

METHOD

Single-cell gene regulatory network inference at scale: The Inferelator 3.0

Claudia Skok Gibbs^{1†}, Christopher A Jackson^{2,3†}, Giuseppe-Antonio Saldi^{2,3†}, Aashna Shah¹, Andreas Tjärnberg^{2,3}, Aaron Watters¹, Nicholas De Veaux¹, Konstantine Tchourine⁶, Ren Yi⁴, Tymor Hamamsy⁵, Dayanne M Castro^{2,3}, Nicholas Carriero⁷, David Gresham^{2,3}, Emily R Miraldi^{8,9} and Richard Bonneau^{1,2,3,4,5*}

*Correspondence: rb133@nyu.edu

² Center For Genomics and Systems Biology, NYU, 10008, New York, USA

Full list of author information is available at the end of the article

[†]These authors contributed equally and are listed alphabetically by last name. They have agreed to change the order of equal contributing authors to list themselves first as a convenience when necessary.

Abstract

Gene regulatory networks define regulatory relationships between transcription factors and target genes within a biological system, and reconstructing them is essential for understanding cellular growth and function. In this work, we present the Inferelator 3.0, which has been significantly updated to integrate data from distinct cell types to learn context-specific regulatory networks and aggregate them into a shared regulatory network, while retaining the functionality of the previous versions. The Inferelator 3.0 reliably learns informative networks from the model organisms *Bacillus subtilis* and *Saccharomyces cerevisiae*. We demonstrate its capabilities by learning networks for multiple distinct neuronal and glial cell types in the developing *Mus musculus* brain at E18 from a large (1.3 million) single-cell gene expression dataset with paired single-cell chromatin accessibility data.

Keywords: Gene Regulation; Network Inference; Transcription Factors; Transcription Factor Activity

1 Background

Gene expression is tightly regulated at multiple levels in order to control growth, development, and response to environmental conditions (Figure 1A). Transcriptional regulation is principally controlled by Transcription Factors (TFs) that bind to DNA and effect chromatin remodeling [1] or directly modulate the output of RNA polymerases [2]. Three percent of *Saccharomyces cerevisiae* genes are TFs [3], and more than six percent of human genes are believed to be TFs or cofactors [4]. Connections between TFs and genes combine to form a transcriptional Gene Regulatory Network (GRN) that can be represented as a directed graph (Figure 1B). Learning the true regulatory network that connects regulatory TFs to target genes is a key problem in biology [5, 6]. Determining the valid GRN is necessary to explain how mutations that cause gene dysregulation lead to complex disease states [7], how variation at the genetic level leads to selectable phenotypic variation [8, 9], and how to re-engineer organisms to efficiently produce industrial chemicals and enzymes [10].

Learning genome-scale networks relies on genome-wide expression measurements, initially captured with microarray technology [11], but today typically measured by RNA-sequencing (RNA-seq) [12, 13]. A major difficulty is that biological systems

19 have large numbers of both regulators and targets; there is poor network identifi-
20 ability because many plausible networks can explain observed expression data and
21 the regulation of gene expression in an organism [14]. Designing experiments to
22 produce data that increases network identifiability is possible [15], but most data is
23 collected for specific projects and repurposed for network inference as a consequence
24 of the cost of data collection. Large-scale experiments in which a perturbation is
25 made and dynamic data is collected over time is exceptionally useful for learning
26 GRNs but systematic studies that collect this data are rare [16].

27 Measuring the expression of single cells using single-cell RNA-sequencing (scRNA-
28 seq) is an emerging and highly scalable technology. Microfluidic-based single-cell
29 techniques [17, 18, 19] allow for thousands of measurements in a single experiment.
30 Split-pool barcoding techniques [20] are poised to increase single-cell throughput
31 by an order of magnitude. These techniques have been successfully applied to gener-
32 ate multiplexed gene expression data from pools of barcoded cell lines with loss-
33 of-function TF mutants [21, 22], enhancer perturbations [23], and disease-causing
34 oncogene variants [24]. Individual cell measurements are sparser and noisier than
35 measurements generated using traditional RNA-seq, although in aggregate the gene
36 expression profiles of single-cell data match RNA-seq data well [25], and techniques
37 to denoise single-cell data have been developed [26, 27].

38 The *seurat* [28] and *scanpy* [29] bioinformatics toolkits are established tools for
39 single-cell data analysis, but pipelines for inferring GRNs from single-cell data are
40 still nascent. The SCENIC [30] GRN inference pipeline is based around the GENIE3
41 method that uses ensemble regression trees [31] to estimate the importance of TFs
42 in explaining gene expression profiles. CellOracle [32] has been recently proposed
43 as a pipeline to integrate single-cell ATAC and expression data using a motif-based
44 search for potential regulators, followed by Bayesian ridge regression to enforce
45 sparsity in the output GRN. SCODE [33] infers GRNs by solving linear ordinary
46 differential equations using time-course single-cell data. GRN inference is compu-
47 tationally challenging, and the most scalable of these GRN pipelines has learned
48 GRNs from 50,000 cells of gene expression data [30].

49 Here we describe the Inferelator 3.0 pipeline for single-cell GRN inference, based
50 on regularized linear regression [34]. The Inferelator 2.0 [35] has performed well
51 inferring networks from *Bacillus subtilis* [36], human Th17 cells [37, 38], mouse
52 Lymphocytes [39], *Saccharomyces cerevisiae* [40], and *Oryza sativa* [41]. We have
53 re-implemented the Inferelator 3.0 with new functionality in python to learn GRNs
54 from single-cell gene expression data. Specifically, this new package provides scala-
55 bility, allowing millions of cells to be analyzed together, as well as integrated support
56 for multi-task GRN inference, while retaining the ability to utilize bulk gene ex-
57 pression data. As a demonstration of these capabilities, we learn GRNs from several
58 model organisms, and generate a mouse neuronal GRN from a publicly available
59 dataset containing 1.3 million cells.

60 **2 Results**

61 **2.1 The Inferelator 3.0**

62 In the 11 years since the last major release of the Inferelator [35], the scale of data
63 collection in biology has accelerated enormously. We have therefore rewritten the

64 Inferelator as a python package to take advantage of the concurrent advances in data
65 processing. For inference from small scale gene expression datasets ($< 10^4$ observa-
66 tions), the Inferelator 3.0 uses native python multiprocessing to run on individual
67 computers. For inference from extremely large scale gene expression datasets ($10^4 -$
68 10^7 observations) that are increasingly available from scRNA-seq experiments, the
69 Inferelator 3.0 takes advantage of the Dask analytic engine [42] for deployment to
70 high-performance clusters (Figure 1C), or for deployment as a kubernetes image to
71 the Google cloud computing infrastructure.

72 2.2 Network Inference using Bulk RNA-Seq Expression Data

73 We have incorporated several network inference model selection methods into the
74 Inferelator (Figure 2A). In order to evaluate the network inference performance of
75 these methods, we test on the prokaryotic model *Bacillus subtilis* and the eukaryotic
76 model *Saccharomyces cerevisiae*. Both *B. subtilis* [36, 43] and *S. cerevisiae* [40, 16]
77 have large bulk RNA-seq and microarray gene expression datasets, in addition to a
78 relatively large number of experimentally determined TF-target gene interactions
79 that can be used as a gold standard for assessing network inference.

80 Using two independent datasets for each organism, we find that the model se-
81 lection methods Bayesian Best Subset Regression (BBSR) [44] and Stability Ap-
82 proach to Regularization Selection for Least Absolute Shrinkage and Selection Oper-
83 ator (StARS-LASSO) [38] perform equivalently (Figure 2B). The datasets separate
84 on the first principal component (Supplemental Figure 1A), indicating that there
85 are substantial batch-specific effects between these independent datasets. These
86 dataset-specific batch effects make combining this data for network inference diffi-
87 cult; conceptually, each dataset is in a separate space, and must be mapped into a
88 shared space if they are to be combined. We take a different approach to address-
89 ing the batch effects between datasets by treating them as separate learning tasks
90 [45] and then combining network information into a unified GRN. This results in
91 a considerable improvement in network inference performance over either dataset
92 individually (Figure 2C). The best performance is obtained with Adaptive Mul-
93 tiple Sparse Regression (AMuSR) [45], a multi-task learning method that shares
94 information between tasks during regression. The GRN learned with AMuSR ex-
95 plains the variance in the expression data better than learning networks from each
96 dataset individually with BBSR or StARS-LASSO and then combining them (Sup-
97 plemental Figure 1B). There is a high overlap in the number of GRN edges learned
98 from each dataset, showing that there is a common network across different tasks
99 (Supplemental Figure 1C).

100 2.3 Generating Prior Networks from Chromatin Data and Transcription Factor Motifs

101 The Inferelator produces an inferred network from a combination of gene expression
102 data and a prior network constructed from existing knowledge about known gene
103 regulation. Curated databases of regulator-gene interactions culled from domain-
104 specific literature are an excellent source for prior networks. While some model
105 systems have excellent databases of known interactions, these resources are unavail-
106 able for most organisms or cell types. In these cases, using chromatin accessibility
107 determined by a standard Assay for Transposase-Accessible Chromatin (ATAC) in

108 combination with the known DNA-binding preferences for TFs to identify putative
109 target genes is a viable alternative [38].

110 To generate these prior networks we have developed the inferelator-prior acces-
111 sory package that uses TF motif position-weight matrices to score TF binding within
112 gene regulatory regions and build sparse prior networks (Figure 3A). These gene
113 regulatory regions can be identified by ATAC, by existing knowledge from TF Chro-
114 matin Immunoprecipitation (ChIP) experiments, or from known databases (e.g. EN-
115 CODE [46]). Here, we compare the inferelator-prior tool to the CellOracle package
116 [32] that also constructs motif-based networks that can be constrained to regula-
117 tory regions, in *Saccharomyces cerevisiae* by using sequences 200bp upstream and
118 50bp downstream of each gene TSS as the gene regulatory region. The inferelator-
119 prior and CellOracle methods produce networks that are similar when measured by
120 Jaccard index but are dissimilar to the YEASTRACT literature-derived network
121 (Figure 3B). These motif-derived prior networks from both the inferelator-prior and
122 CellOracle methods perform well as prior knowledge for GRN inference using the
123 Inferelator pipeline (Figure 3C). The source of the motif library has a significant ef-
124 fect on network output, as can be seen with the well-characterized TF GAL4. GAL4
125 has a canonical $CGGN_{11}CGG$ binding site; different motif libraries have different
126 annotated binding sites (Supplemental Figure 2A) and yield different motif-derived
127 networks with the inferelator-prior pipeline (Supplemental Figure 2B-C).

128 2.4 Network Inference using Single-Cell Expression Data

129 Single-cell data is undersampled and noisy, but large numbers of observations are
130 collected in parallel and count data derived from unique molecular identifiers have
131 some intrinsic advantages. In order to quantitatively evaluate network inference per-
132 formance, we apply the Inferelator to *Saccharomyces cerevisiae* single-cell expression
133 data [22, 47], and score the model performance based on a previously-defined yeast
134 gold standard [40]. This data is split into 15 separate tasks, based on labels that
135 correspond to experimental conditions from the original works (Figure 4A). A net-
136 work is learned for each task separately using the YEASTRACT literature-derived
137 prior, and aggregated into a final network for scoring on held-out genes from the
138 gold standard. We test a combination of several preprocessing options with three
139 network inference model selection methods (Figure 4B-D).

140 We find that network inference is generally sensitive to the preprocessing op-
141 tions chosen, and that these differences due to preprocessing generally outweigh the
142 differences between different model selection methods (Figure 4B-D). A standard
143 Freeman-Tukey or \log_2 pseudocount transformation on raw count data yields the
144 best performance, with notable decreases in recovery of the gold standard when
145 count data is depth-normalized (such that each cell has the same total transcript
146 counts). The performance of the randomly generated Noise control (N) is higher
147 than the performance of the shuffled (S) control when counts per cell are not nor-
148 malized, suggesting that total counts per cell provides additional information during
149 inference.

150 Different model performance metrics, like AUPR, Matthews Correlation Coef-
151 ficient (MCC), and F1 score correlate very well and identify the same optimal
152 hyperparameters (Supplemental Figure 3). We apply StARS-LASSO to data that

153 has been Freeman-Tukey transformed to generate a final network without holding
154 out genes for cross-validation (Figure 4E). While we use AUPR as a metric for
155 evaluating model performance, selecting a threshold for including edges in a GRN
156 by precision or recall requires a target precision or recall to be chosen arbitrarily.
157 Alternatively, MCC and F1 score allow a threshold to be determined that maxi-
158 mizes similarity to the known prior or gold standard GRN. Choosing the Inferelator
159 confidence score threshold to include edges in a final network that maximizes MCC
160 is a simple heuristic way to select the size of a learned network that maximizes the
161 overlap of the gold standard while minimizing links not in the gold standard (Fig-
162 ure 4F and section 5.6). Using the maximum F1 score is also an option, but gives
163 a less conservative GRN as true negatives are not considered and will not diminish
164 the score. Both metrics balance similarity to the gold standard with overall network
165 size, and therefore represent straightforward heuristics that do not rely on arbitrary
166 thresholds.

167 2.5 Large-scale Single-Cell Mouse Neuron Network Inference

168 To show scalability, we apply the Inferelator to a large-scale (1.3 million cells of
169 scRNA-seq expression data) publicly available dataset of mouse brain cells (10x ge-
170 nomics) that is accompanied by 15,000 single-cell ATAC (scATAC) measurements.
171 By using Dask to parallelize network inference, we are able to distribute work across
172 multiple computational nodes, allowing networks to be rapidly learned from $\sim 10^5$
173 cells (Figure 4A). We separate the expression and scATAC data into broad cate-
174 gories; Excitatory neurons, Interneurons, Glial cells and Vascular cells (Figure 5A-
175 E). After initial quality control, filtering, and cell type assignment, 766,402 scRNA-
176 seq and 7,751 scATAC observations remain (Figure 5F, Supplemental Figure 4B-D).
177 scRNA-seq data is further clustered within broad categories into clusters (Figure
178 5B) that are assigned to specific cell types based on marker expression (Figure
179 5C, Supplemental Figure 5). scATAC data is aggregated into chromatin accessibil-
180 ity profiles for Excitatory neurons, Interneurons, and Glial cells (Figure 5D) based
181 on accessibility profiles (Figure 5E), which are then used with the TRANSFAC
182 mouse motif position-weight matrices to construct prior knowledge networks with
183 the inferelator-prior pipeline. Most scRNA-seq cell type clusters have thousands of
184 cells, however a few clusters of rare cell types have as few as 42 (Figure 5G)

185 After processing scRNA-seq into 36 cell type clusters and scATAC data into 3
186 broad (Excitatory neurons, Interneurons, and Glial) priors, we used the Inferelator
187 to learn an aggregate mouse brain GRN. Each of the 36 clusters was assigned the
188 most appropriate of the three prior networks and learned as a separate task using
189 the AMuSR model selection framework. The resulting aggregate network contains
190 20,991 TF - gene regulatory edges, selected from the highest confidence predictions
191 to maximize MCC (Figure 6A-B). 1,909 of these network edges are present in every
192 task-specific network, implying that they are a common regulatory core (Figure
193 6C). Task-specific networks from similar cell types tend to be highly similar, as
194 measured by Jaccard index (Figure 6D). We learn very similar GRNs from each
195 excitatory neuron task, and very similar GRNs from each interneuron task, although
196 each of these broad categories yields different regulatory networks. There are also
197 notable examples where glial and vascular tasks produce GRNs that are distinctively

198 different from other glial and vascular GRNs. Finally, we can examine specific TFs
199 and compare networks between cell type categories (Supplemental Figure 6). The
200 TFs *Egr1* and *Atf4* are expressed in all cell types and *Egr1* is known to have an
201 active role at embryonic day 18 (E18) [48]. In our learned network, *Egr1* targets
202 103 genes, of which 20 are other TFs (Figure 6E-G). Half of these targets (49) are
203 common to both neurons and glial cells, while 38 target genes are specific to neuronal
204 GRNs and 16 target genes are specific to glial GRNs. We identify 14 targets for
205 *Atf4* (Figure 6H), the majority of which (8) are common to both neurons and glial
206 cells, with only 1 target gene specific only to neuronal GRNs and 5 targets specific
207 only to glial GRNs.

