

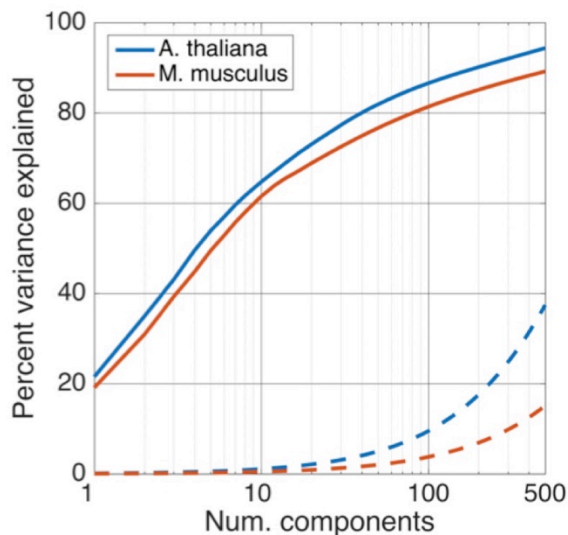
# Tradict enables high fidelity reconstruction of the eukaryotic transcriptome from 100 marker genes

## Supplemental Information

Surojit Biswas, Konstantin Kerner, Paulo José Pereira Lima Teixeira, Jeffery L. Dangl, Vladimir Jojic, Philip A. Wigge

### Supplemental Note 1 - Our training transcriptomes are reflective of biology and are of high technical quality.

We manually annotated metadata for 1,626 (62.6%, *A. thaliana*) and 6,682 (32.1%, *M. musculus*) of the training transcriptomes for both organisms, and found that the major drivers of variation were tissue and developmental stage (Figure 1a-b, main text). The first three principal components of our training collection explained a substantial proportion of expression variation for each organism (43.1% *A. thaliana*, 39.3% *M. musculus*). For *A. thaliana* PC1 was primarily aligned with the physical axis of the plant, with above ground, photosynthetic tissues having lower PC1 scores and below ground, root tissues having higher PC1 scores. Interestingly, samples found intermediate to the major below- and above-ground tissue clusters consisted of seedlings grown in constant darkness or mutant seedlings (e.g. *det1*, *pif*, *phy*) compromised for photomorphogenesis. Thus, PC1 can also be considered to align with light perception and signaling. By contrast, PC2 represented a developmental axis, with more embryonic tissues (seeds, endosperms) having lower PC2 scores, and more developed tissues having higher PC2 scores (Figure 1a, main text).

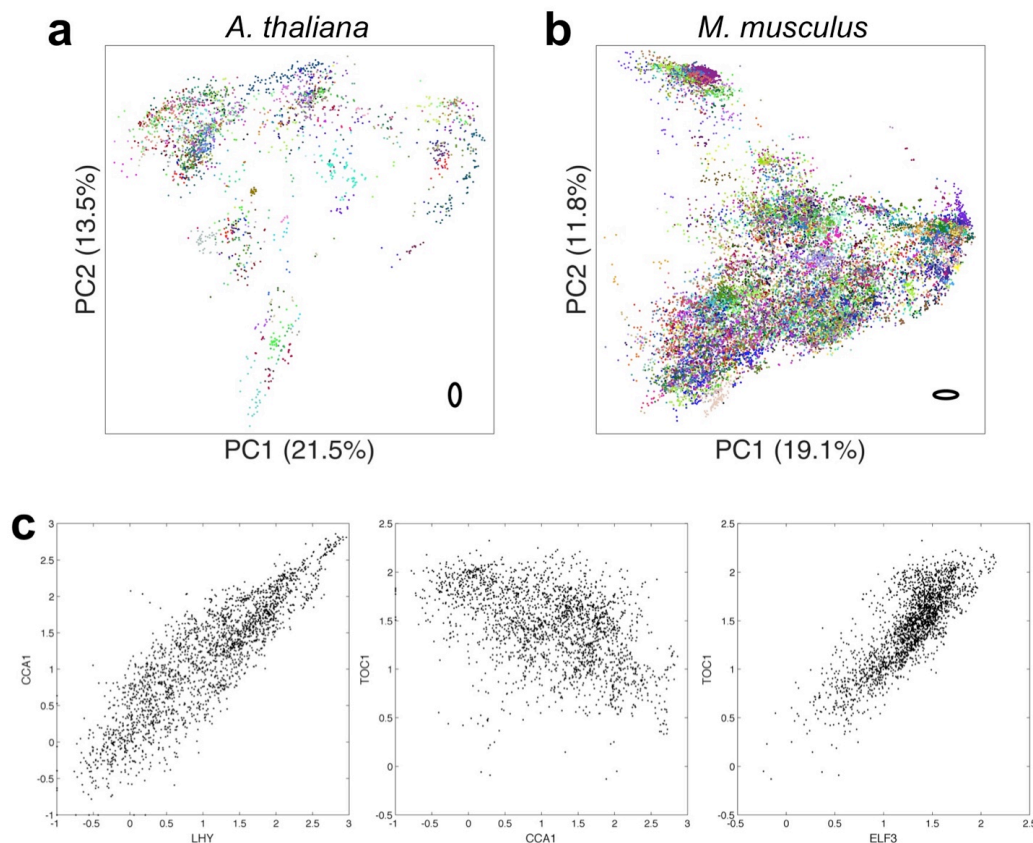


**Figure S1. The eukaryotic transcriptome is compressible.** The transcriptome is of low dimensionality, with 100 principal components able to explain 80% or more of expression variation. Dotted lines illustrate cumulative expression variation explained on a null model realization, where each gene's expression vector was permuted to break correlative ties to other genes.

