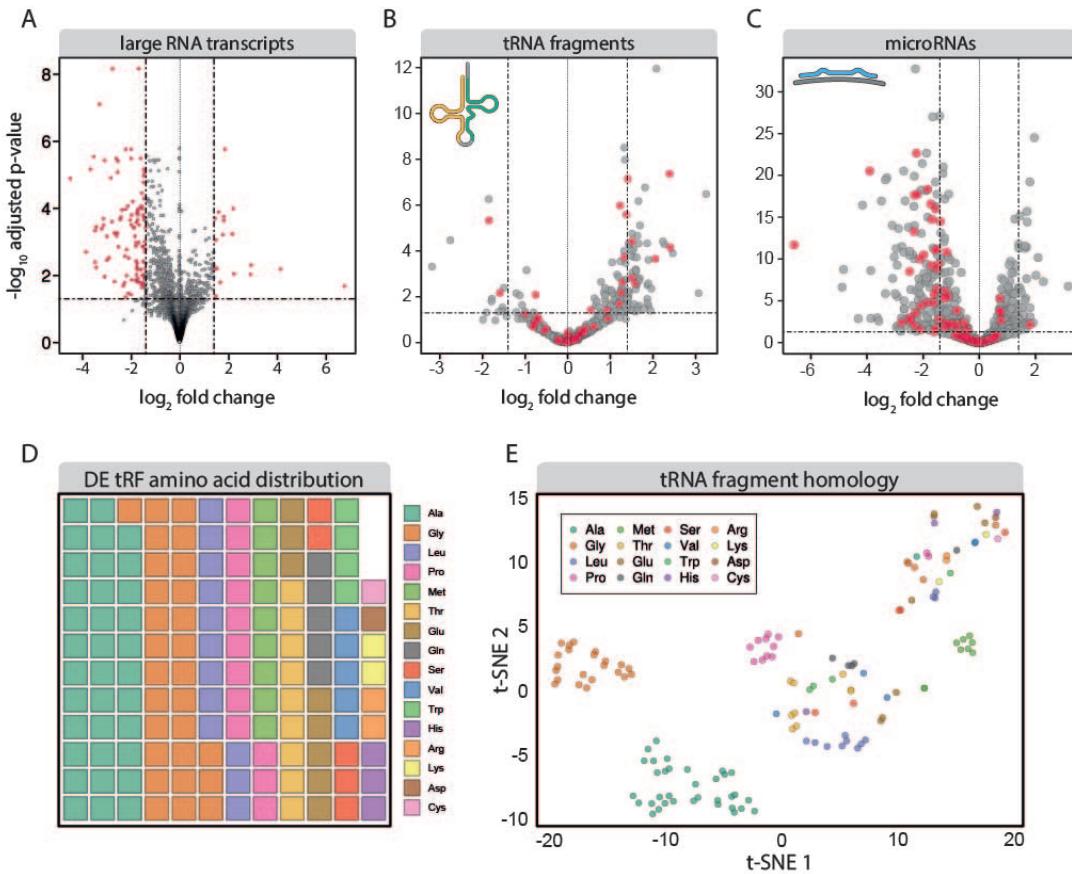


Figure 4.2: Small and Large RNA Differential Expression and tRF Properties. A) Differential expression analysis reveals multiple large RNA transcripts changed in patient blood after stroke. The majority of differentially expressed (DE) transcripts above a \log_2 fold change threshold of 1.4 (red) are down-regulated. B) Blood-borne tRNA fragments (tRFs) also change after stroke. Unlike large RNA and miRNAs, the majority of DE tRFs are up-regulated. Cholinergic-associated (CA) tRFs in red. C) Blood-borne miRNAs are heavily influenced by the events following stroke. Like large RNA transcripts, miRNAs are also overwhelmingly down-regulated. CA miRNAs in red. D) The distribution of amino acid origin among the DE tRFs is non-random and biased towards the amino acids alanine, glycine, leucine, proline, and methionine. Each square represents one DE tRF, colour denotes amino acid origin. E) t-SNE of pairwise fragment homology by local Smith-Waterman alignment shows clustering of the dominant amino acid groups of tRFs. Clear clusters can be observed for tRFs derived from tRNA carrying alanine, glycine, leucine, proline, and methionine.

of cholinergic-related large transcripts. We do not assess the question whether this small RNA also targets other systems equally, or even if it targets cholinergic transcripts with greater likelihood than a random selection of genes. For this reason, we select a fairly high threshold for the definition of a »cholinergic« smRNA, which is the targeting of at least 5 cholinergic-related transcripts (above 80% on the empirical cumulative density function of cholinergic targeting, see Section 4.1.10).

Following this definition, we detected 52 CA miRNAs (90% down-regulated, 5 up and 47 down), and 18 CA tRFs (83% up-regulated, 15 up and 3 down). Above a threshold of \log_2 -fold change of 1.4, we found 33 CA miRNAs (97% down-regulated, 1 up and 32 down), and 9 CA tRFs (78% up-regulated, 7 up and 2 down). CA smRNAs are marked in red in Figure 4.2 B&C.

amino acids?

4.3 BLOOD COMPARTMENTS OF CHOLINERGIC SYSTEMS AND SMALL RNA SPECIES

To address the shortcomings of whole-blood RNA sequencing, which is more representative of the clinical setting, but less specific regarding cellular compartments, we consulted third party datasets to assess RNA species distribution in different cellular and non-cellular compartments of the blood. For small RNA species, we re-analysed a published dataset of small RNA-seq of 450 human samples from various blood tissues;¹⁷¹ for large RNA transcripts, we utilised the tissue specificity of Marbach's regulatory circuits.⁹⁸ The large RNA information was used to identify blood cell types with cholinergic transcriptional activity, which was then used to zoom in to small RNA expression subsets related to cholinergic processes.

4.3.1 Large RNA Regulatory Circuits in Tissues of the Blood

To evaluate the cell type distribution of cholinergic genes in blood tissue types, we utilised the expression patterns derived from cumulative transcription factor activity of Marbach's regulatory circuits.⁹⁸ As shown by the authors, the cumulative activity of all transcription factors towards one gene describe well the actual expression of that gene in the respective tissue type. To maximise comparability to the parallel analyses of small RNA species (Section 4.3.2), blood cell types (i.e., »regulatory circuits«) were selected to reflect the cell type selection of Juzenais *et al.*¹⁷¹ based on similar markers of the »cluster of differentiation« family of genes. These were: CD4-positive T-helper cells, CD8-positive cytotoxic T-cells, CD14-positive monocytes, CD15-positive neutrophils, CD19-positive B-cells, CD56-positive natural killer cells, and, for comparison, whole blood. For the sake of simplicity, genes were considered »present« in each blood tissue type if at least one TF showed activity towards the gene.

4.3.2 An Atlas of Small RNA Expression in Cell Types of the Blood

To evaluate the cell type distribution of our small RNA molecules, we analysed a dataset deposited by Juzenais *et al.*,¹⁷¹ who separated and sequenced 450 samples comprising seven types of individual blood cell types (characterised by »cluster of differentiation«-type membrane-bound receptors), serum, exosomes, and whole blood. The individual blood cell types comprised CD4-positive T-helper cells, CD8-positive cytotoxic T-cells, CD14-positive monocytes, CD15-positive neutrophils, CD19-positive B-cells, CD56-positive natural killer cells, and CD235a-positive erythrocytes (the only distinct cell type not available in Marbach's regulatory circuits, since mature erythrocytes do not transcribe). Starting from the raw data deposited on NCBI GEO, we controlled the quality, applied quality based filtering, and aligned the 450 samples to miRNA and tRF sequences, as described above. The original publication did not offer statistical analyses because of a failure in the spike-in procedure, and defined presence of a small RNA by a measure of at least five counts in 85% of samples. However, since this definition relies heavily on sequencing depth, and depth can vary widely even in methodically robust sequencing experiments depending on a large number of variables (see Figure 3.6 C), we defined our own test for descriptive analysis of presence or absence of lowly expressed small RNAs in each of the sample types:

Definition of Presence and Absence of Lowly Expressed smRNA Molecules

This definition comprises estimation of a log-normal distribution from a small RNA expression profile, and a statistical test to refute the null hypothesis that the distribution is in fact log-normal. The danger of evaluating true expression of lowly expressed smRNA molecules by a count-based threshold is the possibility of random reads resulting from degradation products of highly expressed RNA with similar sequence, and the amplification of noise. Both problems are exacerbated by an increase in sequencing depth. In today's RNA-seq technology, most chips can accommodate only a limited amount of samples compared to the amount of reads that can be generated. While this is not as problematic in cases of longer inserts and paired design, which is usually employed in large RNA-seq, in small RNA-seq this can lead to enormous overheads of reads. It is not uncommon to receive tens of millions of reads for each sample, which exceeds the recommended amount (of at least one million) by large margins.

Thus, there is the need to distinguish between degradation products of highly expressed RNA molecules or amplified noise and legitimate lowly expressed smRNA molecules. The central assumption for our proposed method is: The expression pattern of legitimate smRNA molecules follows, as is common in biology, a normal distribution of some kind, or, for the discrete case, a normal poisson distribution. On the other hand, degradation products or noise would rather follow other, »non-biological« distributions, such as the uniform distribution or a monotonously decreasing power-law distribution such as the Pareto distribution. Thus, we chose to statistically test each smRNA in each tissue type for the adherence to this criterion, by comparing the measured counts with a distribution function estimated based on the mean and standard deviation of the measured counts. During testing, we found the log-normal distribution to give the best classification results.

The distribution mean and standard deviation of the expression values per cell type and smRNA were estimated using the *fitdist* function of the R/*fitdistrplus* package.(cite) The count distribution was then tested against a log-normal distribution with the estimated mean and standard deviation via the R implementation of the Kolmogorov-Smirnov test, with a cutoff of 0.1. The small RNA was defined as present if the test failed to reject the null hypothesis (see Appendix ?? for numerous examples).

Analysis of Expression Patterns and Establishment of Virtual Tissues

The distribution of smRNA expression across the different cell types was used to assign 8 functional compartments (i.e., »virtual tissues«) to the entirety of detected fragments such that each smRNA was sorted into one of the tissue classes. Ideally, these classes would be unambiguous, i.e., there would be no overlap of smRNA molecules between the classes. Eight classes were created via hierarchical clustering of miRNA and tRF expression separately (Figure 4.4), and then used in combination with t-SNE applied to the entire expression matrix, to visualise the compartmentalisation of smRNAs in these virtual tissues. The samples taken from stroke patients in the PREDICT study were sequenced from whole blood, which precludes direct information about tissue distribution. Thus, the two-dimensional maps from t-SNE visualisation were used to, first, explore the tissue association of smRNAs differentially expressed in stroke patients' whole blood samples, and second, examine the potential impact of cholinergic-associated smRNAs in these tissues.