208 **3 Discussion**

209 We have developed the Inferelator v3.0 software package to address several key
210 needs in single-cell gene regulatory network inference that have remained difficult
211 to meet with existing solutions. First, this package is well-documented and straight-
212 forward to install and run on an individual computer, in the cloud, or on a large
213 cluster. The Inferelator workflow can be scaled to match the size of the network
214 inference problem and has no organism-specific requirements that preclude easy
215 application to non-standard organisms. Second, different model selection methods
216 can be compared with identical pre- and post-processing methods, including arbi-
217 trary methods implemented through the common scikit-learn estimator framework.
218 Model baselines can be easily established by setting flags to shuffle labels or gen-
219 erate noised data sets, and cross-validation and scoring on holdout genes is built
220 directly into the pipeline. We believe this is particularly important, as many of the
221 performance differences between gene regulatory network inference methods are not
222 due to clever methods for model selection, but are instead the result of differences in
223 data cleaning and preprocessing. Third, we have suggested a principled method for
224 selecting regulatory edges to retain in a GRN. Many GRNs have been inferred by
225 applying a collection of arbitrary heuristics to potential regulatory edges; here we
226 propose ranking regulatory edges by the amount of target gene variance that they
227 explain, and then selecting a threshold for inclusion that maximizes the MCC when
228 scored against a known prior or gold standard network. Finally, we have evaluated
229 the network inference performance on several model organisms that have been well-
230 studied, and for which a reasonable gold standard ground truth GRN can be created
231 from experimental literature. Complex eukaryotes (e.g. mice or humans) lack a gold
232 standard ground truth that can be used to determine real-world network inference
233 performance. Many GRN inference methods instead benchmark on simulated or toy
234 data, with limited experimental validation of a carefully selected tiny subset of an
235 inferred real-world GRNs, making method comparisons difficult.

236 Multi-task modeling is also a key advantage for single-cell GRN inference. Unlike
237 traditional RNA-seq that effectively measures the average gene expression of large
238 number of cells, scRNA-seq can yield individual measurements for many different
239 cell types that are implementing distinct regulatory programs. Learning GRNs from
240 each of these cell types as a separate learning task in a multi-task framework al-
241 lows cell type differences to be retained, while still taking advantage of the common
242 regulatory programs. We demonstrate the use of this multi-task approach to simul-
243 taneously learn regulatory GRNs for a variety of mouse neuronal cell types from a

244 very large (10^6) single-cell data set. This includes learning GRNs for rare cell types;
245 by sharing information between cell types during regression, we are able to learn
246 a core regulatory network while also retaining cell type specific interactions. As an
247 example, the TFs *Egr1* and *Atf4* are broadly expressed and have multiple functions
248 from memory formation and post synaptic development in neurons to cell migration
249 and genome methylation in many cell types [48, 49]. We find a number of target
250 genes regulated by *Egr1* across all neuronal and glial cell types, like the RNA-binding
251 protein *Nova2* that regulates alternative splicing and axonal development [50] (Fig-
252 ure 6F). The stress response TF *Atf4* [51] is known to regulate neuronal GABA_B
253 receptor trafficking [52], and we identify it as a regulator of *Rnf166*, a RING-finger
254 protein that promotes apoptotic cell death in neurons [53]. We also determine that
255 *Atf4* regulates the highly conserved cytosolic sulfotransferase *Sult4a1*, which mod-
256 ulates neuronal branching complexity and dendritic spine formation, and has been
257 linked to neurodevelopmental disorders [54]. As the GRNs that have been learned
258 for each cell type are sparse and consist of the highest-confidence regulatory edges,
259 they are very amenable to exploration and experimental validation.

260 A number of limitations remain that impact our ability to accurately predict
261 gene expression and cell states. Most important is a disconnect between the linear
262 modeling that we use to learn GRNs and the non-linear biophysical models that
263 incorporate both transcription and RNA decay. Modeling strategies that more ac-
264 curately reflect the underlying biology will improve GRN inference directly, and
265 will also allow prediction of useful latent parameters (e.g. RNA half-life) that are
266 experimentally difficult to access. It is also difficult to determine if regulators are
267 activating or repressing specific genes [32], complicated further by biological com-
268 plexity that allows TFs to switch between activation and repression [55]. Improving
269 prediction of the directionality of network edges, and if directionality is stable in
270 different contexts would also be a major advance. Many TFs bind cooperatively as
271 protein complexes, or antagonistically via competitive binding, and explicit model-
272 ing of these TF-TF interactions would also improve GRN inference and make novel
273 biological predictions. Finally, we note that core regulatory networks are likely to
274 be conserved between related species, and further work to develop multi-species
275 inference techniques can leverage evolutionary relationships to improve GRN infer-
276 ence [56]. The modular Inferelator 3.0 framework will allow us to further explore
277 these open problems in regulatory network inference without having to repeatedly
278 reinvent and reimplement existing work.

279 **4 Conclusion**

280 The Inferelator 3.0 is a state-of-the-art, easily deployable, and highly scalable net-
281 work inference tool that is generally applicable to learning GRNs from both single-
282 cell and traditional RNA-seq experiments in any organism of interest. With its
283 accessory software packages, genome-wide expression data of any type can be inte-
284 grated with chromatin accessibility data to construct and explore cell type-specific
285 GRNs. We have established the reliability of this tool by benchmarking on real-
286 world data in model organisms *Bacillus subtilis* and *Saccharomyces cerevisiae* with
287 known gold standard GRNs, and demonstrated how it could be applied to large-
288 scale network inference on many different cell types in the developing mouse brain.

289 We expect this to be a valuable tool to build biologically-relevant GRNs for exper-
290 imental follow-up, as well as a baseline for further development of computational
291 methods in the network inference field.

292 **5 Methods**

293 **5.1 TF Motif-Based Connectivity Matrix (inferelator-prior)**

294 A prior knowledge matrix consists of a signed or unsigned connectivity matrix
295 between regulatory transcription factors (TFs) and target genes. This matrix can
296 be obtained experimentally or by mining regulatory databases. Scanning genomic
297 sequence near promoter regions for TF motifs allows for the construction of motif-
298 derived priors which can be further constrained experimentally by incorporating
299 information about chromatin accessibility [38].

300 We have further refined the generation of prior knowledge matrices with the
301 python inferelator-prior package, which takes as input a gene annotation GTF file,
302 a genomic FASTA file, and a TF motif file, and generates an unsigned connec-
303 tivity matrix. It has dependencies on the common scientific computing packages
304 NumPy [57], SciPy [58], and scikit-learn [59]. In addition, it uses the BEDTools
305 kit [60] and associated python interface pybedtools [61]. The inferelator-prior pack-
306 age (v0.3.0 was used to generate the networks in this manuscript) is available on
307 github (<https://github.com/flatironinstitute/inferelator-prior>) and can
308 be installed through the python package manager pip.

309 *5.1.1 Motif Databases*

310 DNA binding motifs were obtained from published databases. CISBP [62] mo-
311 tifs were obtained from CIS-BP (<http://cisbp.ccbbr.utoronto.ca/>; Build 2.00;
312 Downloaded 11/25/2020) and processed into a MEME-format file with the PWM-
313 toMEME module of inferelator-prior. JASPAR [63] motifs were obtained as MEME
314 files from JASPAR (<http://jaspar.genereg.net/>; 8th Release; Downloaded
315 11/25/2020) . TRANSFAC [64] motifs were licensed from geneXplain (<http://genexplain.com/transfac/>; Version 2020.1; Downloaded 09/13/2020) and pro-
316 cessed into a MEME-format file with the inferelator-prior motif parsing tools. A
317 network of literature-curated network edges was obtained as a gold standard from
318 the YEASTRACT database [65, 66] (<http://www.yeasttract.com/>; Downloaded
319 07/13/2019).

321 *5.1.2 Motif Scanning*

322 Genomic regions of interest are identified by locating annotated Transcription Start
323 Sites (TSS) and opening a window that is appropriate for the organism. For micro-
324 bial species with a compact genome (e.g. yeast), regions of interest are defined as
325 1000bp upstream and 100bp downstream of the TSS. For complex eukaryotes with
326 large intergenic regions (e.g. mammals), regions of interest are defined as 50000bp
327 upstream and 2500bp downstream of the TSS. This is further constrained by inter-
328 secting the genomic regions of interest with a user-provided BED file, which can be
329 derived from a chromatin accessibility experiment (ATAC-seq) or any other method
330 of identifying chromatin of interest. Within these regions of interest, motif locations
331 are identified using the Find Original Motif Occurrences (FIMO) [67] tool from the

MEME suite [68], called in parallel on motif chunks to speed up processing. Each motif hit identified by FIMO is then scored for information content (IC) [69]. IC_i , ranging between 0 and 2 bits, is calculated for each base i in the binding site, where $p_{b,i}$ is the probability of the base b at position i of the motif and $p_{b,bg}$ is the background probability of base b in the genome (Equation 1). Effective information content (EIC) (Equation 2) is the sum of all motif at position i is IC_i penalized with the ℓ_2 -norm of the hit IC_i and the consensus motif base at position i , $IC_{i,consensus}$.

$$IC_i = p_{b,i} \log_2 \left(\frac{p_{b,i}}{p_{b,bg}} \right) \quad (1)$$

$$EIC = \sum_i IC_i - |IC_i - IC_{i,consensus}|_2^2 \quad (2)$$

5.1.3 Connectivity Matrix

A TF-gene binding score is calculated separately for each TF and gene. Each motif hit for a TF within the region of interest around the gene is identified. Overlapping motif hits are resolved by taking the maximum IC for each overlapping base, penalized with the ℓ_2 -norm of differences from the motif consensus sequence. To account for cooperative TF binding effects, any motif hits within 100 bases (25 bases for yeast) are combined, and their EIC scores are summed. The TF-gene binding score is the maximum TF EIC after accounting for overlapping and adjacent TF motifs, and all TF-gene scores are assembled into a Genes x TFs score matrix.

This unfiltered TF-gene score matrix is not sparse as motifs for many TFs are expected to occur often by chance, and TF-gene scores for each TF are not comparable to scores for other TFs as motif position-weight matrices have differing information content. Scores for each TF are clustered using the density-based k-nearest neighbors algorithm DBSCAN [70] (MinPts = 0.001 * number of genes, eps = 1). The cluster of TF-gene edges with the highest score values, and any high-score outliers, are retained in the connectivity matrix, and other TF-gene edges are discarded.

5.1.4 CellOracle Connectivity Matrix

CellOracle [32] was cloned from github (v0.6.5; <https://github.com/morris-lab/CellOracle>; a0da790). CellOracle was provided a BED file with promoter locations for each gene (200bp upstream of transcription start site to 50bp downstream of transcription start site) and the appropriate MEME file for each motif database. Connectivity matrices were predicted using a false positive rate of 0.02 and a motif score threshold of 6. The inferelator-prior pipeline was run using the same promoter locations and MEME files so that the resulting networks are directly comparable, and the Jaccard index between each network and the YEASTRACT network was calculated. Each motif-based network was used as a prior for inferelator network inference on *Saccharomyces cerevisiae*, with the same 2577 genome-wide expression microarray measurements [40]. 20% of the genes were held out of the prior networks and used for scoring the resulting network inference. The motif-based network files have been included in Supplemental Data 1.

369 5.2 Network Inference (The Inferelator)

370 The Inferelator modeling of gene regulatory networks relies on three main modeling
371 assumptions. First, because many transcription factors (TFs) are post transcription-
372 ally controlled and their expression level may not reflect their underlying biological
373 activity, we assume that the activity of a TF can be estimated using expression
374 levels of known targets from prior interactions data [36, 71]. Second, we assume
375 that gene expression can be modeled as a weighted sum of the activities of TFs
376 [34, 45]. Finally, we assume that each gene is regulated by a small subset of TFs
377 and regularize the linear model to enforce sparsity.

378 The Inferelator was initially developed and distributed as an R package [34, 44,
379 72, 73]. We have rewritten it as a python package with dependencies on the common
380 scientific computing packages NumPy [57], SciPy [58], pandas [74], AnnData [29],
381 and scikit-learn [59]. Scaling is implemented either locally through python or as a
382 distributed computation with the Dask [42] parallelization library. The inferelator
383 package (v0.5.4 was used to generate the networks in this manuscript) is available
384 on github (<https://github.com/flatironinstitute/inferelator>) and can be
385 installed through the python package manager pip. The Inferelator takes as in-
386 put gene expression data and prior information on network structure, and outputs
387 ranked regulatory hypotheses of the relative strength and direction of each interac-
388 tion with an associated confidence score.

389 5.3 Transcription Factor Activity

390 The expression level of a TF is often not suitable to describe its activity [75].
391 Transcription factor activity (TFA) is an estimate of the latent activity of a TF that
392 is inducing or repressing transcription of its targets in a sample. A gene expression
393 dataset (\mathbf{X}) is a matrix where $X_{i,j}$ is the observed mRNA expression level ($i \in$
394 Samples and $j \in$ Genes), measured either by microarray, RNA-seq, or single cell
395 RNA sequencing (scRNA-seq).

$$X_{i,j} = \sum_k A_{i,k} P_{k,j} \quad (3)$$

396 We estimate TFA by solving (Equation 3) for activity ($A_{i,k}$), where $k \in$ TFs, and
397 P is a prior connectivity matrix where $P_{k,j}$ is 1 if gene j is regulated by TF k
398 and 0 if it is not. In matrix notation, $\mathbf{X} = \mathbf{A}\mathbf{P}$, and $\hat{\mathbf{A}}$ is estimated by minimizing
399 $\|\hat{\mathbf{A}}\mathbf{P} - \mathbf{X}\|_2^2$. This is calculated by taking the pseudoinverse of \mathbf{P} and solving
400 $\hat{\mathbf{A}} = \mathbf{X}\mathbf{P}^{-1}$. The resulting $\hat{\mathbf{A}}$ is a matrix where $\hat{A}_{i,k}$ is the estimated latent TFA
401 for sample i and TF k . In cases where all values in \mathbf{P} for a TF are 0, that TF is
402 removed from \mathbf{P} and the expression \mathbf{X} of that TF is used in place of activity.