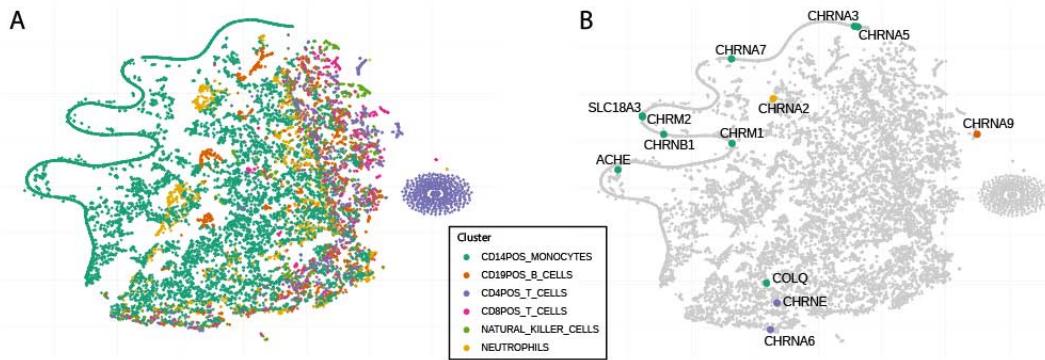


Figure 4.3: Large RNA Expression Patterns in Blood-Borne Cells. Expression derived from transcriptional activity in blood-borne cell types in the Marbach dataset⁹⁸ was visualised via t-SNE. The input matrix comprised all 15 032 detected genes in 6 types of blood-borne immune cells. Genes were plotted on the first two t-SNE dimensions and coloured by the cell type of their highest expression, i.e., the highest cumulative transcriptional activity of all active TFs. **A)** Complete t-SNE shows a gradient of expression across the different cell types, with much expression in CD14-positive monocytes and T-cells. CD4-positive T-cells possess a cluster of genes that distinguishes itself from the bulk of the other genes. **B)** Highlighting of cholinergic core genes reveals an enrichment in close compartments of CD14-positive monocytes. The vesicular ACh-transporter SLC18A3 can serve as substitute for the main cholinergic marker, CHAT, as discussed in Section 2.2.3.

4.3.3 LARGE RNA EXPRESSION PATTERNS IDENTIFY CHOLINERGIC SYSTEMS IN CD14-POSITIVE MONOCYTES

4.3.4 IDENTIFICATION OF FUNCTIONAL ENRICHMENT OF smRNA EXPRESSION IN BLOOD-BORNE CELLS

To date, there is no comprehensive expression catalogue of smRNA species expression in the tissue types of the human body that is comparable to what has been achieved in the description of large RNA. To classify the detected smRNAs in a manner specific to tissues in human blood, we utilised a dataset published by Juzenas *et al.*,¹⁷¹ who describe miRNA expression in a variety of blood tissues (see Section 4.3.2 for details). We re-analysed the publicly deposited data for miRNA and tRF expression, and developed our own method of defining »presence« of the smRNA in each tissue type based on the evaluation of a log-normal distribution model (instead of using a simple count threshold).

Using this presence/absence data, we first utilised hierarchical clustering to establish »virtual tissues« that could be assigned to each smRNA (Figure 4.4) for later evaluation in the stroke patient sequencing. Both miRNAs and tRFs showed a number of clearly associated smRNAs with several compartments, whereas other compartments and smRNAs were distributed in a more complex manner. The ten tissue types of the Juzenas *et al.*¹⁷¹ study were equally parted into two five-tissue superclusters by the expression patterns of both smRNA species (Figures 4.4 A&B, x-axis). These two clusters distinguish immune from non-immune compartments in the blood, but for one notable exception: while the immune cluster comprises monocytes, T-cells, B-cells, and NK-cells, the non-immune cluster contains neutrophils in addition to erythrocytes and the non-cellular tissues serum, exosomes, and whole blood. Notably, the neutrophil samples cluster closest to the whole blood compartment in both smRNA species.

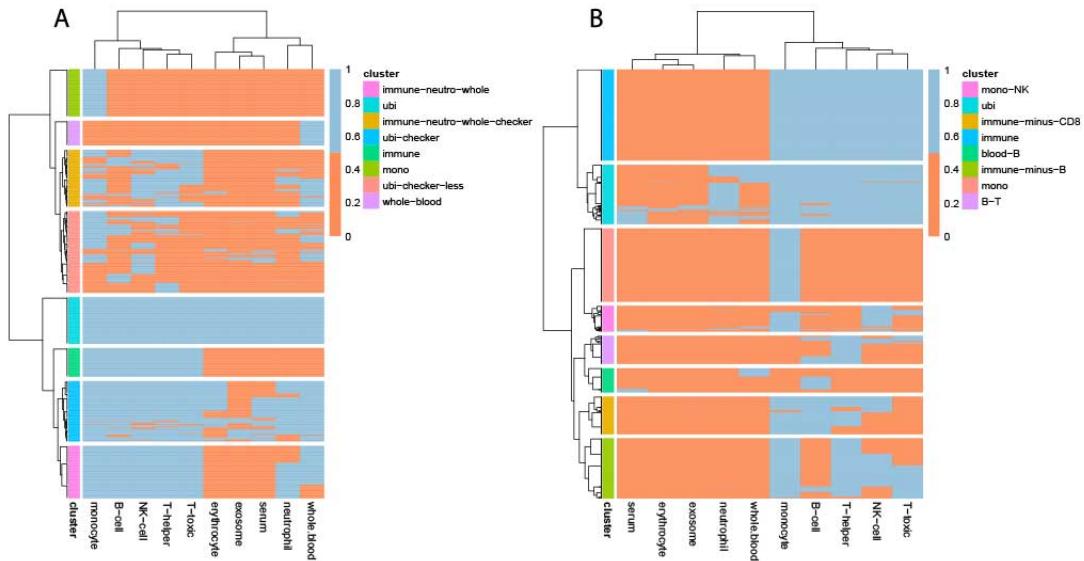


Figure 4.4: Functional Characterisation of Hierarchical Clusters in Blood Cell Small RNA Expression. Information on presence/absence of miRNAs and tRFs in the tissue types analysed in Juzenas *et al.*¹⁷¹ were hierarchically clustered into 8 clusters using the Ward method,² and plotted on a heatmap (single smRNAs on the y-axis, tissue types on the x-axis). To assign meaning to these clusters, manual inspection was followed by annotation of enrichment in tissue types. Complex combinations were approximated by their most prominent features. **A)** Clusters of miRNA presence/absence in blood cell compartments. Clearest cluster association was shown by miRNAs expressed only in monocytes (»mono«), in all blood-borne immune cells except neutrophils (»immune«), ubiquitously without exception (»ubi«), and only in whole blood (i.e., in none of the single compartments (»whole-blood«)). **B)** Clusters of tRF presence/absence in blood cell compartments. Clearest cluster association was shown by tRFs expressed only in monocytes (»mono«), and in all blood-borne immune cells except neutrophils (»immune«). The other tissue-related clusters were not as clear as in the miRNA expression data, indicating a looser association to cell type of tRNA-derived smRNAs.

Two distinct virtual tissues showed high consistency in both smRNA species: a virtual tissue containing only CD14-positive monocytes and another tissue comprising all studied cellular immune components except neutrophils (i.e., monocytes, B-cells, CD4- and CD8-positive T-cells, and NK-cells). miRNAs (Figure 4.4 A), in addition, yield clear clusters for whole blood expressed miRNAs, and for miRNAs expressed ubiquitously without exception. In tRFs (Figure 4.4 B), the general picture is more complicated, as the clusters are often mixed.

4.3.5 EXPRESSION PATTERNS OF DIFFERENTIALLY EXPRESSED AND CHOLINERGIC-ASSOCIATED smRNAs

4.4 REGULATORY CIRCUITS OF SMALL RNA AND TRANSCRIPTION FACTORS IN CD14-POSITIVE MONOCYTES

4.4.1 Comprehensive Circuit Network Creation

The comprehensive transcriptomic network in CD14-positive monocytes was created in a two-step process of *miRNeo* targeting. First, the complete TF→gene network was created from the targeting data derived from Marbach *et al.*⁹⁸, yielding a CD14-specific network comprising XX TFs with activity towards XX transcripts. Second, this network was then subjected to successive *miRNeo* targeting of all transcripts in the network by miRNAs and



Figure 4.5: Small RNA Expression Patterns in Blood-Borne Cells. Two-dimensional expression maps were created using t-SNE on the full numeric expression data derived from re-analysis of the Juzenas *et al.* ¹⁷¹ data set for miRNAs and tRFs separately. Single smRNAs (points) were coloured by the virtual tissues derived from the cluster heatmap analysis (Figure 4.4). Node size reflects absolute count change in C, D, E, and F). Shown are full data, differentially expressed (DE) smRNAs, and cholinergic-associated (CA) smRNAs for each species. **A)** Full t-SNE visualisation of miRNA expression. The largest 2D-associative clusters are comprised of the clearest presence/absence virtual tissues, monocytes (yellow) and ubiquitously expressed (orange). Smaller clusters can be identified for the tissue of all immune cells except neutrophils (green) and the complex cluster of immune cells including neutrophils and whole blood (turquoise). **B)** Full t-SNE visualisation of tRF. The largest 2D-associative clusters are, as in miRNAs, comprised of the clearest presence/absence virtual tissues, monocytes (brown) and immune cells except neutrophils (pink). A smaller cluster can be identified for ubiquitously expressed tRFs (orange). **C)** miRNAs DE after stroke are ubiquitously expressed in all virtual tissues. Highest differential expression is seen in the »ubi« cluster. **D)** Likewise, tRFs DE after stroke are ubiquitously expressed in all virtual tissues, and highest differential expression is seen in the »ubi« cluster. **E)** CA miRNAs are enriched in the lower quadrants of the 2D map, particularly in the clusters associated with ubiquitous expression (»ubi«, »ubi-checker«). **F)** CA tRFs show a similar distribution, skewed towards virtual tissues with ubiquitous expression. This indicates covariation of detection with broadness of expression (see text).

tRFs.

For each node fulfilling an active role in this network (i.e., miRNAs, tRFs, and TFs), an activity parameter was computed. The activity of each node is hereby defined as the sum of all scores of each of its targeting relationships. In the case of miRNAs, the score is the summary score introduced in Section 2.2.4, for tRFs, it is the score calculated with the BL-PCT method (see Section 2.2.6), and for TFs, it is the transcriptional activity given by Marbach *et al.*⁹⁸ Activities were normalised, for each biotype separately, by scaling the calculated values v onto a range between 0 and 1, using

$$v_{i,\text{norm}} = \frac{v_i}{\max(v)}$$

, with $\max(v)$ being the maximum of all scores in this biotype category, and all $v > 0$. The activity of each relationship determined the weight of the edge between the two connected nodes.

The network was visualised in gephi,¹⁴⁰ omitting all non-TF genes, and using ForceAtlas2 to generate a force-directed 2D map of smRNA→TF interactions in CD14-positive monocytes. Network modularity was calculated using the function included in gephi, with a resolution of 2.0, to yield 2 distinct modularity classes. The module association of TFs to the tRF- and miRNA-associated modules were used to perform subsequent analyses of the distinct modules.

4.4.2 Gene Ontology Analyses of TF→Gene Networks of CD14-positive Monocytes

The TF→gene networks of each of the two modules derived from smRNA species association (miRNAs versus tRFs) were analysed using topGO¹³⁹ essentially as described in Section 3.7.2. Genes were ordered according to the cumulative activity of TF targeting of each gene in CD14-positive monocytes. To display a range of top genes, transcript background was iterated in five equal steps from 1000 transcripts to the maximum size of target transcripts in each network (12 927 for miRNA-targeted TFs, 12 904 for tRF-targeted TFs). The test set was the top 10% of transcripts for each background size.

4.4.3 DICHOTOMY OF SMALL RNA TARGETING OF TRANSCRIPTION FACTORS IN CD14-POSITIVE MONOCYTES

Organisation of the smRNA→TF network via a force-directed algorithm resulted in visible clustering of two distinct subnetworks, that are governed by miRNAs and tRFs, respectively (Figure 4.6). Inside this network, 10 TFs were found DE in patient blood after stroke (Figure 4.6 A). Calculation of modularity clearly divided the network into TFs primarily influenced by miRNAs and TFs primarily influenced by tRFs (Figure 4.6 B). Based on these two sets of TFs, two distinct TF→gene networks were created: 289 miRNA-biased TFs with 152 649 unique TF→gene targeting relationships, and 280 tRF-biased TFs with 163 641 unique TF→gene targeting relationships. It is notable that, although the graph shows clear segregation between miRNA-targeted and tRF-targeted transcripts, merely XX TFs are targeted by only one of the two smRNA species (XX only by miRNAs and XX only by tRFs).