403 5.4 Inferelator Network Inference

404 Linear models (Equation 4) are separately constructed for each gene j .

$$X_i = \sum_k \hat{A}_{i,k} \beta_k \quad (4)$$

405 In addition to the model selection methods described here, we have implemented a
406 module which takes any scikit-learn regression object (for example, elastic net [76]).
407 Model selection and regularization techniques are applied to enforce the biological
408 property of sparsity. If the coefficient $\beta_{j,k}$ is non-zero, it is evidence for a regulatory
409 relationship between TF k and gene j .

$$S_{j,k} = 1 - \frac{\sigma_{allTFs}^2}{\sigma_{TF_k\text{leaveout}}^2} \quad (5)$$

410 For each gene j , the amount of variance explained by each regulatory TF k is
411 calculated as the ratio between the variance of the residuals in the full model and
412 the variance of the residuals when the linear model is refit by ordinary least squares
413 (OLS) and k is left out (Equation 5).

414 In order to mitigate the effect of outliers and sampling error, model selection is
415 repeated multiple times using input expression data \mathbf{X} that has been bootstrapped
416 (resampled with replacement). Predicted TF-gene interactions are ranked for each
417 bootstrap by amount of variance explained and then rank-combined into a unified
418 network prediction. Confidence scores are assigned based on the combined rank
419 for each interaction, and the overall network is compared to a gold standard and
420 performance is evaluated by area under the precision-recall curve.

421 The effects of setting hyperparameters can be tested by cross-validation on the
422 prior and gold standard networks. This strategy holds out a subset of genes (rows)
423 from the prior knowledge network \mathbf{P} . Network inference performance is then eval-
424 uated on only those held-out genes, using the gold standard network.

425 5.4.1 Model Selection: Bayesian Best Subset Regression

426 Bayesian Best Subset Regression (BBSR) is a model selection method described in
427 detail in [73]. Initial feature selection for this method is necessary as best subset re-
428 gression on all possible combinations of hundreds of TF features is computationally
429 intractable. We therefore select ten TF features with the highest context likelihood
430 of relatedness between expression of each gene and activity of each TF. This method
431 is described in detail in [72].

432 First, gene expression and TF activity are discretized into equal-width bins ($n=10$)
433 and mutual information is calculated based on their discrete probability distribu-
434 tions (Equation 6) to create a mutual information matrix \mathbf{M}^{dyn} .

$$M_{j,k}^{dyn} = p(X_j, \hat{A}_k) \log \frac{p(X_j, \hat{A}_k)}{p(X_j)p(\hat{A}_k)} \quad (6)$$

$$M_{k_1,k_2}^{stat} = p(\hat{A}_{k_1}, \hat{A}_{k_2}) \log \frac{p(\hat{A}_{k_1}, \hat{A}_{k_2})}{p(\hat{A}_{k_1})p(\hat{A}_{k_2})} \quad (7)$$

435 Mutual information is also calculated between activity of each TF (Equation 7) to
436 create a mutual information matrix \mathbf{M}^{stat} .

$$z_{j,k}^{dyn} = \frac{M_{j,k}^{dyn} - \sum_j \frac{M_{j,k}^{dyn}}{n_i}}{\sigma_k^{dyn}} \quad (8)$$

$$z_{j,k}^{stat} = \frac{M_{j,k}^{stat} - \sum_j \frac{M_{j,k}^{stat}}{n_i}}{\sigma_k^{stat}} \quad (9)$$

$$z_{j,k}^{mixed} = \sqrt{(z_{j,k}^{dyn})^2 + (z_{j,k}^{stat})^2} \quad (10)$$

437 A mixed context likelihood of relatedness score is then calculated as a pseudo-zscore
 438 by calculating \mathbf{Z}^{dyn} (Equation 8) and \mathbf{Z}^{stat} (Equation 9). Any values less than 0
 439 in \mathbf{Z}^{dyn} or \mathbf{Z}^{stat} are set to 0, and then they are combined into a mixed context
 440 likelihood of relatedness matrix \mathbf{Z}^{mixed} (Equation 10). For each gene j , the 10
 441 TFs with the highest mixed context likelihood of relatedness values are selected for
 442 regression.

443 For best subset regression, a linear model is fit with OLS for every combination
 444 of the selected predictor variables.

$$\rho(\beta, \sigma^2 | X_j) = \rho(\beta | X_j, \sigma^2) \rho(\sigma^2 | X_i) \quad (11)$$

$$\rho(\sigma^2 | X_i) \propto IG\left(\frac{n}{2}, \frac{SSR}{2} + \frac{(\beta_0 - \beta_{OLS}) \mathbf{G} \mathbf{X}' \mathbf{X} \mathbf{G} (\beta_0 - \beta_{OLS})}{2}\right) \quad (12)$$

445 We define β_0 as our null prior for the model parameters (zeros), β_{OLS} as the model
 446 coefficients from OLS, SSR as the sum of squared residuals, and \mathbf{G} as a g -prior
 447 diagonal matrix where the diagonal values represent a weight for each predictor
 448 variable. g -prior weights in \mathbf{G} close to 0 favor β values close to β_0 . Large g -prior
 449 weights favor β values close to β_{OLS} . By default, we select g -prior weights of 1 for
 450 all predictor variables. From the joint posterior distribution (Equation 11) we can
 451 calculate the marginal posterior distribution of σ^2 (Equation 12), where IG is the
 452 inverse gamma distribution. The Bayesian information criterion (BIC) is calculated
 453 for each model, where n is the number of observations and k is the number of
 454 predictors (Equation 13).

$$BIC = n \ln(\sigma^2) - k \ln(n) \quad (13)$$

$$E[\sigma^2] = \frac{\frac{SSR}{2} + \frac{(\beta_0 - \beta_{OLS}) \mathbf{G} \mathbf{X}' \mathbf{X} \mathbf{G} (\beta_0 - \beta_{OLS})}{2}}{\frac{n}{2} - 1} \quad (14)$$

$$E[BIC] = n \left(\ln\left(\frac{SSR}{2} + \frac{(\beta_0 - \beta_{OLS}) \mathbf{G} \mathbf{X}' \mathbf{X} \mathbf{G} (\beta_0 - \beta_{OLS})}{2}\right) - \text{Digamma}\left(\frac{n}{2}\right) \right) - k \ln(n) \quad (15)$$

455 We calculate the expected posterior distribution of σ^2 (Equation 14) for each subset
456 of predictors, and use it to determine the model BIC (Equation 15). We then select
457 the model with the smallest $E[BIC]$. The predictors in the selected subset model
458 for gene j are TFs which regulate its expression.

459 5.4.2 Model Selection: StARS-LASSO

460 Least absolute shrinkage and selection operator (LASSO) [77] combined with the
461 Stability Approach to Regularization Selection (StARS) [78] is a model selection
462 method described in detail in [38]. In short, the StARS-LASSO approach is to
463 select the optimal λ parameter for (Equation 16). N random subsamples of X and
464 \hat{A} without replacement subnetworks $S_{n,\lambda}$ are defined as the non-zero coefficients
465 $\beta_{n,\lambda}$ after LASSO regression. Initially, λ is set large, so that each subnetwork S_n
466 is highly sparse, and is then decreased, resulting in increasingly dense networks.
467 Edge instability is calculated as the fraction of times subnetworks disagree about
468 the presence of an network edge. As λ decreases, the subnetworks are expected
469 to have increasing edge instability initially and then decreasing edge instability as
470 λ approaches 0, as (Equation 16) reduces to OLS and each subnetwork becomes
471 dense.

$$\min_{\beta} \frac{1}{2n} \|X - \hat{A}\beta\|_2^2 - \lambda \|\beta\|_1 \quad (16)$$

472 We choose the largest value of λ such that the edge instability is less than 0.05, which
473 is interpretable as all subnetworks share $> 95\%$ of edges. This selection represents
474 a balance between increasing the network size and minimizing the instability that
475 occurs when data is sampled.

476 5.5 Multiple Task Network Inference

477 We separate biological samples which represent different states into separate tasks,
478 learn networks from these tasks, and then combine task-specific networks into an
479 ensemble network. One method of solving these states is to sequentially apply a
480 single-task method for network inference (i.e. 5.4.1 or 5.4.2). The networks generated
481 for each task are then rank-combined into a unified network. The Adaptive Multiple
482 Sparse Regression (AMuSR) method, described in detail in [45], uses a multi-task
483 learning framework, where each task is solved together.

$$\arg \min_{B,S} \frac{1}{2n} \|X_{d,i} - (S_d + B)\hat{A}_d\|_2^2 + \lambda_s \|S\|_{1,1} + \lambda_b \|B\|_{1,\infty} \quad (17)$$

$$\hat{W} = \hat{B} + \hat{S} \quad (18)$$

484 In (Equation 17), \mathbf{B} is a block-sparse weight matrix in which the weights for any
485 feature are the same across all tasks. \mathbf{S} is a sparse weight matrix in which the weights
486 for features can vary between tasks. The combination \mathbf{W} of \mathbf{B} and \mathbf{S} (Equation 18)
487 are model weights representing regulatory interactions between TFs and genes. In

488 short, this method uses adaptive penalties to favor regulatory interactions shared
489 across multiple tasks in \mathbf{B} , while recognizing dataset specific interactions in \mathbf{S} .
490 Model hyperparameters λ_s and λ_b are identified by grid search, selecting the model
491 that minimizes the extended Bayesian Information Criterion (eBIC) (Equation 19),
492 where D is the number of task datasets, and for dataset d , n_d is the number of
493 observations, $X_i^{(d)}$ is gene expression for gene i , $\hat{A}^{(d)}$ is TF activity estimates, $W_{*,d}$
494 is model weights, k_d is the number of non-zero predictors, and p_d is the total number
495 of predictors. For this work, we choose to set the eBIC parameter γ to 1.

$$eBIC = \frac{1}{D} \sum n_d \ln \frac{1}{n_d} \|X_i^{(d)} - \hat{A}^{(d)T} W_{*,d}\|_2^2 + k_d \ln n_d + 2\gamma \ln \binom{p_d}{k_d} \quad (19)$$

496 5.6 Network Performance Metrics

497 Prior work has used the area under the Precision (Equation 20) - Recall (Equation
498 21) curve to determine performance, by comparing to some known, gold-standard
499 network. Here we add two metrics; Matthews correlation coefficient [79] (MCC)
500 (Equation 22) and F1 score (Equation 23). MCC can be calculated directly from
501 the confusion matrix True Positive (TP), False Positive (FP), True Negative (TN),
502 and False Negative (FN) values.

$$Precision = \frac{TP}{TP + FP} \quad (20)$$

$$Recall = \frac{TP}{TP + FN} \quad (21)$$

$$MCC = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (22)$$

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (23)$$

503 We compute an MCC and F1 score for each cutoff along ranked interactions in
504 order to generate MCC and F1 scores for all possible networks in growing ranked
505 order. The maximum MCC along ranked interactions gives the subnetwork that
506 has maximum similarity to the comparison network, accounting for TP, FP, TN,
507 and FN. The maximum F1 along ranked interactions gives the subnetwork that has
508 maximum similarity to the comparison network accounting for TP, FP, and FN.

509 5.7 Network Inference in *Bacillus subtilis*

510 Microarray expression data for *Bacillus subtilis* was obtained from NCBI GEO;
511 GSE67023 [36] (n=268) and GSE27219 [43] (n=266). The Inferelator (v0.5.4)
512 learned GRNs using each expression dataset separately in conjunction with a known
513 prior network [36] (Supplemental Data 1). Performance was evaluated by AUPR on
514 ten replicates by holding 20% of the genes in the known prior network out, learning
515 the GRN, and then scoring based on the held-out genes. Baseline shuffled controls
516 were performed by randomly shuffling the labels on the known prior network.

517 Multi-task network inference uses the same *B. subtilis* prior for both tasks, with
518 20% of genes held out for scoring. Individual task networks are learned and rank-
519 combined into an aggregate network. Performance was evaluated by AUPR on the
520 held-out genes.

521 5.8 Network Inference in *Saccharomyces cerevisiae*

522 A large microarray dataset was obtained from NCBI GEO and normalized for a pre-
523 vious publication [40] (n=2,577). It is available on zenodo with DOI: 10.5281/zen-
524 odo.3247754. In short, this data was preprocessed with limma [80] and quantile
525 normalized. A second microarray dataset consisting of a large dynamic perturba-
526 tion screen [16] was obtained from NCBI GEO accession GSE142864 (n=1,693).
527 This dataset is \log_2 fold change of an experimental channel over a control channel
528 which is the same for all observations. The Inferelator (v0.5.4) learned GRNs using
529 each expression dataset separately in conjunction with a known YEASTRACT prior
530 network [65, 66] (Supplemental Data 1). Performance was evaluated by AUPR on
531 ten replicates by holding 20% of the genes in the known prior network out, learning
532 the GRN, and then scoring based on the held-out genes in a separate gold standard
533 [40]. Baseline shuffled controls were performed by randomly shuffling the labels on
534 the known prior network.

535 Multi-task network inference uses the same YEASTRACT prior for both tasks,
536 with 20% of genes held out for scoring. Individual task networks are learned and
537 rank-combined into an aggregate network, which is then evaluated by AUPR on the
538 held-out genes in the separate gold standard.

539 5.9 Single-Cell Network Inference in *Saccharomyces cerevisiae*

540 Single-cell expression data for *Saccharomyces cerevisiae* was obtained from NCBI
541 GEO (GSE125162 [22] and GSE144820 [47]). Individual cells (n=44,343) are orga-
542 nized into one of 14 groups based on experimental metadata and used as separate
543 tasks in network inference. Genes were filtered such that any gene with fewer than
544 than 2217 total counts in all cells (1 count per 20 cells) was removed. Data was used
545 as raw, unmodified counts, was Freeman-Tukey transformed ($\sqrt{x+1} + \sqrt{x-1}$), or
546 was \log_2 pseudocount transformed ($\log_2(x+1)$). Data was either not normalized,
547 or depth normalized by scaling so that the sum of all counts for each cell is equal
548 to the median of the sum of counts of all cells. For each set of parameters, network
549 inference is run 10 times, using the YEASTRACT network as prior knowledge with
550 20% of genes held out for scoring. For noise-only controls, gene expression counts
551 are simulated randomly such that for each gene i , $x_i \sim N(\mu_{x_i}, \sigma_{x_i})$ and the sum for
552 each cell is equal to the sum in the observed data. For shuffled controls, the gene
553 labels on the prior knowledge network are randomly shuffled.

554 5.10 Single-Cell Network Inference in *Mus musculus* neurons

555 Single-cell expression data from *Mus musculus* brain samples taken at E18 was ob-
556 tained from 10x genomics [81]. SCANPY was used to preprocess and cluster the
557 scRNAseq dataset. Genes present in fewer than 2% of cells were removed. Cells
558 were filtered out when fewer than 1000 genes were detected, the cell had more than
559 20,000 total gene counts, or the cell had more than 7% of gene counts assigned to

560 mitochondrial transcripts. Transcript counts were then log transformed and nor-
561 malized and scaled. Cells were assigned to mitotic or post mitotic phase based on
562 cell cycle marker genes using `score_genes_cell_cycle` [82]. In order to focus on neu-
563 ronal cells, all 374,369 mitotic cells were removed. Remaining cells were clustered
564 by Leiden clustering (Resolution = 0.5) using the first 300 principal components of
565 the 2000 most highly variable genes. Broad cell types were assigned to each cluster
566 based on the expression of marker genes *Neurod6* for Excitatory neurons, *Gad1* for
567 Interneurons, and *Apoe* for glial cells. Cells from each broad cell type were then
568 re-clustered into clusters based on the 2000 most highly variable genes within the
569 cluster. Specific cell types were assigned to each subcluster based on the expression
570 of marker genes[83]. Ambiguous clusters were discarded, removing 151,765 cells,
571 leaving resulting in 36 specific cell type clusters that consist of 766,402 total cells.

572 Single-cell ATAC data from *Mus musculus* brain samples taken at E18 was ob-
573 tained from 10x genomics; datasets are from samples prepared fresh [84], samples
574 dissociated and cryopreserved [85], and samples flash-frozen [86]. ChromA [87] and
575 SnapATAC [88] were used to process the scATACseq datasets. Consensus peaks
576 were called on the 3 datasets using ChromA. Each dataset was then run through
577 the SnapATAC pipeline using the consensus peaks. Cells were clustered and labels
578 from the scRNAseq object were transferred to the scATAC data. Cells that did not
579 have an assignment score $\geq .5$ were discarded. Assigned barcodes were split by cell
580 class(EXC, IN or GL). ChromA was run again for each cell class generating 3 sets
581 of cell class specific peaks.

582 Aggregated chromatin accessibility profiles were used with TRANSFAC v2020.1
583 motifs and the inferelator-prior (v0.3.0) pipeline to create prior knowledge connec-
584 tivity matrices between TFs and target genes for excitatory neurons, interneurons,
585 and glial cells. Vascular cells were not present in the scATAC data sufficiently to
586 allow construction of a vascular cell prior with this method, and so vascular cells
587 were assigned the glial prior for network inference.

588 GRNs were learned using AMuSR on \log_2 pseudocount transformed count data
589 for each of 36 cell type specific clusters as separate tasks with the appropriate prior
590 knowledge network. An aggregate network was created by rank-summing each cell
591 type GRN. MCC was calculated for this aggregate network based on a comparison
592 to the union of the three prior knowledge networks, and the confidence score which
593 maximized MCC was selected as a threshold to determine the size of the final
594 network. Neuron specific edges were computed by aggregating filtered individual
595 task networks with their respective confidence score to maximize MCC. Each edge
596 that was shared with a glial or vascular network was removed. The remaining neuron
597 specific edges are interneuron specific, excitatory specific or shared.

598 5.11 Inferelator 3.0 Single-Cell Computational Speed Profiling

599 144,682 mouse cells from the mouse neuronal subcluster EXC_IT.1 were used with
600 the mouse excitatory neuron prior knowledge network to determine Inferelator 3.0
601 runtime. To benchmark the Dask engine, the Inferelator was deployed to 5 28-core
602 (Intel® Xeon® E5-2690) nodes for a total of 140 cpu cores. To benchmark the
603 python-based multiprocessing engine, the Inferelator was deployed to a single 28-
604 core (Intel® Xeon® E5-2690) node. Either all 144,682 mouse cells were used, or a

605 subset was randomly selected for each run, and used to learn a single GRN. Runtime
606 was determined by the length of workflow execution, which includes loading data,
607 running all regressions, and producing output files.

608 5.12 Visualization

609 Figures were generated with R [89] and the common ggplot2 [90], umap [91], and
610 tidyverse packages [92]. Additional figures were generated with python using scanpy
611 [29], matplotlib [93], and seaborn [94]. Network diagrams were created with the
612 python package `jp_gene_viz` [95]. Schematic figures were created in Adobe Illustrator,
613 and other figures were adjusted in Illustrator to improve panelling and layout.

6 Declarations

Ethics Approval and Consent to Participate

Not applicable

Competing interests

The authors declare that they have no competing interests.

Availability of Data and Materials

The datasets supporting the conclusions of this article are available in the NCBI GEO repository with accession IDs: GSE125162, GSE144820, GSE67023, GSE27219, GSE142864. A large number of GEO records were compiled and normalized in a previous work [40] into a combined dataset which is available on Zenodo (DOI: 10.5281/zenodo.3247754). Single-cell mouse datasets are publicly available from 10x genomics [81, 84, 85, 86] under a Creative Commons Attribution (CC-BY 4.0) license. Software packages developed for this article are available on github (<https://github.com/flatironinstitute/inferelator> and <https://github.com/flatironinstitute/inferelator-prior>) and have been released as python packages through PyPi (<https://pypi.org/project/inferelator/> and <https://pypi.org/project/inferelator-prior/>). Specific analysis scripts for this work have been included in Supplemental Data 1.

Author's contributions

CSG contributed to Methodology, Software, Validation, Formal Analysis, Writing – Original Draft Preparation, and Visualization. CJ and GS contributed to Conceptualization, Methodology, Software, Validation, Investigation, Resources, Data Curation, Formal Analysis, Writing – Original Draft Preparation, and Visualization. AS contributed to Validation, Data Curation, Formal Analysis, and Visualization. AW contributed to Software and Visualization. AT contributed to Software, Writing – Original Draft Preparation, and Formal Analysis. DC and KT contributed to Software, Data Curation, and Conceptualization. NDV, NC, RY, and TH contributed to Software. DG contributed to Supervision, Project Administration, and Funding Acquisition. EM contributed to Conceptualization, Writing – Original Draft Preparation, and Software. RB contributed to Conceptualization, Writing – Original Draft Preparation, Supervision, Project Administration, and Funding Acquisition.

Funding

RB thanks the following sources for research support: NSF IOS-1546218, NIH R35GM122515, NIH R01HD096770, NIH R01NS116350, and the Simons Foundation. DG thanks the following sources for research support: NSF MCB-1818234, NIH R01GM107466, NIH R01GM134066, and NIH R01AI140766.

Acknowledgements

We thank past and present members of the Gresham, Miraldi, and Bonneau labs for discussions and valuable feedback on this manuscript. We also thank the staff of the Flatiron Institute Scientific Computing Core for their tireless efforts to build and maintain the High Performance Computing resources which we rely on. This work was supported in part through the NYU IT High Performance Computing resources, services, and staff expertise.

Additional Files

Supplemental Data 1 is a .tar.gz file containing the prior knowledge networks used in this work, the gold standard networks used in this work, and the python scripts used to generate the learned networks in this work
Supplemental Data 2 is a .tar.gz file containing the mouse E18 neuronal network learned in Figure 6 of this work
Supplemental Table 1 is a .tsv file containing the crossvalidation performance results from Figure 2 of this work
Supplemental Table 2 is a .tsv file containing the crossvalidation performance results from Figure 3 of this work
Supplemental Table 3 is a .tsv file containing the crossvalidation performance results from Figure 4 of this work

Author details

¹ Flatiron Institute, Center for Computational Biology, Simons Foundation, 10010, New York, USA. ² Center For Genomics and Systems Biology, NYU, 10008, New York, USA. ³ Department of Biology, NYU, 10008, New York, USA. ⁴ Courant Institute of Mathematical Sciences, Computer Science Department, NYU, 10008, New York, USA. ⁵ Center For Data Science, NYU, 10008, New York, USA. ⁶ Department of Systems Biology, Columbia University, 10032, New York USA. ⁷ Flatiron Institute, Scientific Computing Core, Simons Foundation, 10010, New York, USA. ⁸ Divisions of Immunobiology and Biomedical Informatics, Cincinnati Children's Hospital Medical Center, 45229, Cincinnati USA. ⁹ Department of Pediatrics, University of Cincinnati College of Medicine, 45229, Cincinnati USA.

References

1. Zaret, K.S.: Pioneer transcription factors initiating gene network changes. *Annu. Rev. Genet.* (2020)
2. Kadonaga, J.T.: Regulation of RNA polymerase II transcription by sequence-specific DNA binding factors. *Cell* **116**(2), 247–257 (2004)
3. Hahn, S., Young, E.T.: Transcriptional regulation in *Saccharomyces cerevisiae*: transcription factor regulation and function, mechanisms of initiation, and roles of activators and coactivators. *Genetics* **189**(3), 705–736 (2011)
4. Lambert, S.A., Jolma, A., Campitelli, L.F., Das, P.K., Yin, Y., Albu, M., Chen, X., Taipale, J., Hughes, T.R., Weirauch, M.T.: The human transcription factors. *Cell* **172**(4), 650–665 (2018)
5. Thompson, D., Regev, A., Roy, S.: Comparative analysis of gene regulatory networks: from network reconstruction to evolution. *Annu. Rev. Cell Dev. Biol.* **31**, 399–428 (2015)
6. Chasman, D., Fotuhi Siahipirani, A., Roy, S.: Network-based approaches for analysis of complex biological systems. *Curr. Opin. Biotechnol.* **39**, 157–166 (2016)
7. Hu, J.X., Thomas, C.E., Brunak, S.: Network biology concepts in complex disease comorbidities. *Nat. Rev. Genet.* **17**(10), 615–629 (2016)
8. Mehta, T.K., Koch, C., Nash, W., Knaack, S.A., Sudhakar, P., Olbei, M., Bastkowski, S., Penso-Dolfi, L., Korcsmaros, T., Haerty, W., Roy, S., Di-Palma, F.: Evolution of regulatory networks associated with traits under selection in cichlids. *Genome Biol.* **22**(1), 25 (2021)
9. Peter, I.S., Davidson, E.H.: Evolution of gene regulatory networks controlling body plan development. *Cell* **144**(6), 970–985 (2011)
10. Huang, M., Bao, J., Hallström, B.M., Petranovic, D., Nielsen, J.: Efficient protein production by yeast requires global tuning of metabolism. *Nat. Commun.* **8**(1), 1131 (2017)
11. DeRisi, J.L., Iyer, V.R., Brown, P.O.: Exploring the metabolic and genetic control of gene expression on a genomic scale. *Science* **278**(5338), 680–686 (1997)
12. Nagalakshmi, U., Wang, Z., Waern, K., Shou, C., Raha, D., Gerstein, M., Snyder, M.: The transcriptional landscape of the yeast genome defined by RNA sequencing. *Science* **320**(5881), 1344–1349 (2008)
13. Wang, Z., Gerstein, M., Snyder, M.: RNA-Seq: a revolutionary tool for transcriptomics. *Nat. Rev. Genet.* **10**(1), 57–63 (2009)
14. Szederkényi, G., Banga, J.R., Alonso, A.A.: Inference of complex biological networks: distinguishability issues and optimization-based solutions. *BMC Syst. Biol.* **5**, 177 (2011)
15. Ud-Dean, S.M.M., Gunawan, R.: Optimal design of gene knockout experiments for gene regulatory network inference. *Bioinformatics* **32**(6), 875–883 (2016)
16. Hackett, S.R., Baltz, E.A., Coram, M., Wrani, B.J., Kim, G., Baker, A., Fan, M., Hendrickson, D.G., Berndt, M., Mclsaac, R.S.: Learning causal networks using inducible transcription factors and transcriptome-wide time series. *Mol. Syst. Biol.* **16**(3), 9174 (2020)
17. Macosko, E.Z., Basu, A., Satija, R., Nemes, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A.R., Kamitaki, N., Martnersteck, E.M., Trombetta, J.J., Weitz, D.A., Sanes, J.R., Shalek, A.K., Regev, A., McCarroll, S.A.: Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell* **161**(5), 1202–1214 (2015)
18. Zilionis, R., Nainys, J., Veres, A., Savova, V., Zemmour, D., Klein, A.M., Mazutis, L.: Single-cell barcoding and sequencing using droplet microfluidics. *Nat. Protoc.* **12**(1), 44 (2017)
19. Zheng, G.X.Y., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B., Wheeler, T.D., McDermott, G.P., Zhu, J., Gregory, M.T., Shuga, J., Montesclaros, L., Underwood, J.G., Masquelier, D.A., Nishimura, S.Y., Schnall-Levin, M., Wyatt, P.W., Hindson, C.M., Bharadwaj, R., Wong, A., Ness, K.D., Beppu, L.W., Deeg, H.J., McFarland, C., Loeb, K.R., Valente, W.J., Ericson, N.G., Stevens, E.A., Radich, J.P., Mikkelsen, T.S., Hindson, B.J., Bielas, J.H.: Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**, 14049 (2017)
20. Rosenberg, A.B., Roco, C.M., Muscat, R.A., Kuchina, A., Sample, P., Yao, Z., Graybuck, L.T., Peeler, D.J., Mukherjee, S., Chen, W., Pun, S.H., Sellers, D.L., Tasic, B., Seelig, G.: Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding. *Science* **360**(6385), 176–182 (2018)
21. Dixit, A., Parnas, O., Li, B., Chen, J., Fulco, C.P., Jerby-Arnon, L., Marjanovic, N.D., Dionne, D., Burks, T., Raychowdhury, R., Adamson, B., Norman, T.M., Lander, E.S., Weissman, J.S., Friedman, N., Regev, A.: Perturb-Seq: Dissecting molecular circuits with scalable Single-Cell RNA profiling of pooled genetic screens. *Cell* **167**(7), 1853–1866 (2016)
22. Jackson, C.A., Castro, D.M., Saldi, G.-A., Bonneau, R., Gresham, D.: Gene regulatory network reconstruction using single-cell RNA sequencing of barcoded genotypes in diverse environments. *Elife* **9**, 51254 (2020)
23. Schraivogel, D., Gschwind, A.R., Milbank, J.H., Leonce, D.R., Jakob, P., Mathur, L., Korb, J.O., Merten, C.A., Velten, L., Steinmetz, L.M.: Targeted perturb-seq enables genome-scale genetic screens in single cells. *Nat. Methods* **17**(6), 629–635 (2020)
24. Ursu, O., Neal, J.T., Shea, E., Thakore, P.I., Jerby-Arnon, L., Nguyen, L., Dionne, D., Diaz, C., Bauman, J., Mosaad, M., Fagre, C., Giacomelli, A., Ly, S.H., Rozenblatt-Rosen, O., Hahn, W., Aguirre, A., Berger, A., Regev, A., Boehm, J.S.: Massively parallel phenotyping of variant impact in cancer with Perturb-seq reveals a shift in the spectrum of cell states induced by somatic mutations (2020)
25. Svensson, V.: Droplet scRNA-seq is not zero-inflated. *Nat. Biotechnol.* (2020)
26. Aridakesian, C., Poirion, O., Yunits, B., Zhu, X., Garmire, L.X.: DeepImpute: an accurate, fast, and scalable deep neural network method to impute single-cell RNA-seq data. *Genome Biol.* **20**(1), 211 (2019)
27. Tjårnberg, A., Mahmood, O., Jackson, C.A., Saldi, G.-A., Cho, K., Christiaen, L.A., Bonneau, R.A.: Optimal tuning of weighted kNN- and diffusion-based methods for denoising single cell genomics data. *PLoS Comput. Biol.* **17**(1), 1008569 (2021)
28. Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck, W.M. 3rd, Hao, Y., Stoeckius, M., Smibert, P., Satija, R.: Comprehensive integration of Single-Cell data. *Cell* **177**(7), 1888–1902 (2019)