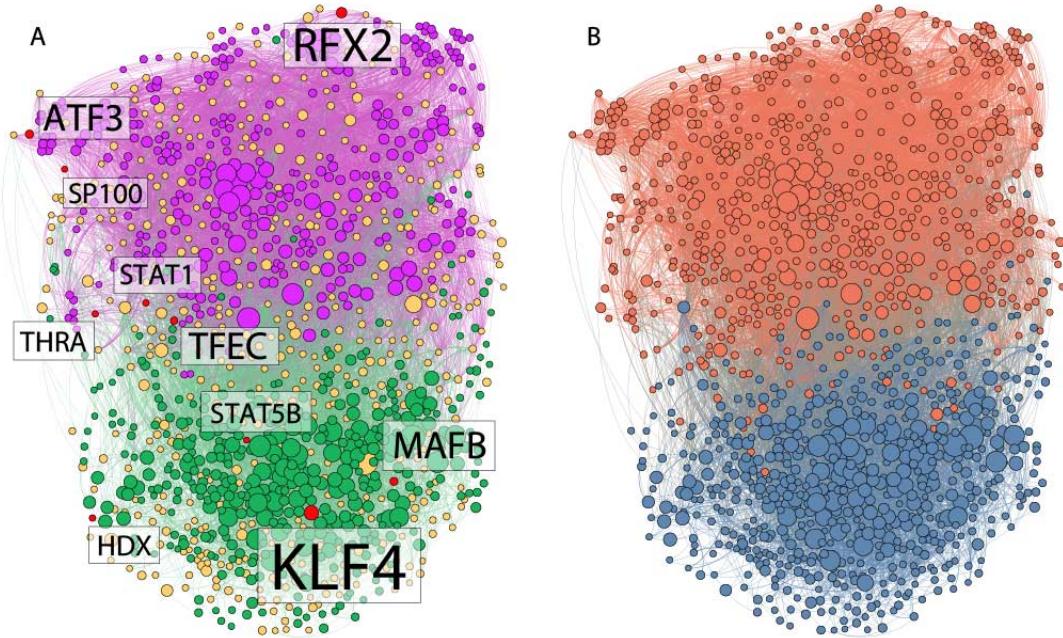


Figure 4.6: Small RNA Targeting of Transcription Factors in CD14-positive Monocytes. The two-dimensional map of all TFs active in CD14-positive monocytes and their smRNA controllers was created via *miRNeo* targeting and visualisation via force-directed algorithm. Node size is determined by activity (see Section 4.4.1). **A)** Nodes coloured by biotype: miRNAs - green, tRFs - purple, TFs - yellow, differentially expressed TFs: red. TFs targeted mainly by tRFs segregate visually from TFs targeted mainly by miRNAs. Both sets contain DE TFs, indicating complementary function. **B)** The network was segregated into two modules by internal network parameters. Node colour denotes modularity class association. Network modularity largely reflects TF targeting of tRFs versus miRNAs.

4.4.4 TRANSCRIPTOMIC FOOTPRINTS OF DICHOTOMOUS TRANSCRIPTION FACTORS IN CD14-POSITIVE MONOCYTES

To determine the putative effect of TF regulation by each smRNA species, we evaluated the potential impact of the TFs most targeted by either miRNAs or tRFs. The top 10% of TF targets in CD14-positive monocytes (derived from Marbach *et al.*⁹⁸) were subjected to iterative GO analysis (see Section 4.4.2). Assuming a general effect of repression in tRF-targeted TFs (because the majority of DE tRFs are up-regulated), and a general de-repression in the set of miRNA-targeted TFs (because most DE miRNAs are down-regulated), the putative functional effects of changes in smRNA levels can be described by GO enrichment analysis of these two test sets.

4.4.5 miRNA-TARGETED TRANSCRIPTION FACTORS CONVEY XX

4.4.6 tRF-TARGETED TRANSCRIPTION FACTORS CONVEY XX

4.5 DIMENSIONALITY REDUCTION AND CORRELATION OF EXPRESSION WITH CLINICAL PARAMETERS

4.5.1 WGCNA

Reduction of dimensionality of high-dimensional data, such as the expression of thousands of small RNA species and their correlation with clinical parameters of patients, can also be achieved by means of Weighted Gene Correlation Network Analysis (WGCNA, also R/WGCNA).¹⁵⁷ The aim of this method is the identification of modules that group similarly-expressed genes. Each of these modules is represented by an eigenvector (called *eigengene*), effectively replacing the individual expression parameters of the genes in the module. The eigengene can then be used to determine correlation of covariates to each module.

To represent expression data in a network-based manner, the expression matrix is transformed into an adjacency matrix representing similarities between nodes, i.e., genes. The co-expression similarity s_{ij} between two nodes i and j is defined as the absolute correlation coefficient between the expression profiles

$$s_{ij} = |cor(x_i, x_j)|$$

with x representing the node expression profile. A signed co-expression measure can be used to keep track of the direction of co-expression. Weighted networks allow the adjacency measures to assume continuous values between 0 and 1, thereby retaining the continuous nature of the underlying expression matrix. This is achieved by implementation of a soft threshold, which is usually chosen by the user to result in a scale-free network topology. The connectivity distribution in a scale-free network follows a power law, i.e., the probability of a node with k connections is, for large values of k :

$$P(k) = k^{-\alpha}$$

With α usually between 2 and 3. Many biological networks are assumed to be scale-free; however, this assumption has very recently been empirically questioned.¹⁷²

Once the network is created, modules are identified via unsupervised clustering of densely connected nodes. Biological significance can then be assigned to each of the modules and genes. Module eigengenes can be calculated (the first principal component of the module) and correlated to traits, such as clinical parameters. In this manner, modules with potential functional implications can be identified.

4.5.2 Co-correlation

The fact that a subset of small RNA-seq samples were also subjected to large RNA-seq enables the application of a secondary correlation between those two approaches. For instance, module eigengenes derived from smRNA analysis can be correlated with the mRNA module eigengenes across the common patients. The module similarity σ thus is defined as the correlation coefficient between eigengene profiles e of modules i and j :

$$\sigma = |cor(e_i, e_j)|$$

Significantly correlating modules ($p < 0.05$) were identified and further analysed downstream.

If the human brain were so simple that we could understand it, we would be so simple that we couldn't.

Emerson M. Pugh

5

Discussion

LOREM IPSUM DOLOR SIT AMET, consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

5.1 METHODS

cell model, chat anomaly, regulation of expression of these two, induction, low vs high control genes Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetur erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit, vel lacinia enim metus eu nunc.

5.2 THE CHOLINERGIC/NEUROKINE INTERFACE

Hypothesis: cholinergic and neurokine systems intermingle significantly in the CNS, affecting physiological as well as pathogenic (pathologic?) processes. Multiple angles reject null (orthogonal evidence)

5.3 SMALL RNA THERAPEUTICS AND PHARMACOLOGY

Extant approaches, methods, diseases, PCSK9, asthma, using small RNA antisense as substitute for single-target small molecules, reduce off-target effects, side effects of a different kind

Transcriptomics as basis for selection and design of antisense therapy, combinatorial, compare dirty drugs from psychiatric disorders, serendipity impossible, determinant is the sequence as opposed to functional groups that can be iteratively modified (only 4 building blocks)

6

Conclusion

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetuer. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetuer erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit, vel lacinia enim metus eu nunc.

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Vestibulum tortor quam, feugiat vitae, ultricies eget, tempor sit amet, ante. Donec eu libero sit amet quam egestas semper. Aenean ultricies mi vitae est. Mauris placerat eleifend leo. Quisque sit amet est et sapien ullamcorper pharetra. Vestibulum erat wisi, condimentum sed, commodo vitae, ornare sit amet, wisi. Aenean fermentum, elit eget tincidunt condimentum, eros ipsum rutrum orci, sagittis tempus lacus enim ac dui. Donec non enim in turpis pulvinar facilisis. Ut felis.

Cras sed ante. Phasellus in massa. Curabitur dolor eros, gravida et, hendrerit ac, cursus non, massa. Aliquam lorem. In hac habitasse platea dictumst. Cras eu mauris. Quisque lacus. Donec ipsum. Nullam vitae sem at nunc pharetra ultricies. Vivamus elit eros, ullamcorper a, adipiscing sit amet, porttitor ut, nibh. Maecenas adipiscing mollis massa. Nunc ut dui eget nulla venenatis aliquet. Sed luctus posuere justo. Cras vehicula varius turpis. Vivamus eros metus, tristique sit amet, molestie dignissim, malesuada et, urna.

Cras dictum. Maecenas ut turpis. In vitae erat ac orci dignissim eleifend. Nunc quis justo. Sed vel ipsum in purus tincidunt pharetra. Sed pulvinar, felis id consectetur malesuada, enim nisl mattis elit, a facilisis tortor nibh quis leo. Sed augue lacus, pretium vitae, molestie eget, rhoncus quis, elit. Donec in augue. Fusce orci wisi, ornare id, mollis vel, lacinia vel, massa.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetur erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit, vel lacinia enim metus eu nunc.

Bibliography

- [1] Dale H H. THE ACTION OF CERTAIN ESTERS AND ETHERS OF CHOLINE, AND THEIR RELATION TO MUSCARINE. *Journal of Pharmacology and Experimental Therapeutics*, 6(2) (1914).
- [2] Loewi O. Über humorale Übertragbarkeit der Herznervenwirkung. *Pflügers Arch. Ges. Physiol.*, 189:239–242 (1921).
- [3] Dale H H & Dudley H W. THE PRESENCE OF HISTAMINE AND ACETYLCHOLINE IN THE SPLEEN OF THE OX AND THE HORSE. *J. Physiol.*, 68:97 (1929).
URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1402860/pdf/jphysiol01676-0019.pdf>
- [4] Mesulam M M, Mufson E J, Levey A I, & Wainer B H. Atlas of cholinergic neurons in the forebrain and upper brainstem of the macaque based on monoclonal choline acetyltransferase immunohistochemistry and acetylcholinesterase histochemistry. *Neuroscience*, 12(3):669–686 (1984). doi:10.1016/0306-4522(84)90163-5.
- [5] Mesulam M M & Geula C. Nucleus basalis (Ch4) and cortical cholinergic innervation in the human brain: Observations based on the distribution of acetylcholinesterase and choline acetyltransferase. *The Journal of Comparative Neurology*, 275(2):216–240 (1988). doi:10.1002/cne.902750205.
URL <http://www.ncbi.nlm.nih.gov/pubmed/3220975> <http://doi.wiley.com/10.1002/cne.902750205>
- [6] Lobentanzer S, Hanin G, Klein J, & Soreq H. Integrative Transcriptomics Reveals Sexually Dimorphic Control of the Cholinergic/Neurokinin Interface in Schizophrenia and Bipolar Disorder. *Cell Reports*, pp. 1–19 (2019). doi:10.1016/j.celrep.2019.09.017.
URL <https://doi.org/10.1016/j.celrep.2019.09.017>
- [7] Mesulam M M & Van Hoesen G W. Acetylcholinesterase-rich projections from the basal forebrain of the rhesus monkey to neocortex. *Brain Research*, 109(1):152–157 (1976). doi:10.1016/0006-8993(76)90385-1.
URL <https://www.sciencedirect.com/science/article/abs/pii/0006899376903851?via%3Dihub> <https://linkinghub.elsevier.com/retrieve/pii/0006899376903851>