29. Wolf, F.A., Angerer, P., Theis, F.J.: SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* **19**(1), 15 (2018)
30. Van de Sande, B., Flerin, C., Davie, K., De Waegeneer, M., Hulselmans, G., Aibar, S., Seurinck, R., Saelens, W., Cannoodt, R., Rouchon, Q., Verbeiren, T., De Maeyer, D., Reumers, J., Saeyns, Y., Aerts, S.: A scalable SCENIC workflow for single-cell gene regulatory network analysis. *Nat. Protoc.* (2020)
31. Huynh-Thu, V.A., Irrthum, A., Wehenkel, L., Geurts, P.: Inferring regulatory networks from expression data using tree-based methods. *PLoS One* **5**(9) (2010)
32. Kamimoto, K., Hoffmann, C.M., Morris, S.A.: CellOracle: Dissecting cell identity via network inference and in silico gene perturbation (2020)
33. Matsumoto, H., Kiryu, H., Furusawa, C., Ko, M.S.H., Ko, S.B.H., Gouda, N., Hayashi, T., Nikaido, I.: SCODE: an efficient regulatory network inference algorithm from single-cell RNA-Seq during differentiation. *Bioinformatics* **33**(15), 2314–2321 (2017)
34. Bonneau, R., Reiss, D.J., Shannon, P., Facciotti, M., Hood, L., Baliga, N.S., Thorsson, V.: The inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets de novo. *Genome Biol.* **7**, 36 (2006)
35. Madar, A., Greenfield, A., Ostrer, H., Vanden-Eijnden, E., Bonneau, R.: The inferelator 2.0: A scalable framework for reconstruction of dynamic regulatory network models. In: 2009 Annual International Conference of the IEEE Engineering in Medicine and Biology Society, pp. 5448–5451 (2009)
36. Arrieta-Ortiz, M.L., Hafemeister, C., Bate, A.R., Chu, T., Greenfield, A., Shuster, B., Barry, S.N., Gallitto, M., Liu, B., Kacmarczyk, T., Santoriello, F., Chen, J., Rodrigues, C.D.A., Sato, T., Rudner, D.Z., Driks, A., Bonneau, R., Eichenberger, P.: An experimentally supported model of the bacillus subtilis global transcriptional regulatory network. *Mol. Syst. Biol.* **11**(11), 839 (2015)
37. Ciofani, M., Madar, A., Galan, C., Sellars, M., Mace, K., Pauli, F., Agarwal, A., Huang, W., Parkurst, C.N., Muratet, M., Newberry, K.M., Meadows, S., Greenfield, A., Yang, Y., Jain, P., Kirigin, F.K., Birchmeier, C., Wagner, E.F., Murphy, K.M., Myers, R.M., Bonneau, R., Littman, D.R.: A validated regulatory network for th17 cell specification. *Cell* **151**(2), 289–303 (2012)
38. Miraldi, E.R., Pokrovskii, M., Watters, A., Castro, D.M., De Veaux, N., Hall, J.A., Lee, J.-Y., Ciofani, M., Madar, A., Carriero, N., Littman, D.R., Bonneau, R.: Leveraging chromatin accessibility for transcriptional regulatory network inference in T helper 17 cells. *Genome Res.* (2019)
39. Pokrovskii, M., Hall, J.A., Ochayon, D.E., Yi, R., Chaimowitz, N.S., Seelamneni, H., Carriero, N., Watters, A., Waggoner, S.N., Littman, D.R., Bonneau, R., Miraldi, E.R.: Characterization of transcriptional regulatory networks that promote and restrict identities and functions of intestinal innate lymphoid cells. *Immunity* **51**(1), 185–1976 (2019)
40. Tchourine, K., Vogel, C., Bonneau, R.: Condition-Specific modeling of biophysical parameters advances inference of regulatory networks. *Cell Rep.* **23**(2), 376–388 (2018)
41. Wilkins, O., Hafemeister, C., Plessis, A., Holloway-Phillips, M.-M., Pham, G.M., Nicotra, A.B., Gregorio, G.B., Jagadish, S.V.K., Septiningsih, E.M., Bonneau, R., Purugganan, M.: EGRINs (environmental gene regulatory influence networks) in rice that function in the response to water deficit, high temperature, and agricultural environments. *Plant Cell* **28**(10), 2365–2384 (2016)
42. Rocklin, M.: Dask: Parallel computation with blocked algorithms and task scheduling. In: Proceedings of the 14th Python in Science Conference. Proceedings of the Python in Science Conference, pp. 126–132. SciPy, ??? (2015)
43. Nicolas, P., Mäder, U., Dervyn, E., Rochat, T., Leduc, A., Pigeonneau, N., Bidnenko, E., Marchadier, E., Hoebeke, M., Aymerich, S., Becher, D., Biscicchia, P., Botella, E., Delumeau, O., Doherty, G., Denham, E.L., Fogg, M.J., Fromion, V., Goelzer, A., Hansen, A., Härtig, E., Harwood, C.R., Homuth, G., Jarmer, H., Jules, M., Klipp, E., Le Chat, L., Lecointe, F., Lewis, P., Liebermeister, W., March, A., Mars, R.A.T., Nannapaneni, P., Noone, D., Pohl, S., Rinn, B., Rügheimer, F., Sappa, P.K., Samson, F., Schaffer, M., Schwikowski, B., Steil, L., Stülke, J., Wiegert, T., Devine, K.M., Wilkinson, A.J., van Dijl, J.M., Hecker, M., Völker, U., Bessières, P., Noirot, P.: Condition-dependent transcriptome reveals high-level regulatory architecture in bacillus subtilis. *Science* **335**(6072), 1103–1106 (2012)
44. Greenfield, A., Madar, A., Ostrer, H., Bonneau, R.: DREAM4: Combining genetic and dynamic information to identify biological networks and dynamical models. *PLoS One* **5**(10), 13397 (2010)
45. Castro, D.M., de Veaux, N.R., Miraldi, E.R., Bonneau, R.: Multi-study inference of regulatory networks for more accurate models of gene regulation. *PLoS Comput. Biol.* **15**(1), 1006591 (2019)
46. ENCODE Project Consortium, Moore, J.E., Purcaro, M.J., Pratt, H.E., Epstein, C.B., Shores, N., Adrian, J., Kawli, T., Davis, C.A., Dobin, A., Kaul, R., Halow, J., Van Nostrand, E.L., Freese, P., Gorkin, D.U., Shen, Y., He, Y., Mackiewicz, M., Pauli-Behn, F., Williams, B.A., Mortazavi, A., Keller, C.A., Zhang, X.-O., Elhajjajay, S.I., Huey, J., Dickel, D.E., Snetkova, V., Wei, X., Wang, X., Rivera-Mulia, J.C., Rozowsky, J., Zhang, J., Chhetri, S.B., Zhang, J., Victorson, A., White, K.P., Visel, A., Yeo, G.W., Burge, C.B., Lécuyer, E., Gilbert, D.M., Dekker, J., Rinn, J., Mendenhall, E.M., Ecker, J.R., Kellis, M., Klein, R.J., Noble, W.S., Kundaje, A., Guigó, R., Farnham, P.J., Cherry, J.M., Myers, R.M., Ren, B., Graveley, B.R., Gerstein, M.B., Pennacchio, L.A., Snyder, M.P., Bernstein, B.E., Wold, B., Hardison, R.C., Gingeras, T.R., Stamatoyannopoulos, J.A., Weng, Z.: Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**(7818), 699–710 (2020)
47. Jariani, A., Vermeersch, L., Cerulus, B., Perez-Samper, G., Voordeckers, K., Van Brussel, T., Thienpont, B., Lambrechts, D., Verstrepen, K.J.: A new protocol for single-cell RNA-seq reveals stochastic gene expression during lag phase in budding yeast. *Elife* **9** (2020)
48. Sun, Z., Xu, X., He, J., Murray, A., Sun, M.-a., Wei, X., Wang, X., McCoig, E., Xie, E., Jiang, X., Li, L., Zhu, J., Chen, J., Morozov, A., Pickrell, A.M., Theus, M.H., Xie, H.: Egr1 recruits tet1 to shape the brain methylome during development and upon neuronal activity. *Nature Communications* **10**(1), 3892 (2019). doi:10.1038/s41467-019-11905-3
49. Liu, J., Pasini, S., Shelanski, M.L., Greene, L.A.: Activating transcription factor 4 (ATF4) modulates

- post-synaptic development and dendritic spine morphology. *Front. Cell. Neurosci.* **8**, 177 (2014)
50. Saito, Y., Miranda-Rottmann, S., Ruggiu, M., Park, C.Y., Fak, J.J., Zhong, R., Duncan, J.S., Fabella, B.A., Junge, H.J., Chen, Z., Araya, R., Fritsch, B., Hudspeth, A.J., Darnell, R.B.: NOVA2-mediated RNA regulation is required for axonal pathfinding during development. *Elife* **5** (2016)
 51. Wortel, I.M.N., van der Meer, L.T., Kilberg, M.S., van Leeuwen, F.N.: Surviving stress: Modulation of atf4-mediated stress responses in normal and malignant cells. *Trends in Endocrinology & Metabolism* **28**(11), 794–806 (2017). doi:10.1016/j.tem.2017.07.003
 52. Corona, C., Pasini, S., Liu, J., Amar, F., Greene, L.A., Shelanski, M.L.: Activating transcription factor 4 (atf4) regulates neuronal activity by controlling gababr trafficking. *Journal of Neuroscience* **38**(27), 6102–6113 (2018). doi:10.1523/JNEUROSCI.3350-17.2018. <https://www.jneurosci.org/content/38/27/6102.full.pdf>
 53. Oh, C.-K., Choi, Y.K., Hwang, I.-Y., Ko, Y.U., Chung, I.K., Yun, N., Oh, Y.J.: RING-finger protein 166 plays a novel pro-apoptotic role in neurotoxin-induced neurodegeneration via ubiquitination of XIAP. *Cell Death Dis.* **11**(10), 939 (2020)
 54. Culotta, L., Scalmani, P., Vinci, E., Terragni, B., Sessa, A., Broccoli, V., Mantegazza, M., Boeckers, T., Verpelli, C.: Sult4a1 modulates synaptic development and function by promoting the formation of psd-95/nmdar complex. *Journal of Neuroscience* **40**(37), 7013–7026 (2020). doi:10.1523/JNEUROSCI.2194-19.2020. <https://www.jneurosci.org/content/40/37/7013.full.pdf>
 55. Papatsenko, D., Levine, M.S.: Dual regulation by the hunchback gradient in the drosophila embryo. *Proc. Natl. Acad. Sci. U. S. A.* **105**(8), 2901–2906 (2008)
 56. Lam, K.Y., Westrick, Z.M., Müller, C.L., Christiaen, L., Bonneau, R.: Fused regression for multi-source gene regulatory network inference. *PLoS Comput. Biol.* **12**(12), 1005157 (2016)
 57. Harris, C.R., Millman, K.J., van der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N.J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M.H., Brett, M., Haldane, A., Del Río, J.F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., Oliphant, T.E.: Array programming with NumPy. *Nature* **585**(7825), 357–362 (2020)
 58. Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S.J., Brett, M., Wilson, J., Millman, K.J., Mayorov, N., Nelson, A.R.J., Jones, E., Kern, R., Larson, E., Carey, C.J., Polat, İ., Feng, Y., Moore, E.W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E.A., Harris, C.R., Archibald, A.M., Ribeiro, A.H., Pedregosa, F., van Mulbregt, P., Vijaykumar, A., Bardelli, A.P., Rothberg, A., Hilboll, A., Kloeckner, A., Scopatz, A., Lee, A., Rokem, A., Woods, C.N., Fulton, C., Masson, C., Häggström, C., Fitzgerald, C., Nicholson, D.A., Hagen, D.R., Pasechnik, D.V., Olivetti, E., Martin, E., Wieser, E., Silva, F., Lenders, F., Wilhelm, F., Young, G., Price, G.A., Ingold, G.-L., Allen, G.E., Lee, G.R., Audren, H., Probst, I., Dietrich, J.P., Silterra, J., Webber, J.T., Slavič, J., Nothman, J., Buchner, J., Kulick, J., Schönberger, J.L., de Miranda Cardoso, J.V., Reimer, J., Harrington, J., Rodríguez, J.L.C., Nunez-Iglesias, J., Kuczynski, J., Tritz, K., Thoma, M., Newville, M., Kümmerer, M., Bolingbroke, M., Tartre, M., Pak, M., Smith, N.J., Nowaczyk, N., Shebanov, N., Pavlyk, O., Brodtkorb, P.A., Lee, P., McGibbon, R.T., Feldbauer, R., Lewis, S., Tygier, S., Sievert, S., Vigna, S., Peterson, S., More, S., Pudlik, T., Oshima, T., Pingel, T.J., Robitaille, T.P., Spura, T., Jones, T.R., Cera, T., Leslie, T., Zito, T., Krauss, T., Upadhyay, U., Halchenko, Y.O., Vázquez-Baeza, Y., SciPy 1.0 Contributors: SciPy 1.0: fundamental algorithms for scientific computing in python. *Nat. Methods* (2020)
 59. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, É.: Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.* **12**(Oct), 2825–2830 (2011)
 60. Quinlan, A.R., Hall, I.M.: BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**(6), 841–842 (2010)
 61. Dale, R.K., Pedersen, B.S., Quinlan, A.R.: Pybedtools: a flexible python library for manipulating genomic datasets and annotations. *Bioinformatics* **27**(24), 3423–3424 (2011)
 62. Lambert, S.A., Yang, A.W.H., Sasse, A., Cowley, G., Albu, M., Caddick, M.X., Morris, Q.D., Weirauch, M.T., Hughes, T.R.: Similarity regression predicts evolution of transcription factor sequence specificity. *Nat. Genet.* **51**(6), 981–989 (2019)
 63. Fornes, O., Castro-Mondragon, J.A., Khan, A., van der Lee, R., Zhang, X., Richmond, P.A., Modi, B.P., Correard, S., Gheorghe, M., Baranašić, D., Santana-García, W., Tan, G., Chèneby, J., Ballester, B., Parcy, F., Sandelin, A., Lenhard, B., Wasserman, W.W., Mathelier, A.: JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **48**(D1), 87–92 (2020)
 64. Matys, V., Kel-Margoulis, O.V., Fricke, E., Liebich, I., Land, S., Barre-Dirrie, A., Reuter, I., Chekmenev, D., Krull, M., Hornischer, K., Voss, N., Stegmaier, P., Lewicki-Potapov, B., Saxel, H., Kel, A.E., Wingender, E.: TRANSFAC and its module TRANSCOMP: transcriptional gene regulation in eukaryotes. *Nucleic Acids Res.* **34**(Database issue), 108–110 (2006)
 65. Teixeira, M.C., Monteiro, P.T., Palma, M., Costa, C., Godinho, C.P., Pais, P., Cavalheiro, M., Antunes, M., Lemos, A., Pedreira, T., Sá-Correia, I.: YEASTRACT: an upgraded database for the analysis of transcription regulatory networks in *saccharomyces cerevisiae*. *Nucleic Acids Res.* **46**(D1), 348–353 (2018)
 66. Monteiro, P.T., Oliveira, J., Pais, P., Antunes, M., Palma, M., Cavalheiro, M., Galocha, M., Godinho, C.P., Martins, L.C., Bourbon, N., Mota, M.N., Ribeiro, R.A., Viana, R., Sá-Correia, I., Teixeira, M.C.: YEASTRACT+: a portal for cross-species comparative genomics of transcription regulation in yeasts. *Nucleic Acids Res.* **48**(D1), 642–649 (2020)
 67. Grant, C.E., Bailey, T.L., Noble, W.S.: FIMO: scanning for occurrences of a given motif. *Bioinformatics* **27**(7), 1017–1018 (2011)
 68. Bailey, T.L., Boden, M., Buske, F.A., Frith, M., Grant, C.E., Clementi, L., Ren, J., Li, W.W., Noble, W.S.: MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.* **37**(Web Server issue), 202–8 (2009)
 69. Kim, J.T., Martinetz, T., Polani, D.: Bioinformatic principles underlying the information content of transcription factor binding sites. *J. Theor. Biol.* **220**(4), 529–544 (2003)