- [8] Berrios G E. Alzheimer's disease: A conceptual history. *International Journal of Geriatric Psychiatry*, 5(6):355–365 (1990). doi:10.1002/gps.930050603.
URL <http://doi.wiley.com/10.1002/gps.930050603>
- [9] Miech R A, Breitner J C, Zandi P P, Khachaturian A S, Anthony J C, & Mayer L. Incidence of AD may decline in the early 90s for men, later for women: The Cache County study. *Neurology*, 58(2):209–218 (2002). doi:10.1212/WNL.58.2.209.
- [10] Bartus R, Dean R, Beer B, & Lippa A. The cholinergic hypothesis of geriatric memory dysfunction. *Science*, 217(4558):408–414 (1982). doi:10.1126/science.7046051.
URL <http://www.sciencemag.org/cgi/doi/10.1126/science.7046051>
- [11] Braak H & Braak E. Staging of alzheimer's disease-related neurofibrillary changes. *Neurobiology of Aging*, 16(3):271–278 (1995). doi:10.1016/0197-4580(95)00021-6.
URL <https://www.sciencedirect.com/science/article/abs/pii/0197458095000216?via%3Dihub>
- [12] Schmitz T W & Nathan Spreng R. Basal forebrain degeneration precedes and predicts the cortical spread of Alzheimer's pathology. *Nature Communications*, 7:1–13 (2016). doi:10.1038/ncomms13249.
- [13] Schmitz T W, Mur M, Aghourian M, Bedard M A, & Spreng R N. Longitudinal Alzheimer's Degeneration Reflects the Spatial Topography of Cholinergic Basal Forebrain Projections. *Cell Reports*, 24(1):38–46 (2018). doi:10.1016/j.celrep.2018.06.001.
URL <https://doi.org/10.1016/j.celrep.2018.06.001>
- [14] Kraepelin E. *Psychiatrie. Ein Lehrbuch für Studirende und Aerzte*. Leipzig Barth, edition 8, edition (1913).
- [15] Bortolato B, Miskowiak K, Vieta E, Köhler C, & Carvalho A F. Cognitive dysfunction in bipolar disorder and schizophrenia: a systematic review of meta-analyses. *Neuropsychiatric Disease and Treatment*, 11:3111 (2015). doi:10.2147/NDT.S76700.
URL <https://www.dovepress.com/cognitive-dysfunction-in-bipolar-disorder-and-schizophrenia-a-systematic-peer-reviewed-review>
- [16] Van Enkhuizen J, Janowsky D S, Olivier B, Minassian A, Perry W, Young J W, & Geyer M A. The catecholaminergic-cholinergic balance hypothesis of bipolar disorder revisited. *European Journal of Pharmacology*, 753:114–126 (2015). doi:10.1016/j.ejphar.2014.05.063.
URL <http://dx.doi.org/10.1016/j.ejphar.2014.05.063>
- [17] Smucny J & Tregellas J R. Targeting neuronal dysfunction in schizophrenia with nicotine: Evidence from neurophysiology to neuroimaging. *Journal of Psychopharmacology*, 31(7):801–811 (2017). doi:10.1177/0269881117705071.

- [18] Gray S L, Anderson M L, Dublin S, Hanlon J T, Hubbard R, Walker R, Yu O, Crane P K, & Larson E B. Cumulative Use of Strong Anticholinergics and Incident Dementia. *JAMA Internal Medicine*, 175(3):401 (2015). doi:10.1001/jamainternmed.2014.7663.
URL <http://archinte.jamanetwork.com/article.aspx?doi=10.1001/jamainternmed.2014.7663>
- [19] Eum S, Hill S K, Rubin L H *et al.* Cognitive burden of anticholinergic medications in psychotic disorders. *Schizophrenia Research*, 190:129–135 (2017). doi:10.1016/j.schres.2017.03.034.
URL <http://www.ncbi.nlm.nih.gov/pubmed/28390849><http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5628100/><http://linkinghub.elsevier.com/retrieve/pii/S0920996417301718><https://doi.org/10.1016/j.schres.2017.03.034>
- [20] Koukouli F, Rooy M, Tziotis D *et al.* Nicotine reverses hypofrontality in animal models of addiction and schizophrenia. *Nature Medicine*, 23(3):347–354 (2017). doi:10.1038/nm.4274.
URL <http://dx.doi.org/10.1038/nm.4274>
- [21] Forget B, Scholze P, Langa F, Mourot A, Faure P, & Maskos U. Article A Human Polymorphism in CHRNA5 Is Linked to Relapse to Nicotine Seeking in Transgenic Rats Article A Human Polymorphism in CHRNA5 Is Linked to Relapse to Nicotine Seeking in Transgenic Rats. *Current Biology*, pp. 1–10 (2018). doi:10.1016/j.cub.2018.08.044.
URL <https://doi.org/10.1016/j.cub.2018.08.044>
- [22] Sacco K A, Bannon K L, & George T P. Nicotinic receptor mechanisms and cognition in normal states and neuropsychiatric disorders. *Journal of Psychopharmacology*, 18(4):457–474 (2004). doi:10.1177/026988110401800403.
URL <http://www.ncbi.nlm.nih.gov/pubmed/15582913><http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1201375/><http://journals.sagepub.com/doi/10.1177/026988110401800403>
- [23] Rowe A R, Mercer L, Casetti V, Sendt K V, Giaroli G, Shergill S S, & Tracy D K. Dementia praecox redux: A systematic review of the nicotinic receptor as a target for cognitive symptoms of schizophrenia. *Journal of Psychopharmacology*, 29(2):197–211 (2015). doi:10.1177/0269881114564096.
URL <http://journals.sagepub.com/doi/10.1177/0269881114564096>
- [24] Lewis A S, van Schalkwyk G I, & Bloch M H. Alpha-7 nicotinic agonists for cognitive deficits in neuropsychiatric disorders: A translational meta-analysis of rodent and human studies. *Progress in Neuro-Psychopharmacology and Biological Psychiatry*, 75:45–53 (2017). doi:10.1016/j.pnpbp.2017.01.001.

- URL <http://www.ncbi.nlm.nih.gov/pubmed/28065843><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC5446073><https://linkinghub.elsevier.com/retrieve/pii/S0278584616304353>
- [25] Higley M J & Picciotto M R. Neuromodulation by acetylcholine: Examples from schizophrenia and depression. *Current Opinion in Neurobiology*, 29:88–95 (2014). doi:10.1016/j.conb.2014.06.004.
URL <http://dx.doi.org/10.1016/j.conb.2014.06.004>
- [26] Değirmenci Y & Keçeci H. Visual Hallucinations Due to Rivastigmine Transdermal Patch Application in Alzheimer's Disease; The First Case Report. *International Journal of Gerontology*, 10(4):240–241 (2016). doi:10.1016/j.ijge.2015.10.010.
- [27] Leger M & Neill J C. A systematic review comparing sex differences in cognitive function in schizophrenia and in rodent models for schizophrenia, implications for improved therapeutic strategies. *Neuroscience & Biobehavioral Reviews*, 68:979–1000 (2016). doi:10.1016/j.neubiorev.2016.06.029.
URL <http://www.ncbi.nlm.nih.gov/pubmed/27344000><https://linkinghub.elsevier.com/retrieve/pii/S0149763415302712>
- [28] de Leon J & Diaz F J. A meta-analysis of worldwide studies demonstrates an association between schizophrenia and tobacco smoking behaviors. *Schizophrenia Research*, 76(2-3):135–157 (2005). doi:10.1016/j.schres.2005.02.010.
URL <https://linkinghub.elsevier.com/retrieve/pii/S0920996405000757>
- [29] Berger M, (Editor) *Psychische Erkrankungen*. Berger, Mathias (2014).
- [30] Anttila V, Bulik-Sullivan B, Finucane H K et al. Analysis of shared heritability in common disorders of the brain. *Science*, 360(6395):eaap8757 (2018). doi:10.1126/science.aap8757.
URL <http://www.sciencemag.org/lookup/doi/10.1126/science.aap8757>
- [31] Gandal M J, Haney J R, Parikh N N et al. Shared molecular neuropathology across major psychiatric disorders parallels polygenic overlap. *Science*, 359(6376):693–697 (2018). doi:10.1126/science.aad6469.
URL <http://science.sciencemag.org/content/359/6376/693><https://www.sciencemag.org/lookup/doi/10.1126/science.aad6469>
- [32] Ruderfer D M, Ripke S, McQuillin A et al. Genomic Dissection of Bipolar Disorder and Schizophrenia, Including 28 Subphenotypes. *Cell*, 173(7):1705–1715.e16 (2018). doi:10.1016/j.cell.2018.05.046.
URL <https://linkinghub.elsevier.com/retrieve/pii/S0092867418306585>

- [33] Harrison P J. Recent genetic findings in schizophrenia and their therapeutic relevance. *Journal of Psychopharmacology*, 29(2):85–96 (2015). doi:10.1177/0269881114553647.
URL <http://www.ncbi.nlm.nih.gov/pubmed/25315827><http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?artid=PMC4361495><http://journals.sagepub.com/doi/10.1177/0269881114553647>
- [34] Henriksen M G, Nordgaard J, & Jansson L B. Genetics of Schizophrenia: Overview of Methods, Findings and Limitations. *Frontiers in Human Neuroscience*, 11:322 (2017). doi:10.3389/fnhum.2017.00322.
URL <http://journal.frontiersin.org/article/10.3389/fnhum.2017.00322/full>
- [35] Kanazawa T, Bousman C A, Liu C, & Everall I P. Schizophrenia genetics in the genome-wide era: a review of Japanese studies. *npj Schizophrenia*, 3(1):27 (2017). doi:10.1038/s41537-017-0028-2.
URL <http://www.nature.com/articles/s41537-017-0028-2>
- [36] Malhi G S, Tanius M, Das P, Coulston C M, & Berk M. Potential Mechanisms of Action of Lithium in Bipolar Disorder. *CNS Drugs*, 27(2):135–153 (2013). doi:10.1007/s40263-013-0039-0.
URL <http://link.springer.com/10.1007/s40263-013-0039-0>
- [37] Fujii T, Mashimo M, Moriwaki Y, Misawa H, Ono S, Horiguchi K, & Kawashima K. Physiological functions of the cholinergic system in immune cells. *Journal of Pharmacological Sciences*, 134(1):1–21 (2017). doi:10.1016/j.jphs.2017.05.002.
URL <https://linkinghub.elsevier.com/retrieve/pii/S1347861317300695>
- [38] Pavlov V A & Tracey K J. Neural regulation of immunity: Molecular mechanisms and clinical translation. *Nature Neuroscience*, 20(2):156–166 (2017). doi:10.1038/nn.4477.
- [39] Dantzer R. Neuroimmune interactions: From the brain to the immune system and vice versa. *Physiological Reviews*, 98(1):477–504 (2018). doi:10.1152/physrev.00039.2016.
- [40] Lurie D I. An Integrative Approach to Neuroinflammation in Psychiatric disorders and Neuropathic Pain. *Journal of Experimental Neuroscience*, 12:117906951879363 (2018). doi:10.1177/1179069518793639.
URL <http://journals.sagepub.com/doi/10.1177/1179069518793639>
- [41] Fullerton J N & Gilroy D W. Resolution of inflammation: A new therapeutic frontier. *Nature Reviews Drug Discovery*, 15(8):551–567 (2016). doi:10.1038/nrd.2016.39.
URL <http://dx.doi.org/10.1038/nrd.2016.39>

- [42] Newson J, Stables M, Karra E, Arce-Vargas F, Quezada S, Motwani M, Mack M, Yona S, Audzevich T, & Gilroy D W. Resolution of acute inflammation bridges the gap between innate and adaptive immunity. *Blood*, 124(11):1748–1764 (2014). doi: 10.1182/blood-2014-03-562710.
URL <http://www.ncbi.nlm.nih.gov/pubmed/25006125><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4383794><https://ashpublications.org/blood/article/124/11/1748/33037/Resolution-of-acute-inflammation-bridges-the-gap>
- [43] Louveau A, Harris T H, & Kipnis J. Revisiting the Mechanisms of CNS Immune Privilege. *Trends in Immunology*, 36(10):569–577 (2015). doi:10.1016/j.it.2015.08.006.
URL <http://dx.doi.org/10.1016/j.it.2015.08.006>
- [44] Negi N & Das B K. CNS: Not an immunoprivileged site anymore but a virtual secondary lymphoid organ. *International Reviews of Immunology*, 37(1):57–68 (2018). doi:10.1080/08830185.2017.1357719.
URL <http://dx.doi.org/10.1080/08830185.2017.1357719><https://www.tandfonline.com/doi/full/10.1080/08830185.2017.1357719>
- [45] Walsh J, Hendrix S, Boato F *et al.* MHCII-independent CD4+ T cells protect injured CNS neurons via IL-4. *Journal of Clinical Investigation*, 125(2):699–714 (2015). doi:10.1172/JCI76210DS1.
- [46] Meisel C, Schwab J M, Prass K, Meisel A, & Dirnagl U. Central nervous system injury-induced immune deficiency syndrome. *Nature reviews. Neuroscience*, 6(10):775–86 (2005). doi:10.1038/nrn1765.
URL <http://www.ncbi.nlm.nih.gov/pubmed/16163382>
- [47] Bina K G, Rusak B, & Semba K. Localization of cholinergic neurons in the forebrain and brainstem that project to the suprachiasmatic nucleus of the hypothalamus in rat. *Journal of Comparative Neurology*, 335(2):295–307 (1993). doi:10.1002/cne.903350212.
- [48] Xu M, Chung S, Zhang S *et al.* Basal forebrain circuit for sleep-wake control. *Nature Neuroscience*, 18(11):1641–1647 (2015). doi:10.1038/nn.4143.
URL <http://www.ncbi.nlm.nih.gov/pubmed/26457552><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC5776144><http://www.nature.com/articles/nn.4143>
- [49] Van Dort C J, Zachs D P, Kenny J D *et al.* Optogenetic activation of cholinergic neurons in the PPT or LDT induces REM sleep. *Proceedings of the National Academy of Sciences*, 112(2):584–589 (2015). doi:10.1073/pnas.1423136112.
URL <http://www.pnas.org/lookup/doi/10.1073/pnas.1423136112>

- [50] Teles-Grilo Ruivo L M, Baker K L, Conway M W, Kinsley P J, Gilmour G, Phillips K G, Isaac J T, Lowry J P, & Mellor J R. Coordinated Acetylcholine Release in Prefrontal Cortex and Hippocampus Is Associated with Arousal and Reward on Distinct Timescales. *Cell Reports*, 18(4):905–917 (2017). doi:10.1016/j.celrep.2016.12.085.
URL <http://linkinghub.elsevier.com/retrieve/pii/S221124716318071>
- [51] Niwa Y, Kanda G N, Yamada R G *et al.* Muscarinic Acetylcholine Receptors Chrm1 and Chrm3 Are Essential for REM Sleep. *Cell Reports*, 24(9):2231–2247.e7 (2018). doi:10.1016/j.celrep.2018.07.082.
URL <https://doi.org/10.1016/j.celrep.2018.07.082>
- [52] Yamakawa G R, Basu P, Cortese F, MacDonnell J, Whalley D, Smith V M, & Antle M C. The cholinergic forebrain arousal system acts directly on the circadian pacemaker. *Proceedings of the National Academy of Sciences*, 113(47):13498–13503 (2016). doi:10.1073/pnas.1610342113.
URL <http://www.pnas.org/lookup/doi/10.1073/pnas.1610342113>
- [53] Ising M, Lauer C, Holsboer F, & Modell S. The Munich vulnerability study on affective disorders: premorbid neuroendocrine profile of affected high-risk probands. *Journal of Psychiatric Research*, 39(1):21–28 (2005). doi:10.1016/j.jpsychires.2004.04.009.
URL <https://linkinghub.elsevier.com/retrieve/pii/S0022395604000585>
- [54] Wu J C & Bunney W E. The biological basis of an antidepressant response to sleep deprivation and relapse: review and hypothesis. *American Journal of Psychiatry*, 147(1):14–21 (1990). doi:10.1176/ajp.147.1.14.
URL <http://psychiatryonline.org/doi/abs/10.1176/ajp.147.1.14>
- [55] Boonstra T W, Stins J F, Daffertshofer A, & Beek P J. Effects of sleep deprivation on neural functioning: an integrative review. *Cellular and Molecular Life Sciences*, 64(7-8):934–946 (2007). doi:10.1007/s00018-007-6457-8.
URL <http://link.springer.com/10.1007/s00018-007-6457-8>
- [56] Nikanova E V, Gilliland J D, Tanis K Q *et al.* Transcriptional Profiling of Cholinergic Neurons From Basal Forebrain Identifies Changes in Expression of Genes Between Sleep and Wake. *Sleep*, 40(6):16–20 (2017). doi:10.1093/sleep/zsx059.
URL <https://academic.oup.com/sleep/article-lookup/doi/10.1093/sleep/zsx059>
- [57] Balsalobre A. Clock genes in mammalian peripheral tissues. *Cell and Tissue Research*, 309(1):193–199 (2002). doi:10.1007/s00441-002-0585-0.
- [58] King D P, Zhao Y, Sangoram A M *et al.* Positional Cloning of the Mouse Circadian Clock Gene. *Cell*, 89(4):641–653 (1997). doi:10.1016/S0092-8674(00)80245-7.

- URL <http://linkinghub.elsevier.com/retrieve/pii/S0030665708702269https://linkinghub.elsevier.com/retrieve/pii/S0092867400802457>
- [59] Levi-Montalcini R & Booker B. Destruction of the sympathetic ganglia in mammals by an antiserum to a nerve-growth protein. *Proceedings of the National Academy of Sciences*, 46(3):384–391 (1960). doi:10.1073/pnas.46.3.384.
URL <http://www.ncbi.nlm.nih.gov/pubmed/16578497http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC222845http://www.pnas.org/cgi/doi/10.1073/pnas.46.3.384>
- [60] Hefti F. Nerve growth factor promotes survival of septal cholinergic neurons after fimbrial transections. *Journal of Neuroscience*, 6(8):2155–2162 (1986).
- [61] McManaman J L, Crawford F G, Stewart S S, & Appel S H. Purification of a Skeletal Muscle Polypeptide Which Stimulates Choline Acetyltransferase Activity in Cultured Spinal Cord Neurons. *Journal of Biological Chemistry*, 263(12):5890–5897 (1988).
- [62] Rao M S, Patterson P H, & Landis S C. Multiple cholinergic differentiation factors are present in footpad extracts: comparison with known cholinergic factors. *Development (Cambridge, England)*, 116(3):731–44 (1992).
URL <http://www.ncbi.nlm.nih.gov/pubmed/1289063>
- [63] White U A & Stephens J M. The gp130 receptor cytokine family: regulators of adipocyte development and function. *Current pharmaceutical design*, 17(4):340–6 (2011). doi: 10.2174/138161211795164202.
URL <http://www.ncbi.nlm.nih.gov/pubmed/21375496http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3119891>
- [64] Rawlings J S. The JAK/STAT signaling pathway. *Journal of Cell Science*, 117(8):1281–1283 (2004). doi:10.1242/jcs.00963.
URL <http://jcs.biologists.org/cgi/doi/10.1242/jcs.00963>
- [65] Nathanson N M. Regulation of neurokine receptor signaling and trafficking. *Neurochemistry International*, 61(6):874–878 (2012). doi:10.1016/j.neuint.2012.01.018.
URL <https://linkinghub.elsevier.com/retrieve/pii/S0197018612000307>
- [66] Mohamed-Ali V, Goodrick S, Rawesh A, Katz D R, Miles J M, Yudkin J S, Klein S, & Coppack S W. Subcutaneous Adipose Tissue Releases Interleukin-6, But Not Tumor Necrosis Factor- α , in Vivo 1. *The Journal of Clinical Endocrinology & Metabolism*, 82(12):4196–4200 (1997). doi:10.1210/jcem.82.12.4450.
URL <https://academic.oup.com/jcem/article-lookup/doi/10.1210/jcem.82.12.4450http://www.ncbi.nlm.nih.gov/pubmed/9398739>

- [67] Gustafson D, Lissner L, Bengtsson C, Bjorkelund C, & Skoog I. A 24-year follow-up of body mass index and cerebral atrophy. *Neurology*, 63(10):1876–1881 (2004). doi:10.1212/01.WNL.0000141850.47773.5F.
URL <http://www.ncbi.nlm.nih.gov/pubmed/15557505><http://www.neurology.org/cgi/doi/10.1212/01.WNL.0000141850.47773.5F>
- [68] Profenno L A, Porsteinsson A P, & Faraone S V. Meta-Analysis of Alzheimer's Disease Risk with Obesity, Diabetes, and Related Disorders. *Biological Psychiatry*, 67(6):505–512 (2010). doi:10.1016/j.biopsych.2009.02.013.
URL <https://www.sciencedirect.com/science/article/abs/pii/S0006322309002261><https://linkinghub.elsevier.com/retrieve/pii/S0006322309002261>
- [69] Depp C A, Strassnig M, Mausbach B T *et al.* Association of obesity and treated hypertension and diabetes with cognitive ability in bipolar disorder and schizophrenia. *Bipolar Disorders*, 16(4):422–431 (2014). doi:10.1111/bdi.12200.
URL <http://doi.wiley.com/10.1111/bdi.12200>
- [70] Eder K, Baffy N, Falus A, & Fulop A K. The major inflammatory mediator interleukin-6 and obesity. *Inflammation Research*, 58(11):727–736 (2009). doi:10.1007/s00011-009-0060-4.
URL <http://link.springer.com/10.1007/s00011-009-0060-4><http://www.ncbi.nlm.nih.gov/pubmed/19543691>
- [71] Heppner F L, Ransohoff R M, & Becher B. Immune attack: the role of inflammation in Alzheimer disease. *Nature Reviews Neuroscience*, 16(6):358–372 (2015). doi:10.1038/nrn3880.
URL <https://www.nature.com/articles/nrn3880><http://www.nature.com/articles/nrn3880>
- [72] Kirkpatrick B & Miller B J. Inflammation and Schizophrenia. *Schizophrenia Bulletin*, 39(6):1174–1179 (2013). doi:10.1093/schbul/sbt141.
URL <https://academic.oup.com/schizophreniabulletin/article-lookup/doi/10.1093/schbul/sbt141>
- [73] Takao K, Kobayashi K, Hagiwara H *et al.* Deficiency of Schnurri-2, an MHC Enhancer Binding Protein, Induces Mild Chronic Inflammation in the Brain and Confers Molecular, Neuronal, and Behavioral Phenotypes Related to Schizophrenia. *Neuropsychopharmacology*, 38(8):1409–1425 (2013). doi:10.1038/npp.2013.38.
URL <http://www.nature.com/articles/npp201338>
- [74] Hodes G E, Kana V, Menard C, Merad M, & Russo S J. Neuroimmune mechanisms of depression. *Nature Neuroscience*, 18(10):1386–1393 (2015). doi:10.1038/nn.4113.