70. Ester, M., Kriegel, H.-P., Sander, J., Xu, X.: A density-based algorithm for discovering clusters in large spatial databases with noise. In: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining. KDD'96, pp. 226–231. AAAI Press, ??? (1996)
71. Fu, Y., Jarboe, L.R., Dickerson, J.A.: Reconstructing genome-wide regulatory network of *e. coli* using transcriptome data and predicted transcription factor activities. *BMC Bioinformatics* **12**, 233 (2011)
72. Madar, A., Greenfield, A., Vanden-Eijnden, E., Bonneau, R.: DREAM3: Network inference using dynamic context likelihood of relatedness and the inferelator. *PLoS One* **5**(3), 9803 (2010)
73. Greenfield, A., Hafemeister, C., Bonneau, R.: Robust data-driven incorporation of prior knowledge into the inference of dynamic regulatory networks. *Bioinformatics* **29**(8), 1060–1067 (2013)
74. Wes McKinney: Data Structures for Statistical Computing in Python. In: Stéfan van der Walt, Jarrod Millman (eds.) Proceedings of the 9th Python in Science Conference, pp. 56–61 (2010). doi:10.25080/Majora-92bf1922-00a
75. Schacht, T., Oswald, M., Eils, R., Eichmüller, S.B., König, R.: Estimating the activity of transcription factors by the effect on their target genes. *Bioinformatics* **30**(17), 401–7 (2014)
76. Zou, H., Hastie, T.: Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Series B Stat. Methodol.* **67**(2), 301–320 (2005)
77. Zou, H.: The adaptive lasso and its oracle properties. *J. Am. Stat. Assoc.* **101**(476), 1418–1429 (2006)
78. Liu, H., Roeder, K., Wasserman, L.: Stability approach to regularization selection (StARS) for high dimensional graphical models (2010). 1006.3316
79. Matthews, B.W.: Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim. Biophys. Acta* **405**(2), 442–451 (1975)
80. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., Smyth, G.K.: limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**(7), 47 (2015)
81. 10x Genomics: 1.3 Million Brain Cells from E18 Mice. https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.3.0/1M_neurons (2017)
82. Satija, R., Farrell, J.A., Gennert, D., Schier, A.F., Regev, A.: Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**(5), 495–502 (2015)
83. Di Bella, D.J., Habibi, E., Yang, S.-M., Stickels, R.R., Brown, J., Yadollahpour, P., Chen, F., Macosko, E.Z., Regev, A., Arlotta, P.: Molecular Logic of Cellular Diversification in the Mammalian Cerebral Cortex (2020)
84. 10x Genomics: Fresh cortex, hippocampus, and ventricular zone from embryonic mouse brain (E18). https://support.10xgenomics.com/single-cell-atac/datasets/1.2.0/atac_v1_E18_brain_fresh_5k (2019)
85. 10x Genomics: Dissociated and cryopreserved cortex, hippocampus, and ventricular zone cells from embryonic mouse brain (E18). https://support.10xgenomics.com/single-cell-atac/datasets/1.2.0/atac_v1_E18_brain_cryo_5k (2019)
86. 10x Genomics: Flash frozen cortex, hippocampus, and ventricular zone from embryonic mouse brain (E18). https://support.10xgenomics.com/single-cell-atac/datasets/1.2.0/atac_v1_E18_brain_flash_5k (2019)
87. Gabitto, M.I., Rasmussen, A., Wapinski, O., Allaway, K., Carriero, N., Fishell, G.J., Bonneau, R.: Characterizing chromatin landscape from aggregate and single-cell genomic assays using flexible duration modeling. *Nature Communications* **11**(1), 747 (2020). doi:10.1038/s41467-020-14497-5
88. Fang, R., Preissl, S., Li, Y., Hou, X., Lucero, J., Wang, X., Motamedi, A., Shiau, A.K., Zhou, X., Xie, F., Mukamel, E.A., Zhang, K., Zhang, Y., Behrens, M.M., Ecker, J.R., Ren, B.: Comprehensive analysis of single cell atac-seq data with snapatac. *Nature Communications* **12**(1), 1337 (2021). doi:10.1038/s41467-021-21583-9
89. R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2020). R Foundation for Statistical Computing. <https://www.R-project.org/>
90. Wickham, H.: Ggplot2: Elegant Graphics for Data Analysis. Springer, ??? (2016). <https://ggplot2.tidyverse.org>
91. McInnes, L., Healy, J., Melville, J.: UMAP: Uniform manifold approximation and projection for dimension reduction. arXiv:1802.03426 [cs, stat] (2018)
92. Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L., François, R., Grolemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T., Miller, E., Bache, S., Müller, K., Ooms, J., Robinson, D., Seidel, D., Spinu, V., Takahashi, K., Vaughan, D., Wilke, C., Woo, K., Yutani, H.: Welcome to the tidyverse. *J. Open Source Softw.* **4**(43), 1686 (2019)
93. Hunter, J.D.: Matplotlib: A 2D graphics environment. *Computing in Science Engineering* **9**(3), 90–95 (2007)
94. Waskom, M.L.: seaborn: statistical data visualization. *Journal of Open Source Software* **6**(60), 3021 (2021). doi:10.21105/joss.03021
95. Watters, A.: `jp.gene.viz`. https://github.com/simonsfoundation/jp_gene_viz (2019)

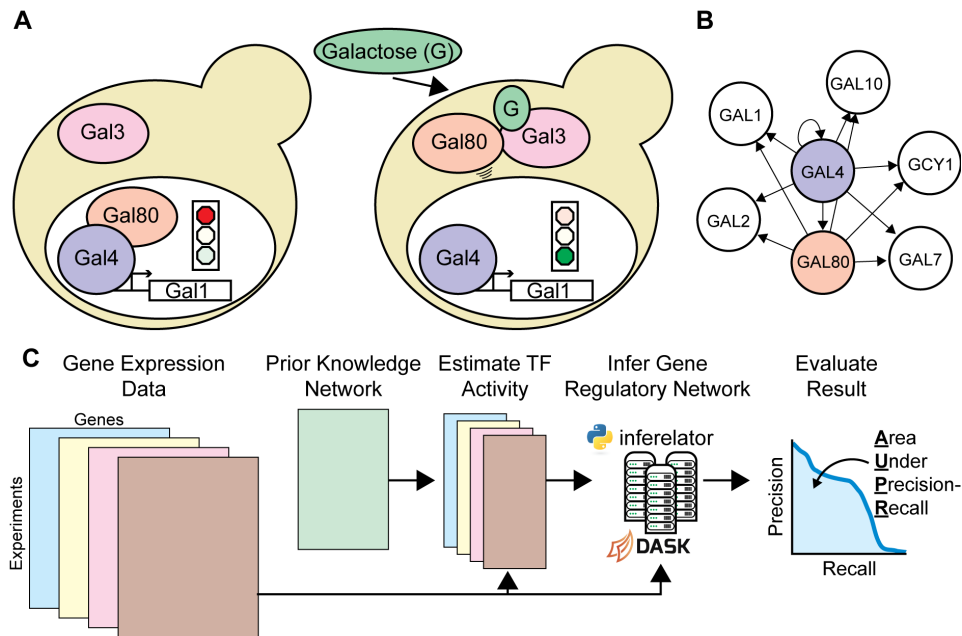


Figure 1: Learning Gene Regulatory Networks with the Inferelator (A) The response to the sugar galactose in *Saccharomyces cerevisiae* is mediated by the Gal4 and Gal80 TFs, a prototypical mechanism for altering cellular gene expression in response to stimuli. (B) Gal4 and Gal80 regulation represented as an unsigned directed graph connecting regulatory TFs to target genes. (C) Genome-wide Gene Regulatory Networks (GRNs) are inferred from gene expression data and prior knowledge about network connections using the Inferelator, and the resulting networks are scored by comparison with a gold standard of known interactions.

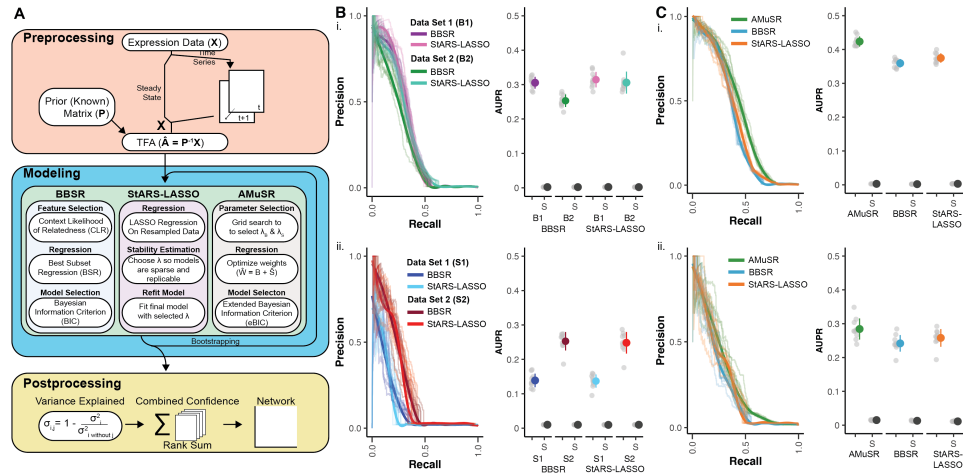


Figure 2: Network Inference Performance on Multiple Model Organism Datasets (A) Schematic of Inferelator workflow and a brief summary of the differences between GRN model selection methods (B) Results from 10 replicates of GRN inference for each modeling method on (i) *Bacillus subtilis* GSE67023 (B1), GSE27219 (B2) and (ii) *Saccharomyces cerevisiae* GSE142864 (S1), and [40] (S2). Precision-recall curves are shown for replicates where 20% of genes are held out of the prior and used for evaluation, with a smoothed consensus curve. AUPR is plotted for each cross-validation result in gray, with mean \pm standard deviation in color. Experiments labeled with (S) are shuffled controls, where the labels on the prior adjacency matrix have been randomly shuffled. 10 shuffled replicates are shown as gray dots, with mean \pm standard deviation in black. (C) Results from 10 replicates of GRN inference using two datasets as two network inference tasks on (i) *Bacillus subtilis* and (ii) *Saccharomyces cerevisiae*. AMuSR is a multi-task learning method; BBSR and StARS-LASSO are run on each task separately and then combined into a unified GRN. Precision-recall curves and AUPR are plotted as in B.

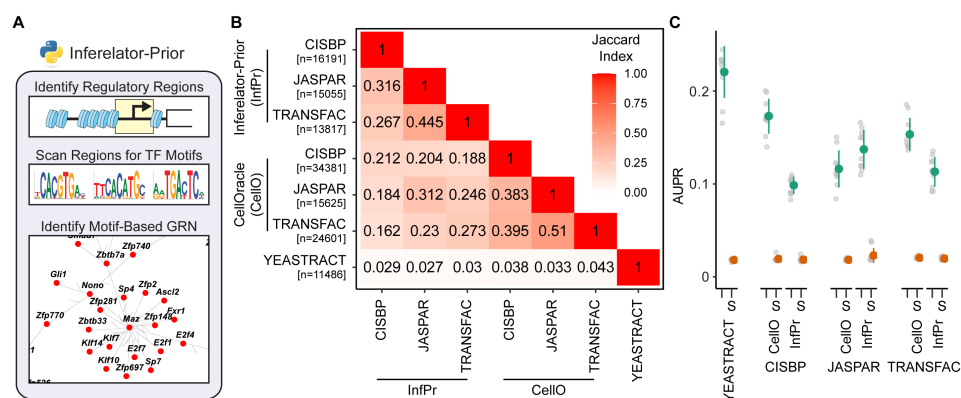


Figure 3: Construction and Performance of Network Connectivity Priors Using TF Motif Scanning **(A)** Schematic of inferelator-prior workflow, scanning identified regulatory regions (e.g. by ATAC) for TF motifs to construct adjacency matrices **(B)** Jaccard similarity index between *Saccharomyces cerevisiae* prior adjacency matrices generated by the inferelator-prior package, by the CellOracle package, and obtained from the YEASTRACT database. Prior matrices were generated using TF motifs from the CIS-BP, JASPAR, and TRANSFAC databases with each pipeline (n is the number of edges in each prior adjacency matrix). **(C)** The performance of Inferelator network inference using each motif-derived prior. Performance is evaluated by AUPR, scoring against genes held out of the prior adjacency matrix, based on inference using 2577 genome-wide microarray experiments. Experiments labeled with (S) are shuffled controls, where the labels on the prior adjacency matrix have been randomly shuffled.