- [75] Stangl H, Springorum H R, Muschter D, Grässle S, & Straub R H. Catecholaminergic-to-cholinergic transition of sympathetic nerve fibers is stimulated under healthy but not under inflammatory arthritic conditions. *Brain, Behavior, and Immunity*, 46:180–191 (2015). doi: 10.1016/j.bbi.2015.02.022.
URL <http://dx.doi.org/10.1016/j.bbi.2015.02.022>
- [76] Babu M M, Luscombe N M, Aravind L, Gerstein M, & Teichmann S A. Structure and evolution of transcriptional regulatory networks. *Current Opinion in Structural Biology*, 14(3):283–291 (2004). doi:10.1016/j.sbi.2004.05.004.
URL <https://linkinghub.elsevier.com/retrieve/pii/S0959440X04000788>
- [77] Lee R C, Feinbaum R L, & Ambros V. The *C. elegans* heterochronic gene lin-4 encodes small RNAs with antisense complementarity to lin-14. *Cell*, 75(5):843–854 (1993). doi: 10.1016/0092-8674(93)90529-Y.
URL <https://linkinghub.elsevier.com/retrieve/pii/009286749390529Y>
- [78] Rodriguez A. Identification of Mammalian microRNA Host Genes and Transcription Units. *Genome Research*, 14(10a):1902–1910 (2004). doi:10.1101/gr.2722704.
URL <http://www.genome.org/cgi/doi/10.1101/gr.2722704>
- [79] Kozomara A, Birgaoanu M, & Griffiths-Jones S. miRBase: from microRNA sequences to function. *Nucleic Acids Research*, 47(D1):D155–D162 (2019). doi:10.1093/nar/gky1141.
URL <https://academic.oup.com/nar/article/47/D1/D155/5179337>
- [80] Ambros V, Bartel B, Bartel D P *et al.* A uniform system for microRNA annotation. *RNA (New York, N.Y.)*, 9(3):277–9 (2003). doi:10.1261/rna.2183803.
URL <http://www.ncbi.nlm.nih.gov/pubmed/12592000><http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?artid=PMC1370393>
- [81] Salta E & De Strooper B. microRNA-132: a key noncoding RNA operating in the cellular phase of Alzheimer’s disease. *The FASEB Journal*, 31(2):424–433 (2017). doi:10.1096/fj.201601308.
URL <http://www.fasebj.org/doi/10.1096/fj.201601308>
- [82] Lu L F, Gasteiger G, Yu I S *et al.* A Single miRNA-mRNA Interaction Affects the Immune Response in a Context- and Cell-Type-Specific Manner. *Immunity*, 43(1):52–64 (2015). doi:10.1016/j.jimmuni.2015.04.022.
URL <http://www.ncbi.nlm.nih.gov/pubmed/26163372><http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?artid=PMC4529747><https://linkinghub.elsevier.com/retrieve/pii/S1074761315002551>

- [83] Borek E, Baliga B S, Gehrke C W, Kuo C W, Belman S, Troll W, & Waalkes T P. High Turnover Rate of Transfer RNA in Tumor Tissue. *CANCER RESEARCH*, 37:3362–3366 (1977). URL <https://cancerres.aacrjournals.org/content/37/9/3362.full-text.pdf>
- [84] Speer J, Gehrke C W, Kuo K C, Waalkes T P, & Borek E. tRNA breakdown products as markers for cancer. *Cancer*, 44(6):2120–2123 (1979). doi:10.1002/1097-0142(197912)44:6<2120::AID-CNCR2820440623>3.0.CO;2-6. URL <http://www.ncbi.nlm.nih.gov/pubmed/509391><http://doi.wiley.com/10.1002/1097-0142%28197912%2944%3A6%3C2120%3A%3AAID-CNCR2820440623%3E3.0.CO%3B2-6>
- [85] Cole C, Sobala A, Lu C, Thatcher S R, Bowman A, Brown J W, Green P J, Barton G J, & Hutvagner G. Filtering of deep sequencing data reveals the existence of abundant Dicer-dependent small RNAs derived from tRNAs. *RNA*, 15(12):2147–2160 (2009). doi:10.1261/rna.1738409. URL <http://www.ncbi.nlm.nih.gov/pubmed/19850906><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2779667><http://rnajournal.cshlp.org/cgi/doi/10.1261/rna.1738409>
- [86] Lee Y S, Shibata Y, Malhotra A, & Dutta A. A novel class of small RNAs: tRNA-derived RNA fragments (tRFs). *Genes & development*, 23(22):2639–49 (2009). doi:10.1101/gad.1837609. URL <http://www.ncbi.nlm.nih.gov/pubmed/19933153><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2779758>
- [87] Godoy P M, Bhakta N R, Barczak A J *et al.* Large Differences in Small RNA Composition Between Human Biofluids. *Cell Reports*, 25(5):1346–1358 (2018). doi:10.1016/j.celrep.2018.10.014. URL <https://doi.org/10.1016/j.celrep.2018.10.014><https://linkinghub.elsevier.com/retrieve/pii/S221124718315778>
- [88] Yamasaki S, Ivanov P, Hu G f, & Anderson P. Angiogenin cleaves tRNA and promotes stress-induced translational repression. *The Journal of Cell Biology*, 185(1):35–42 (2009). doi:10.1083/JCB.200811106. URL <http://jcb.rupress.org/content/185/1/35.long>
- [89] Ivanov P, Emara M M, Villen J, Gygi S P, & Anderson P. Angiogenin-Induced tRNA Fragments Inhibit Translation Initiation. *Molecular Cell*, 43(4):613–623 (2011). doi:10.1016/j.molcel.2011.06.022. URL <https://www.sciencedirect.com/science/article/pii/S1097276511005247?via%3Dihub><https://linkinghub.elsevier.com/retrieve/pii/S1097276511005247>

- [90] Burroughs A M, Ando Y, de Hoon M L, Tomaru Y, Suzuki H, Hayashizaki Y, & Daub C O. Deep-sequencing of human Argonaute-associated small RNAs provides insight into miRNA sorting and reveals Argonaute association with RNA fragments of diverse origin. *RNA Biology*, 8(1):158–177 (2011). doi:10.4161/rna.8.1.14300.
URL <http://www.tandfonline.com/doi/abs/10.4161/rna.8.1.14300>
- [91] Kumar P, Anaya J, Mudunuri S B, & Dutta A. Meta-analysis of tRNA derived RNA fragments reveals that they are evolutionarily conserved and associate with AGO proteins to recognize specific RNA targets. *BMC Biology*, 12(1):78 (2014). doi:10.1186/s12915-014-0078-0.
URL <http://www.ncbi.nlm.nih.gov/pubmed/25270025> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4203973> <http://bmcbiol.biomedcentral.com/articles/10.1186/s12915-014-0078-0>
- [92] Huang B, Yang H, Cheng X *et al.* tRF/miR-1280 Suppresses Stem Cell-like Cells and Metastasis in Colorectal Cancer. *Cancer Research*, 77(12):3194–3206 (2017). doi:10.1158/0008-5472.CAN-16-3146.
URL <http://cancerres.aacrjournals.org/> <http://cancerres.aacrjournals.org/lookup/doi/10.1158/0008-5472.CAN-16-3146>
- [93] Gebetsberger J, Wyss L, Mleczko A M, Reuther J, & Polacek N. A tRNA-derived fragment competes with mRNA for ribosome binding and regulates translation during stress. *RNA Biology*, 14(10):1364–1373 (2017). doi:10.1080/15476286.2016.1257470.
URL <https://www.tandfonline.com/action/journalInformation?journalCode=krnb20> <https://www.tandfonline.com/doi/full/10.1080/15476286.2016.1257470>
- [94] Goodarzi H, Liu X, Nguyen H C, Zhang S, Fish L, & Tavazoie S F. Endogenous tRNA-Derived Fragments Suppress Breast Cancer Progression via YBX1 Displacement. *Cell*, 161(4):790–802 (2015). doi:10.1016/j.cell.2015.02.053.
URL <https://www.sciencedirect.com/science/article/pii/S0092867415003189?via%3Dihub> <https://linkinghub.elsevier.com/retrieve/pii/S0092867415003189>
- [95] Kim H K, Fuchs G, Wang S *et al.* A transfer-RNA-derived small RNA regulates ribosome biogenesis. *Nature*, 552(7683):57 (2017). doi:10.1038/nature25005.
URL <http://www.nature.com/doifinder/10.1038/nature25005>
- [96] Parisien M, Wang X, & Pan T. Diversity of human tRNA genes from the 1000-genomes project. *RNA Biology*, 10(12):1853–1867 (2013). doi:10.4161/rna.27361.
URL <http://www.ncbi.nlm.nih.gov/pubmed/24448271> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3917988> <http://www.tandfonline.com/doi/abs/10.4161/rna.27361>

- [97] Loher P, Telonis A G, & Rigoutsos I. MINTmap: fast and exhaustive profiling of nuclear and mitochondrial tRNA fragments from short RNA-seq data. *Scientific Reports*, 7(1):41184 (2017). doi:10.1038/srep41184.
 URL <http://dx.doi.org/10.1038/srep41184><http://www.nature.com/articles/srep41184>
- [98] Marbach D, Lamparter D, Quon G, Kellis M, Kutalik Z, & Bergmann S. Tissue-specific regulatory circuits reveal variable modular perturbations across complex diseases. *Nature Methods*, 13(4):366–370 (2016). doi:10.1038/nmeth.3799.
 URL <http://www.ncbi.nlm.nih.gov/pubmed/26950747><http://www.nature.com/articles/nmeth.3799><http://regulatorycircuits.org>
- [99] Nowakowski T J, Rani N, Golkaram M *et al.* Regulation of cell-type-specific transcriptomes by microRNA networks during human brain development. *Nature Neuroscience*, 21(12):1784–1792 (2018). doi:10.1038/s41593-018-0265-3.
 URL <http://www.ncbi.nlm.nih.gov/pubmed/30455455><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC6312854><http://www.nature.com/articles/s41593-018-0265-3>
- [100] Londin E, Loher P, Telonis A G *et al.* Analysis of 13 cell types reveals evidence for the expression of numerous novel primate- and tissue-specific microRNAs. *Proceedings of the National Academy of Sciences*, 112(10):E1106–E1115 (2015). doi:10.1073/pnas.1420955112.
 URL <http://www.pnas.org/lookup/doi/10.1073/pnas.1420955112>
- [101] Dweep H & Gretz N. miRWalk2.0: a comprehensive atlas of microRNA-target interactions. *Nature Methods*, 12(8):697–697 (2015). doi:10.1038/nmeth.3485.
 URL <http://www.nature.com/doifinder/10.1038/nmeth.3485>
- [102] Chaudhuri S, Narasayya V, & Ramamurthy R. Estimating Progress of Execution for SQL Queries (2004).
 URL <https://www.microsoft.com/en-us/research/publication/estimating-progress-of-execution-for-sql-queries/?from=https%3A%2F%2Fresearch.microsoft.com%2Fapps%2Fpubs%2F%3Fid%3D76556>
- [103] Hon C c, Ramiłowski J A, Harshbarger J *et al.* An atlas of human long non-coding RNAs with accurate 5' ends. *Nature Publishing Group*, 543(7644):199–204 (2017). doi:10.1038/nature21374.
 URL <http://dx.doi.org/10.1038/nature21374>
- [104] Dweep H & Gretz N. miRWalk2 web page.
- [105] Karagkouni D, Paraskevopoulou M D, Chatzopoulos S *et al.* DIANA-TarBase v8: a decade-long collection of experimentally supported miRNA–gene interactions. *Nucleic Acids Re-*

- search*, 46(D1):D239–D245 (2018). doi:10.1093/nar/gkx1141.
URL <http://academic.oup.com/nar/article/46/D1/D239/4634010>
- [106] Chou C H, Shrestha S, Yang C D *et al.* miRTarBase update 2018: a resource for experimentally validated microRNA-target interactions. *Nucleic Acids Research*, 46(D1):D296–D302 (2018). doi:10.1093/nar/gkx1067.
URL <http://www.ncbi.nlm.nih.gov/pubmed/29126174> <http://www.ncbi.nlm.nih.gov/pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5753222> <http://academic.oup.com/nar/article/46/D1/D296/4595852>
- [107] Yue D, Liu H, & Huang Y. Survey of Computational Algorithms for MicroRNA Target Prediction. *Current Genomics*, 10(7):478–492 (2009). doi:10.2174/138920209789208219.
URL <http://www.ncbi.nlm.nih.gov/pubmed/20436875> <http://www.ncbi.nlm.nih.gov/pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2808675> <http://www.eurekaselect.com/openurl/content.php?genre=article&issn=1389-2029&volume=10&issue=7&spage=478>
- [108] Witkos T M, Koscianska E, & Krzyzosiak W J. Practical Aspects of microRNA Target Prediction. *Current molecular medicine*, 11(2):93–109 (2011). doi:10.2174/156652411794859250.
- [109] Friedman R C, Farh K K H, Burge C B, & Bartel D P. Most mammalian mRNAs are conserved targets of microRNAs. *Genome Research*, 19(1):92–105 (2009). doi:10.1101/gr.082701.108.
URL <http://www.ncbi.nlm.nih.gov/pubmed/18955434> <http://www.ncbi.nlm.nih.gov/pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2612969> <http://genome.cshlp.org/cgi/doi/10.1101/gr.082701.108>
- [110] Alexiou P, Maragkakis M, Papadopoulos G L, Reczko M, & Hatzigeorgiou A G. Lost in translation: an assessment and perspective for computational microRNA target identification. *Bioinformatics*, 25(23):3049–3055 (2009). doi:10.1093/bioinformatics/btp565.
URL <http://www.ncbi.nlm.nih.gov/pubmed/19789267> <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btp565>
- [111] Agarwal V, Bell G W, Nam J W, & Bartel D P. Predicting effective microRNA target sites in mammalian mRNAs. *eLife*, 4 (2015). doi:10.7554/eLife.05005.
URL <https://elifesciences.org/articles/05005>
- [112] Soreq H. Checks and balances on cholinergic signaling in brain and body function. *Trends in Neurosciences*, 38(7):448–458 (2015). doi:10.1016/j.tins.2015.05.007.
URL <http://dx.doi.org/10.1016/j.tins.2015.05.007>
- [113] Smith T & Waterman M. Identification of common molecular subsequences. *Journal of Molecular Biology*, 147(1):195–197 (1981). doi:10.1016/0022-2836(81)90087-5.

- URL [https://www.sciencedirect.com/science/article/pii/0022283681900875?
via%}3Dihub](https://www.sciencedirect.com/science/article/pii/0022283681900875?via%}3Dihub)<https://linkinghub.elsevier.com/retrieve/pii/0022283681900875>
- [114] Needleman S B & Wunsch C D. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48(3):443–453 (1970). doi:10.1016/0022-2836(70)90057-4.
- URL [https://www.sciencedirect.com/science/article/pii/0022283670900574?
via%}3Dihub](https://www.sciencedirect.com/science/article/pii/0022283670900574?via%}3Dihub)
- [115] Raichle M E & Gusnard D A. Appraising the brain's energy budget. *Proceedings of the National Academy of Sciences*, 99(16):10237–10239 (2002). doi:10.1073/pnas.172399499.
- URL <http://www.pnas.org/cgi/doi/10.1073/pnas.172399499>
- [116] Bohn K A, Adkins C E, Mittapalli R K, Terrell-Hall T B, Mohammad A S, Shah N, Dolan E L, Nounou M I, & Lockman P R. Semi-automated rapid quantification of brain vessel density utilizing fluorescent microscopy. *Journal of Neuroscience Methods*, 270:124–131 (2016). doi:10.1016/j.jneumeth.2016.06.012.
- URL <http://www.ncbi.nlm.nih.gov/pubmed/27321229><http://www.ncbi.nlm.nih.gov/entrez/efetch.fcgi?db=PubMed&id=PMC4981522><https://linkinghub.elsevier.com/retrieve/pii/S0165027016301339>
- [117] Lobentanzer S & Klein J. Zentrales und Peripheres Nervensystem. In Wichmann & Fromme, (Editors) *Handbuch für Umweltmedizin*, chapter 11. ecomed Medizin, erg. lfg. edition (2019).
- [118] Darmanis S, Sloan S A, Zhang Y, Enge M, Caneda C, Shuer L M, Hayden Gephart M G, Barres B A, & Quake S R. A survey of human brain transcriptome diversity at the single cell level. *Proceedings of the National Academy of Sciences*, 112(23):201507125 (2015). doi:10.1073/pnas.1507125112.
- URL <http://www.pnas.org/content/112/23/7285.abstract>
- [119] Zeisel a, Manchado a B M, Codeluppi S *et al.* Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science*, 347(6226):1138–42 (2015). doi:10.1126/science.aaa1934.
- URL <http://science.sciencemag.org/docelec.univ-lyon1.fr/content/347/6226/1138.abstract>
- [120] Tasic B, Menon V, Nguyen T N T *et al.* Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nature Neuroscience*, advance on(January):1–37 (2016). doi:10.1038/nn.4216.
- URL <http://dx.doi.org/10.1038/nn.4216>
- [121] Habib N, Li Y, Heidenreich M, Swiech L, Avraham-David I, Trombetta JJ, Hession C, Zhang F, & Regev A. Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn

- neurons. *Science*, 353(6302):925–928 (2016). doi:10.1126/science.aad7038.
URL <http://wwwsciencemag.org/lookup/doi/10.1126/science.aad7038>
- [122] Zeisel A, Hochgerner H, Lönnerberg P *et al.* Molecular Architecture of the Mouse Nervous System. *Cell*, 174(4):999–1014.e22 (2018). doi:10.1016/j.cell.2018.06.021.
URL <http://www.ncbi.nlm.nih.gov/pubmed/30096314><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC6086934><https://linkinghub.elsevier.com/retrieve/pii/S009286741830789X>
- [123] Murtagh F & Legendre P. Ward’s Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward’s Criterion? *Journal of Classification*, 31(3):274–295 (2014). doi:10.1007/s00357-014-9161-z.
URL <http://link.springer.com/10.1007/s00357-014-9161-z>
- [124] Bray J R & Curtis J T. An Ordination of the Upland Forest Communities of Southern Wisconsin. *Ecological Monographs*, 27(4):325–349 (1957). doi:10.2307/1942268.
URL <http://doi.wiley.com/10.2307/1942268>
- [125] Biedler J L, Roffler-Tarlov S, Schachner M, & Freedman L S. Multiple Neurotransmitter Synthesis by Human Neuroblastoma Cell Lines and Clones. *Cancer Res.*, 38(11_Part_1):3751–3757 (1978).
URL <http://cancerres.aacrjournals.org/content/38/11{ }Part{ }1/3751.short>
- [126] Biedler J L, Helson L, & Spengler B A. Morphology and Growth, Tumorigenicity, and Cytogenetics of Human Neuroblastoma Cells in Continuous Culture. *Cancer Research*, 33(11):2643–2652 (1973).
- [127] Seeger R C, Rayner S A, Laug W E, Neustein H B, & Benedict W F. Morphology, Growth, Chromosomal Pattern, and Fibrinolytic Activity of two New Human Neuroblastoma Cell Lines. *Cancer Research*, 37(5):1364–1371. (1977).
- [128] Seeger R C, Danon Y L, Rayner S A, & Hoover F. Definition of a Thy-1 determinant on human neuroblastoma, glioma, sarcoma, and teratoma cells with a monoclonal antibody. *Journal of immunology (Baltimore, Md. : 1950)*, 128(2):983–9 (1982).
URL <http://www.ncbi.nlm.nih.gov/pubmed/6172518>
- [129] Hill D P & Robertson K a. Characterization of the cholinergic neuronal differentiation of the human neuroblastoma cell line LA-N-5 after treatment with retinoic acid. *Developmental Brain Research*, 102(1):53–67 (1997). doi:10.1016/S0165-3806(97)00076-X.
URL <http://www.ncbi.nlm.nih.gov/pubmed/9298234><https://linkinghub.elsevier.com/retrieve/pii/S016538069700076X>

- [130] McManaman J L & Crawford F G. Skeletal Muscle Proteins Stimulate Cholinergic Differentiation of Human Neuroblastoma Cells. *Journal of neurochemistry*, pp. 258–266 (1991).
- [131] Sun M, Liu H, Min S, Wang H, & Wang X. Ciliary neurotrophic factor-treated astrocyte-conditioned medium increases the intracellular free calcium concentration in rat cortical neurons. *Biomedical Reports*, 4(4):417–420 (2016). doi:10.3892/br.2016.602.
URL <https://www.spandidos-publications.com/https://www.spandidos-publications.com/10.3892/br.2016.602>
- [132] Roehr J T, Dieterich C, & Reinert K. Flexbar 3.0 – SIMD and multicore parallelization. *Bioinformatics*, 33(18):2941–2942 (2017). doi:10.1093/bioinformatics/btx330.
URL <https://academic.oup.com/bioinformatics/article/33/18/2941/3852078>
- [133] Wang W C, Lin F M, Chang W C, Lin K Y, Huang H D, & Lin N S. miRExpress: analyzing high-throughput sequencing data for profiling microRNA expression. *BMC Bioinformatics*, 10:328 (2009). doi:10.1186/1471-2105-10-328.
URL <https://www.ncbi.nlm.nih.gov/pubmed/19821977>
- [134] Love M I, Huber W, & Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology*, 15(12):550 (2014). doi:10.1186/s13059-014-0550-8.
URL <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0550-8>
- [135] Wald A. Contributions to the Theory of Statistical Estimation and Testing Hypotheses. *The Annals of Mathematical Statistics*, 10(4):299–326 (1939). doi:10.1098/rsta.1937.0005.
URL <http://rsta.royalsocietypublishing.org/cgi/doi/10.1098/rsta.1937.0005>
- [136] Bullard J H, Purdom E, Hansen K D, & Dudoit S. Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics*, 11(1):94 (2010). doi:10.1186/1471-2105-11-94.
URL <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-11-94>
- [137] Chen Z, Liu J, Ng H, Nadarajah S, Kaufman H L, Yang J Y, & Deng Y. Statistical methods on detecting differentially expressed genes for RNA-seq data. *BMC Systems Biology*, 5(Suppl 3):S1 (2011). doi:10.1186/1752-0509-5-S3-S1.
URL <http://www.ncbi.nlm.nih.gov/pubmed/22784615>
<http://www.ncbi.nlm.nih.gov/pubmed/22784615>
<http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3287564>
<http://bmcsystbiol.biomedcentral.com/articles/10.1186/1752-0509-5-S3-S1>
- [138] Zhu A, Ibrahim J G, & Love M I. Heavy-tailed prior distributions for sequence count data: removing the noise and preserving large differences. *Bioinformatics*, 35(12):2084–2092 (2019).