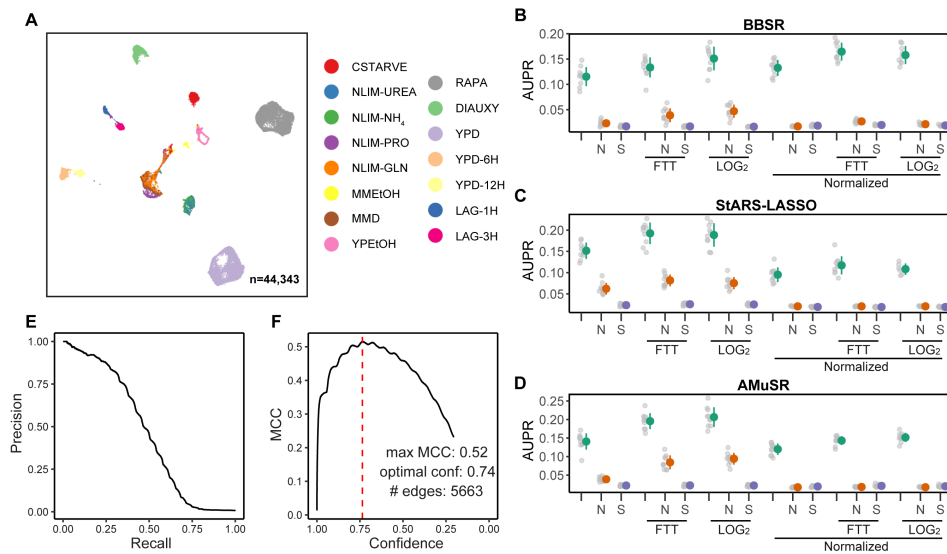


Figure 4: Network Inference Performance Using Single-Cell *Saccharomyces cerevisiae* Expression Data (A) Uniform Manifold Approximation and Projection (UMAP) plot of single-cell yeast data, colored by the experimental grouping of individual cells (tasks). (B) The effect of preprocessing methods on network inference using BBSR model selection on 14 task-specific expression datasets, as measured by AUPR. Colored dots represent mean \pm standard deviation of all replicates. Data is either untransformed (raw counts), transformed by Freeman-Tukey Transform (FTT), or transformed by $\log_2(x_1)$ pseudocount. Non-normalized data is compared to data normalized so that all cells have identical count depth. Network inference performance is compared to two baseline controls; data which has been replaced by Gaussian noise (N) and network inference using shuffled labels in the prior network (S). (C) Performance evaluated as in B on StARS-LASSO model selection. (D) Performance evaluated as in B on AMuSR model selection. (E) Precision-recall of the recovery of the prior on a final network constructed using FTT-transformed, non-normalized AMuSR model selection. (F) Matthews Correlation Coefficient (MCC) of the same network as in E

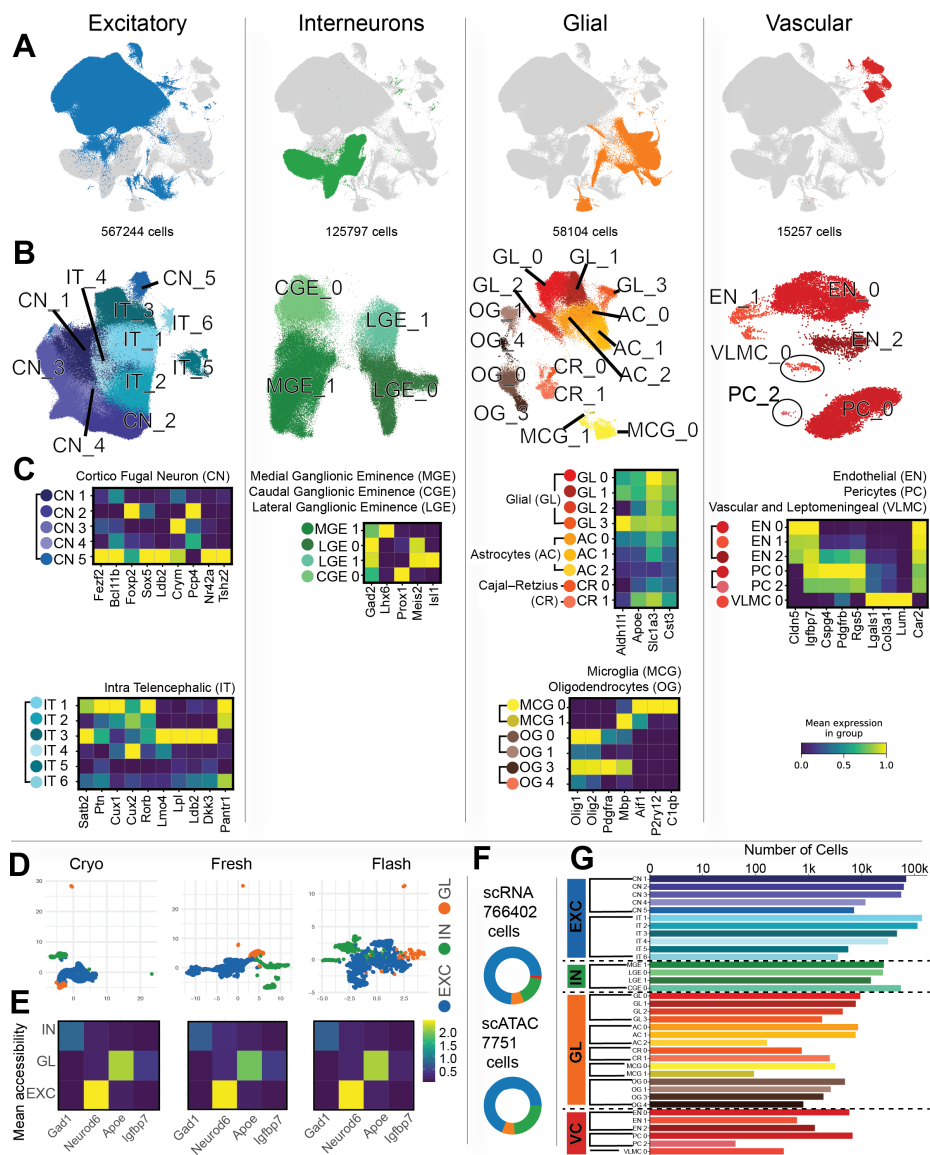


Figure 5: Processing Large Single-Cell Mouse Brain Data for Network Inference (A) UMAP plot of all mouse brain scRNA-seq data with Excitatory neurons, Interneurons, Glial cells and Vascular cells colored. (B) UMAP plot of cells from each broad category colored by Louvain clusters and labeled by cell type. (C) Heatmap of normalized gene expression for marker genes that distinguish cluster cell types within broad categories. (D) UMAP plot of mouse brain scATAC data with Excitatory neurons, Interneurons, and Glial cells colored. (E) Heatmap of normalized mean gene accessibility for marker genes that distinguish broad categories of cells. (F) The number of scRNA-seq and scATAC cells in each of the broad categories. (G) The number of scRNA-seq cells in each cell type specific cluster.

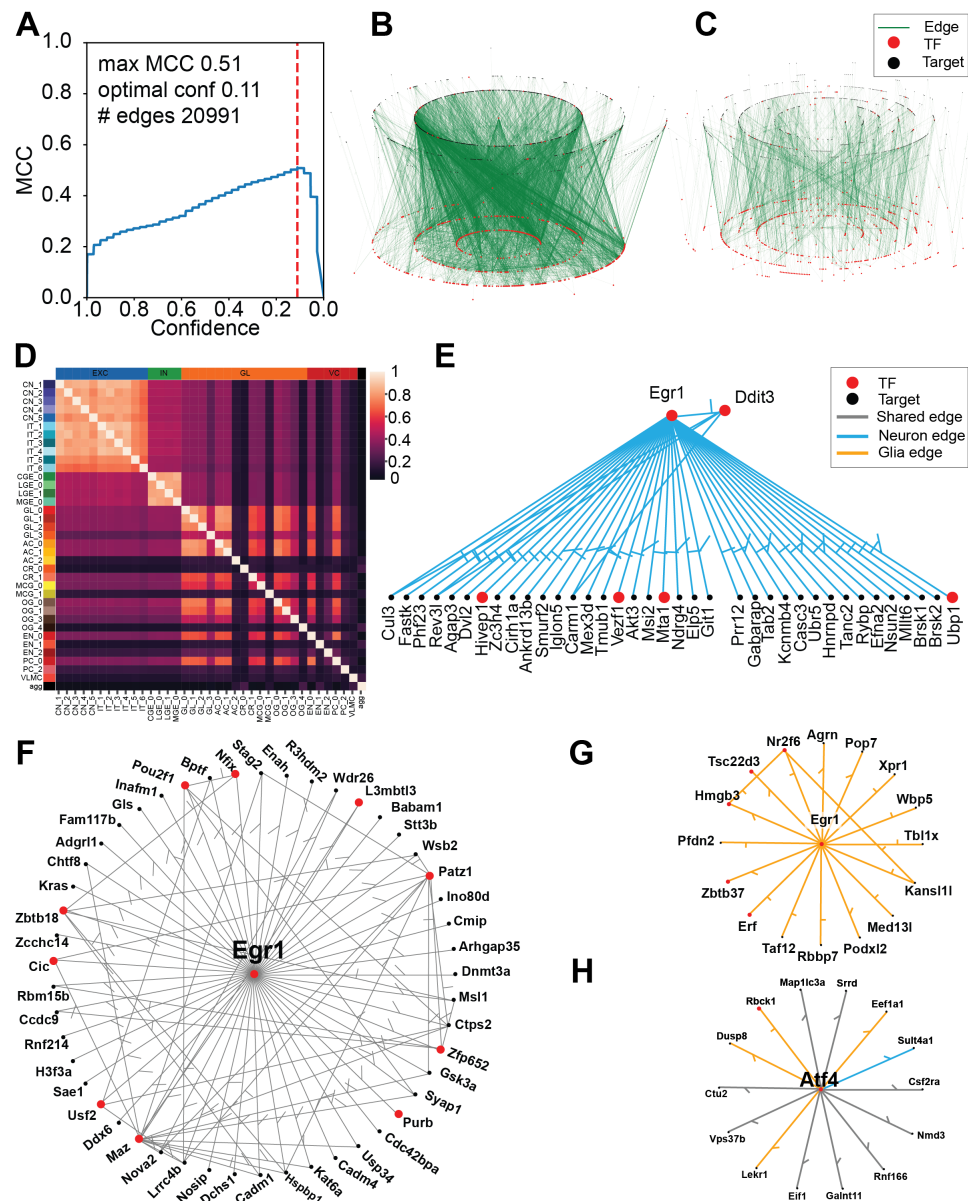
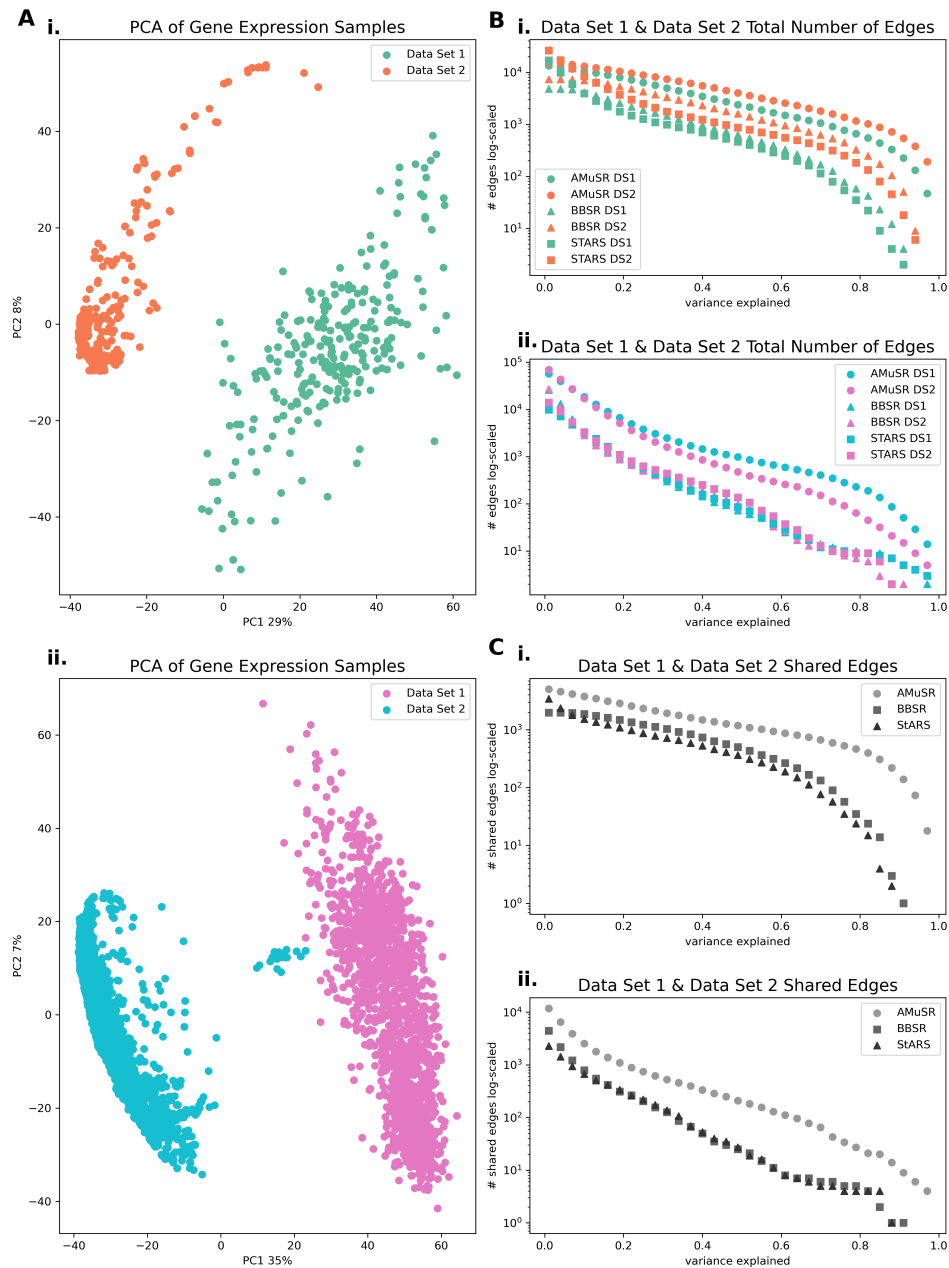
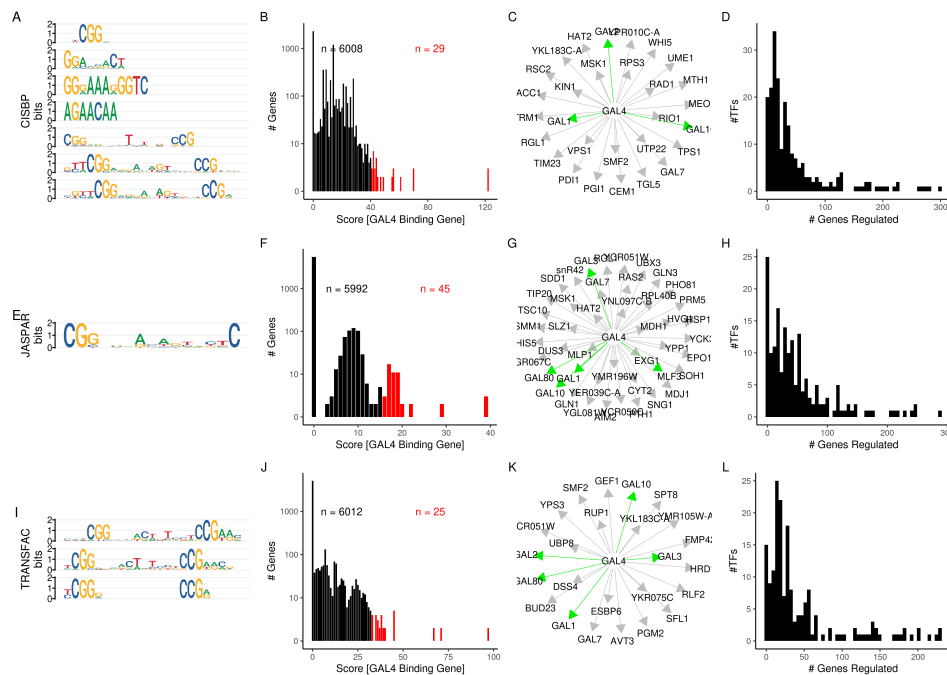


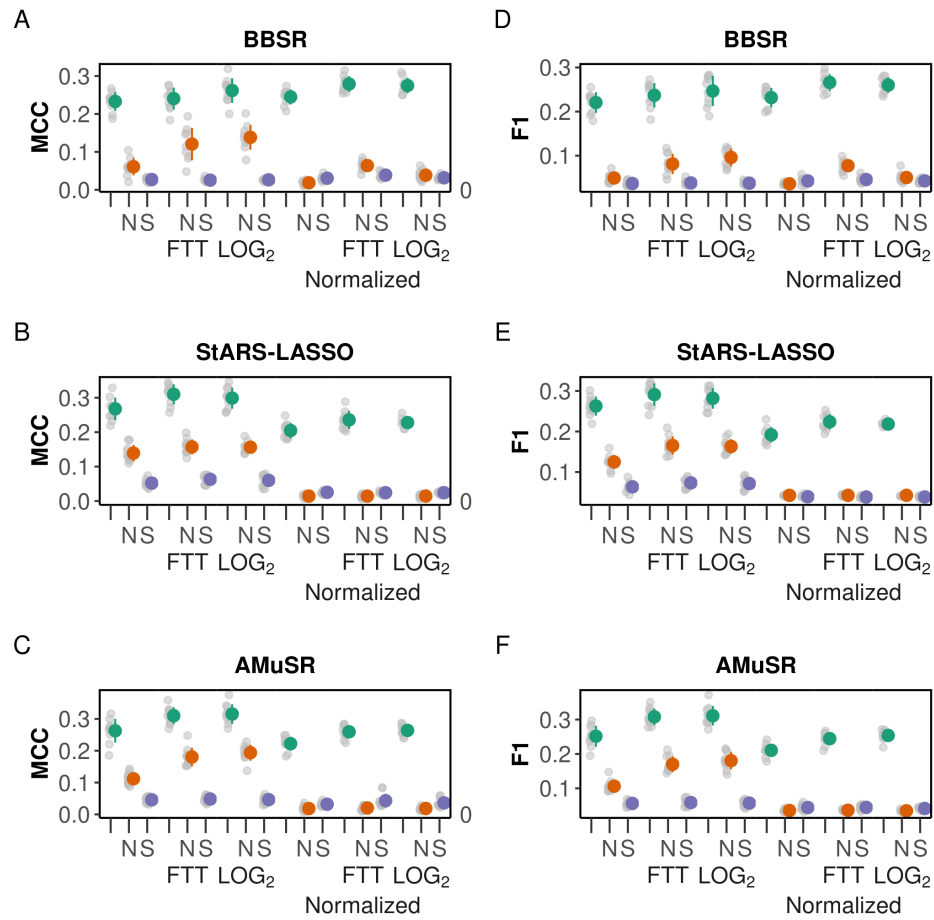
Figure 6: Learned GRN For The Mouse Brain (A) MCC for the aggregate network based on Inferelator prediction confidence. The dashed line shows the confidence score which maximizes MCC. Network edges at and above this line are retained in the final network. (B) Aggregate GRN learned. (C) Network edges which are present in every individual task. (D) Jaccard similarity index between each task network (E) Network targets of the *EGR1* TF in neurons. (F) Network targets of the *EGR1* TF in both neurons and glial cells. (G) Network targets of the *EGR1* TF in glial cells. (H) Network of the *ATF4* TF where blue edges are neuron specific, orange edges are glial specific, and black edges are present in both categories.



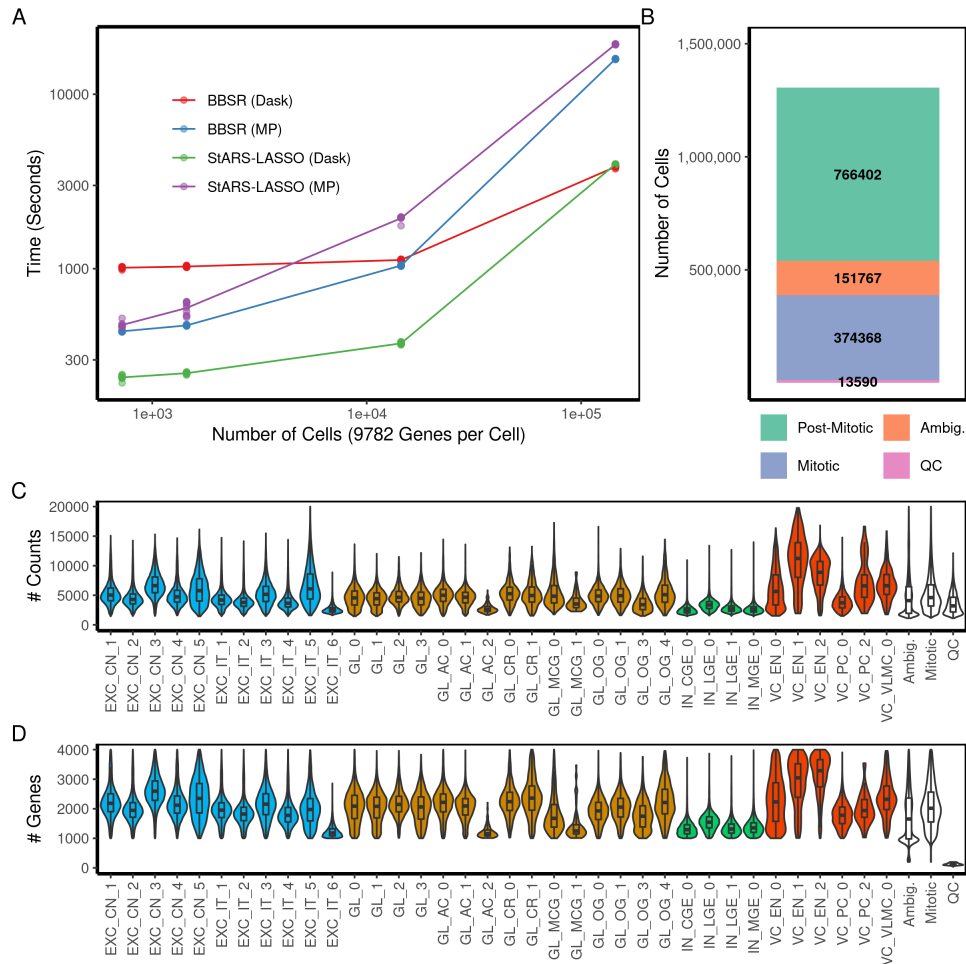
Supplemental Figure 1: Learning *Bacillus subtilis* and *Saccharomyces cerevisiae* networks by tasks. (A) PCA depicts batch effects between datasets for both (i) *Bacillus subtilis* and (ii) *Saccharomyces cerevisiae*. Learning networks by treating the independently collected datasets as separate tasks allows for sharing regulatory commonalities while respecting experimental variance. (B) The number of shared edges between the two datasets, for both model organisms (i) and (ii), shows a high number of overlapping edges. Edges are ranked by their corresponding variance explained for each of the three different model selection approaches: AMuSR, BBSR, and StARS-LASSO. (C) Across the three different model selection approaches, AMuSR learns the highest number of overlapping edges between the respective datasets for model organisms (i) and (ii).



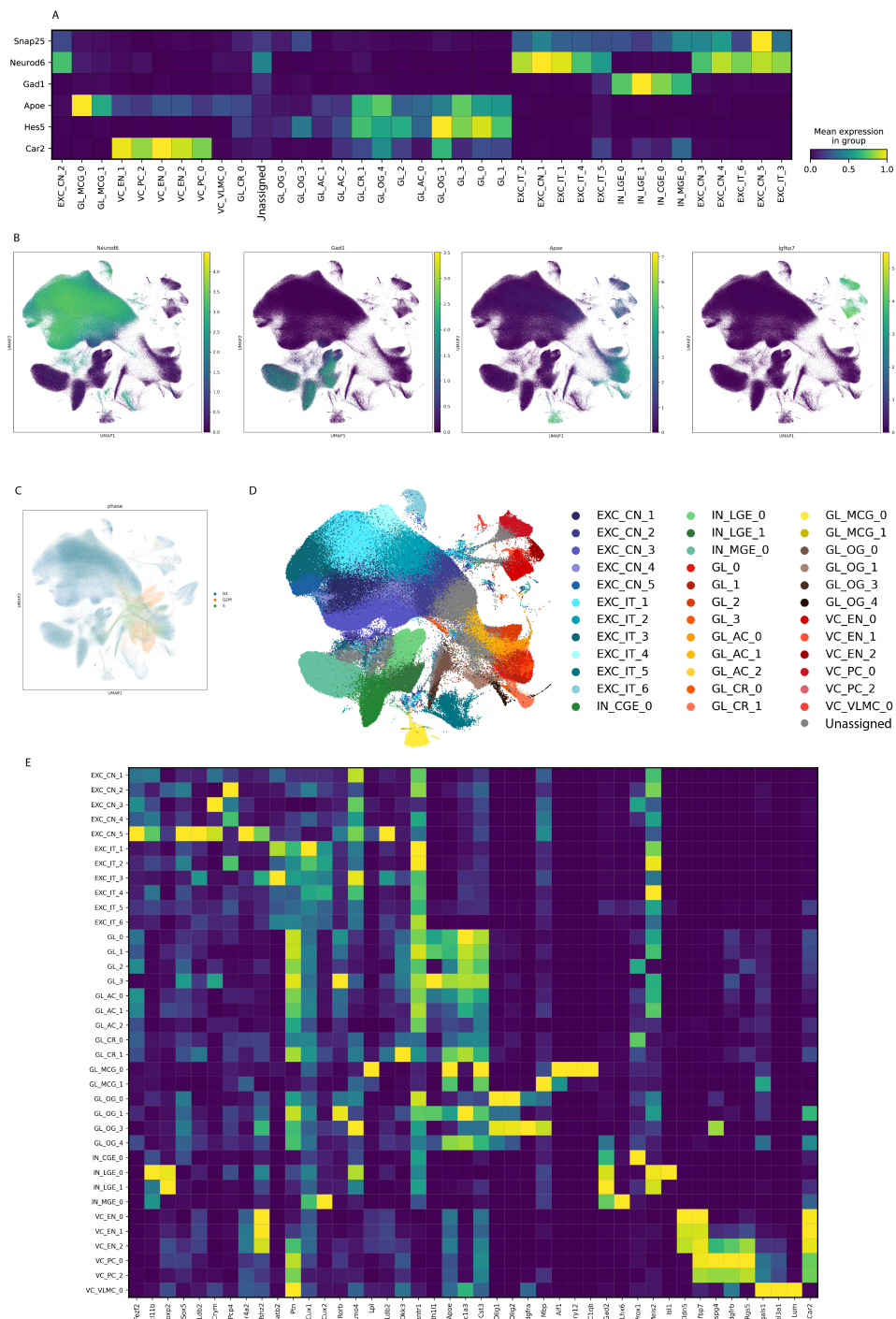
Supplemental Figure 2: Network construction using TF motifs in *Saccharomyces cerevisiae*. (A) Motifs annotated for GAL4 in the CIS-BP motif database. (B) Histogram of scores linking GAL4 to target genes. Genes in black have been omitted from the final connectivity matrix, and genes in red have been included. (C) Network connecting GAL4 and target genes. Green edges are present in the YEASTRACT database. (D) Histogram of out degree for each TF in the complete network. (E-H) Network analysis as A-D for the JASPAR motif database. (I-L) Network analysis as A-D for the TRANSFAC PRO motif database.



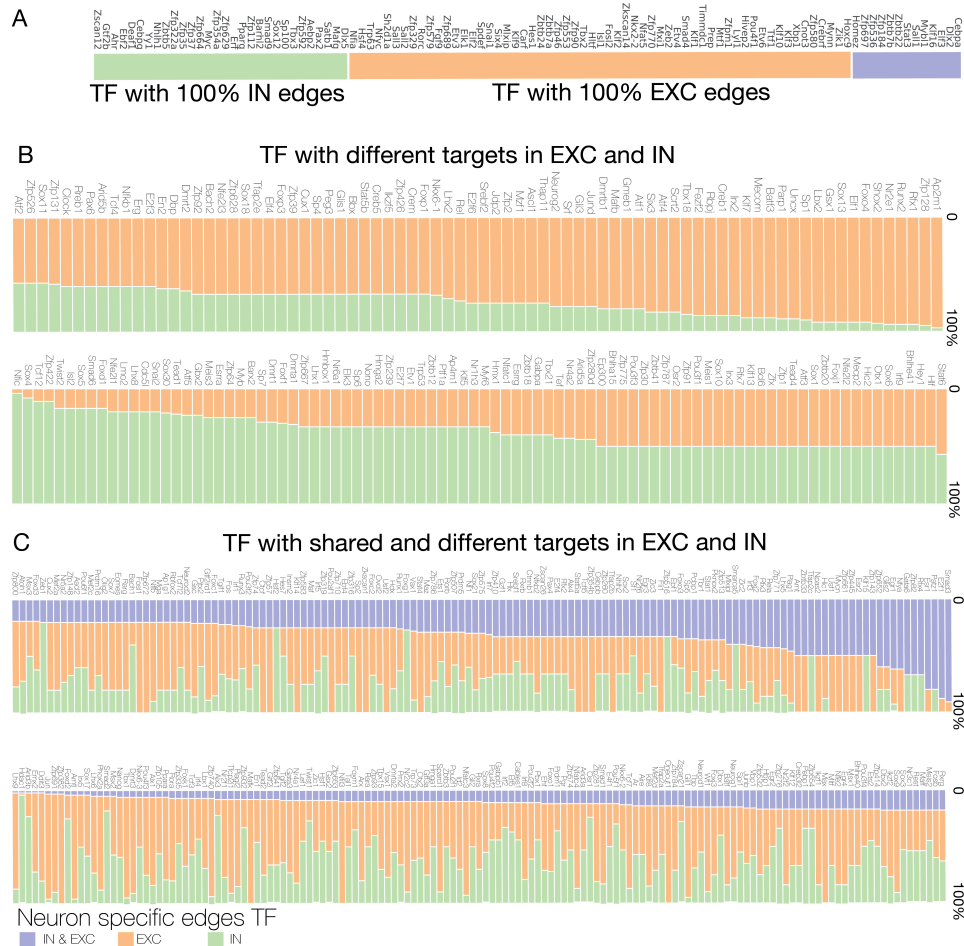
Supplemental Figure 3: Single-cell yeast network inference performance measured by Matthews Correlation Coefficient (MCC) for (A) BBSR, (B) StARS-LASSO, and (C) AMuSR. (D-F) Performance measured by F1 score as A-C



Supplemental Figure 4: (A) The Inferelator 3.0 computational performance as measured by runtime in seconds for BBSR and StARS-LASSO using the Dask engine (140 cpu cores) or using the python-based multiprocessing (MP) engine (28 cpu cores). Expression data is sampled from 144,000 mouse cells and 9,782 genes are modeled for network inference. Runtime is shown for 10 replicate runs for each quantity of cells. (B) Number of cells removed during preprocessing for Quality Control (QC), as Mitotic, and as Ambiguous by neuronal marker. Post-mitotic, non-ambiguous cells are retained and clustered. (C) Number of single-cell counts per cell in each of 36 cell type-specific groups, and in the groups removed during preprocessing. (D) Number of genes per cell in each of 36 cell type-specific groups, and in the groups removed during preprocessing



Supplemental Figure 5: (A) Cell class marker expression for each annotated subcluster in mouse single-cell brain data. (B) UMAP of 766,402 mouse brain cells colored by cell class marker expression. (C) UMAP of 1.3M mouse brain cells colored by the assigned cell cycle phase. (D) UMAP of 766,402 mouse brain cells colored by 36 assigned subcluster. (E) Cell type marker expression by assigned subcluster.



Supplemental Figure 6: (A) List of TFs that have identical target genes in GRNs for both Excitatory neurons (EXC) and Interneurons (IN), that have only target genes in Excitatory neurons, and that have only target genes in Interneurons. (B) List of TFs that have no shared target genes in GRNs for Excitatory neurons and in GRNs for interneurons. (C) TFs that have some shared target genes in GRNs for Excitatory neurons and interneurons, but also have some target genes specific to Excitatory neurons or interneurons.