- doi:10.1093/bioinformatics/bty895.
URL <https://academic.oup.com/bioinformatics/article/35/12/2084/5159452>
- [139] Alexa A, Rahnenfuhrer J, & Lengauer T. Improved scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics*, 22(13):1600–1607 (2006). doi:10.1093/bioinformatics/btl140.
URL <http://www.ncbi.nlm.nih.gov/pubmed/16606683> <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btl140>
- [140] Jacomy M, Venturini T, Heymann S, & Bastian M. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. *PLoS ONE*, 9(6):1–12 (2014). doi:10.1371/journal.pone.0098679.
- [141] Chen C, Cheng L, Grennan K, Pibiri F, Zhang C, Badner J A, Members of the Bipolar Disorder Genome Study (BiGS) Consortium, Gershon E S, & Liu C. Two gene co-expression modules differentiate psychotics and controls. *Molecular psychiatry*, 18(12):1308–14 (2013). doi:10.1038/mp.2012.146.
URL <http://www.ncbi.nlm.nih.gov/pubmed/23147385> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4018461>
- [142] Iwamoto K, Bundo M, & Kato T. Altered expression of mitochondria-related genes in postmortem brains of patients with bipolar disorder or schizophrenia, as revealed by large-scale DNA microarray analysis. *Human molecular genetics*, 14(2):241–53 (2005). doi:10.1093/hmg/ddi022.
URL <http://www.ncbi.nlm.nih.gov/pubmed/15563509>
- [143] Lanz T A, Joshi J J, Reinhart V, Johnson K, Grantham L E, & Volkson D. STEP levels are unchanged in pre-frontal cortex and associative striatum in post-mortem human brain samples from subjects with schizophrenia, bipolar disorder and major depressive disorder. *PloS one*, 10(3):e0121744 (2015). doi:10.1371/journal.pone.0121744.
URL <http://www.ncbi.nlm.nih.gov/pubmed/25786133> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4364624>
- [144] Maycox P R, Kelly F, Taylor A *et al.* Analysis of gene expression in two large schizophrenia cohorts identifies multiple changes associated with nerve terminal function. *Molecular psychiatry*, 14(12):1083–94 (2009). doi:10.1038/mp.2009.18.
URL <http://www.ncbi.nlm.nih.gov/pubmed/19255580>
- [145] Narayan S, Tang B, Head S R, Gilman T J, Sutcliffe J G, Dean B, & Thomas E A. Molecular profiles of schizophrenia in the CNS at different stages of illness. *Brain research*, 1239:235–48 (2008). doi:10.1016/j.brainres.2008.08.023.

- URL <http://www.ncbi.nlm.nih.gov/pubmed/18778695><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC2783475>
- [146] Ryan M M, Lockstone H E, Huffaker S J, Wayland M T, Webster M J, & Bahn S. Gene expression analysis of bipolar disorder reveals downregulation of the ubiquitin cycle and alterations in synaptic genes. *Molecular psychiatry*, 11(10):965–78 (2006). doi:10.1038/sj.mp.4001875.
URL <http://www.ncbi.nlm.nih.gov/pubmed/16894394>
- [147] Ramaker R C, Bowling K M, Lasseigne B N *et al.* Post-mortem molecular profiling of three psychiatric disorders. *Genome Medicine*, 9(1):1–12 (2017). doi:10.1186/s13073-017-0458-5.
- [148] Hoffman G E, Hartley B J, Flaherty E, Ladran I, Gochman P, Ruderfer D M, Stahl E A, Rapoport J, Sklar P, & Brennand K J. Transcriptional signatures of schizophrenia in hiPSC-derived NPCs and neurons are concordant with post-mortem adult brains. *Nature Communications*, 8(1) (2017). doi:10.1038/s41467-017-02330-5.
URL <http://dx.doi.org/10.1038/s41467-017-02330-5>
- [149] Webb A, Papp A C, Curtis A *et al.* RNA sequencing of transcriptomes in human brain regions: Protein-coding and non-coding RNAs, isoforms and alleles. *BMC Genomics*, 16(1):1–16 (2015). doi:10.1186/s12864-015-2207-8.
URL <http://dx.doi.org/10.1186/s12864-015-2207-8>
- [150] Fontenot M R, Berto S, Liu Y, Werthmann G, Douglas C, Usu N, Gleason K, Tamminga C A, Takahashi J S, & Konopka G. Novel transcriptional networks regulated by clock in human neurons. *Genes and Development*, 31(21):2121–2135 (2017). doi:10.1101/gad.305813.117.
- [151] Li G, Klein J, & Zimmermann M. Pathophysiological Amyloid Concentrations Induce Sustained Upregulation of Readthrough AChE Mediating Anti-Apoptotic Effects. *Neuroscience*, 240:349–60 (2013). doi:10.1016/j.neuroscience.2013.02.040.
URL <http://www.ncbi.nlm.nih.gov/pubmed/23485809>
- [152] Gulyás-Kovács A, Keydar I, Xia E, Fromer M, Hoffman G, Ruderfer D, Sachidanandam R, & Chess A. Unperturbed expression bias of imprinted genes in schizophrenia. *Nature Communications*, 9(1):2914 (2018). doi:10.1038/s41467-018-04960-9.
URL <https://www.biorxiv.org/content/early/2018/05/24/329748?%3Fcollection=http://www.nature.com/articles/s41467-018-04960-9>
- [153] Du P, Kibbe W A, & Lin S M. lumi: a pipeline for processing Illumina microarray. *Bioinformatics*, 24(13):1547–1548 (2008). doi:10.1093/bioinformatics/btn224.
URL <http://www.ncbi.nlm.nih.gov/pubmed/18467348><https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btn224>

- [154] Gautier L, Cope L, Bolstad B M, & Irizarry R A. affy—analysis of Affymetrix GeneChip data at the probe level. *Bioinformatics*, 20(3):307–315 (2004). doi:10.1093/bioinformatics/btg405.
 URL <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btg405>
- [155] Oldham M C, Langfelder P, & Horvath S. Network methods for describing sample relationships in genomic datasets: application to Huntington’s disease. *BMC Systems Biology*, 6(1):63 (2012). doi:10.1186/1752-0509-6-63.
 URL <http://www.ncbi.nlm.nih.gov/pubmed/22691535><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3441531><http://bmcsystbiol.biomedcentral.com/articles/10.1186/1752-0509-6-63>
- [156] Durinck S, Spellman P T, Birney E, & Huber W. Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt. *Nature Protocols*, 4(8):1184–1191 (2009). doi:10.1038/nprot.2009.97.
 URL <http://www.ncbi.nlm.nih.gov/pubmed/19617889><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3159387><http://www.nature.com/articles/nprot.2009.97>
- [157] Langfelder P & Horvath S. WGCNA: An R package for weighted correlation network analysis. *BMC Bioinformatics*, 9 (2008). doi:10.1186/1471-2105-9-559.
- [158] Pinheiro J, Bates D, DebRoy S, Sarkar D, & R Core Team. *nlme: Linear and Nonlinear Mixed Effects Models* (2019). R package version 3.1-142.
 URL <https://CRAN.R-project.org/package=nlme>
- [159] Shaltiel G, Hanan M, Wolf Y, Barbash S, Kovalev E, Shoham S, & Soreq H. Hippocampal microRNA-132 mediates stress-inducible cognitive deficits through its acetylcholinesterase target. *Brain structure function*, 218(1):1–5 (2013). doi:10.1007/s00429-011-0376-z.
 URL <http://www.ncbi.nlm.nih.gov/pubmed/22246100>
- [160] Hanin G, Yayon N, Tzur Y *et al.* miRNA-132 induces hepatic steatosis and hyperlipidaemia by synergistic multitarget suppression. *Gut*, 67(6):1124–1134 (2018). doi:10.1136/gutjnl-2016-312869.
 URL <http://gut.bmjjournals.org/lookup/doi/10.1136/gutjnl-2016-312869><http://www.ncbi.nlm.nih.gov/pubmed/28381526><http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC5969364>
- [161] Mellios N, Sugihara H, Castro J *et al.* miR-132, an experience-dependent microRNA, is essential for visual cortex plasticity. *Nature Neuroscience*, 14(10):1240–1242 (2011). doi:10.1038/nn.2909.
 URL <http://www.nature.com/articles/nn.2909>

- [162] Shaked I, Meerson A, Wolf Y, Avni R, Greenberg D, Gilboa-Geffen A, & Soreq H. MicroRNA-132 potentiates cholinergic anti-inflammatory signaling by targeting acetylcholinesterase. *Immunity*, 31(6):965–973 (2009). doi:10.1016/j.jimmuni.2009.09.019.
URL <http://www.ncbi.nlm.nih.gov/pubmed/20005135>
- [163] Pichler S, Gu W, Hartl D, Gasparoni G, Leidinger P, Keller A, Meese E, Mayhaus M, Hampel H, & Riemenschneider M. The miRNome of Alzheimer's disease: consistent downregulation of the miR-132/212 cluster. *Neurobiology of Aging*, 50:167.e1–167.e10 (2017). doi:10.1016/j.neurobiolaging.2016.09.019.
URL <https://linkinghub.elsevier.com/retrieve/pii/S0197458016302330>
- [164] Busch S, Auth E, Scholl F, Huenecke S, Koehl U, Suess B, & Steinhilber D. 5-Lipoxygenase Is a Direct Target of miR-19a-3p and miR-125b-5p. *The Journal of Immunology*, 194(4):1646–1653 (2015). doi:10.4049/jimmunol.1402163.
URL <http://www.jimmunol.org/lookup/doi/10.4049/jimmunol.1402163>
- [165] Zhang J, Qu P, Zhou C, Liu X, Ma X, Wang M, Wang Y, Su J, Liu J, & Zhang Y. MicroRNA-125b is a key epigenetic regulatory factor that promotes nuclear transfer reprogramming. *Journal of Biological Chemistry*, 292(38):15916–15926 (2017). doi:10.1074/jbc.M117.796771.
- [166] Soreq H, Ben-Aziz R, Prody C A, Seidman S, Gnatt A, Neville L, Lieman-Hurwitz J, Lev-Lehman E, Ginzberg D, & Lipidot-Lifson Y. Molecular cloning and construction of the coding region for human acetylcholinesterase reveals a G + C-rich attenuating structure. *Proceedings of the National Academy of Sciences*, 87(24):9688–9692 (1990). doi:10.1073/pnas.87.24.9688.
URL <http://www.pnas.org/cgi/doi/10.1073/pnas.87.24.9688>
- [167] Hanin G, Shenhar-Tsarfaty S, Yayon N *et al.* Competing targets of microRNA-608 affect anxiety and hypertension. *Human Molecular Genetics*, 23(17):4569–4580 (2014). doi:10.1093/hmg/ddu170.
URL <https://academic.oup.com/hmg/article-lookup/doi/10.1093/hmg/ddu170> <http://www.ncbi.nlm.nih.gov/pubmed/24722204> <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?artid=PMC4119407>
- [168] Hoffmann S, Harms H, Ulm L *et al.* Stroke-induced immunodepression and dysphagia independently predict stroke-associated pneumonia – The PREDICT study. *Journal of Cerebral Blood Flow & Metabolism*, 37(12):3671–3682 (2017). doi:10.1177/0271678X16671964.
URL <https://www.ncbi.nlm.nih.gov/pubmed/27733675> <http://journals.sagepub.com/doi/10.1177/0271678X16671964>

- [169] Patro R, Duggal G, Love M I, Irizarry R A, & Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. *Nature Methods*, 14(4):417–419 (2017). doi: 10.1038/nmeth.4197.
URL <http://www.nature.com/articles/nmeth.4197>
- [170] van der Maaten L & Hinton G. Visualizing Data using t-SNE. *Journal of Machine Learning research*, 9(1) (2008).
- [171] Juzenas S, Venkatesh G, Hübenthal M *et al.* A comprehensive, cell specific microRNA catalogue of human peripheral blood. *Nucleic Acids Research*, 45(16):9290–9301 (2017). doi: 10.1093/nar/gkx706.
- [172] Broido A D & Clauset A. Scale-free networks are rare. *Nature Communications*, 10(1):1017 (2019). doi:10.1038/s41467-019-08746-5.
URL <http://www.nature.com/articles/s41467-019-08746-5>

A

Transcription Factor Regulatory Circuits - Tissue Types

B

List of Primate-Specific Homologues of Human microRNAs

C

microRNA Differential Expression in LA-N-2 and LA-N-5

D

List of GO Terms from Analysis of
Differentially Expressed Large RNA in
Stroke