For *M. musculus*, PC1 described a hematopoietic-nervous system axis. Cardiovascular, digestive, respiratory, urinary and connective tissues were found intermediate along this axis, and with the exception of liver tissue, were not differentiable along the first three PCs. Interestingly, as observed for *A. thaliana*, PC2 represented a developmental axis, with general "stemness" decreasing with increasing PC2 score. Consistent with this trend, nervous tissue

from embryos and postnatal mice had consistently lower PC2 scores than mature nervous tissue. We did not find any significant correlation between *Xist* expression and any of the top twenty PCs, suggesting that sex was not a major driver of global gene expression relative to tissue and developmental context. This is consistent with findings reported in Crowley *et al.* (2015)<sup>1</sup>.

To understand the compressibility of our training transcriptome collection beyond the first three PCs, we examined the percent of expression variation explained by subsequent components. Strikingly, we found the first 100 principal components were sufficient to explain 86.6% and 81.4% of expression variation in the observed transcriptomes for *A. thaliana* and *M. musculus*, respectively. By contrast, the first 100 principal components of a null model realization, in which the expression vectors for each gene were independently permuted, could only explain 5-10% of expression variation for both organisms (Figure S1). Given the phylogenetic distance spanned by *A. thaliana* and *M. musculus*, this transcriptomic compressibility is likely a shared property of all eukaryotes.

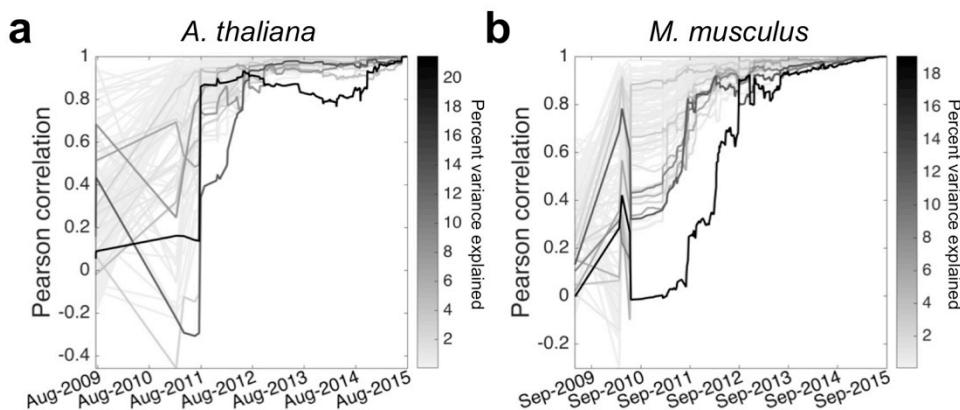


**Figure S2. Our training collection is of high technical quality.** Two dimensional principal components analysis for a) *A. thaliana* and b) *M. musculus*, where each sample is colored by the submission it belongs to. Note that while multiple submissions may have similar colors, each expression cluster contains many submissions. Bold, black ovals in the bottom left of each plot illustrate two standard deviation covariances for the median variance submission. c) Expression of late and early elements of the *A. thaliana* circadian clock matches expectations. Scatter plots of *LHY*, *CCA1*, and *ELF3* expression across all observed transcriptomes. *LHY* and *CCA1* expression is activated by *TOC1*. *CCA1* and *LHY* protein inhibits *TOC1* and *ELF3* transcription.

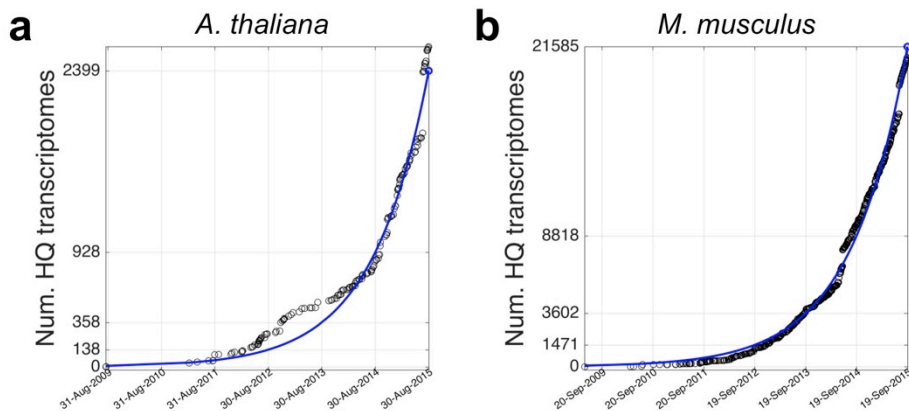
To further assess the quality and representativeness of our training collection, we examined the distribution of SRA submissions across the expression space, compared inter-submission variability within and between tissues, inspected expression correlations among

genes with well established regulatory relationships, and assessed the evolution of the expression space across time. Technical variation due to differences in laboratory procedures across labs is difficult to assess since this requires two different labs to perform the same, equivalently aimed experiment. Nevertheless, for both organisms, each tissue or development specific cluster was supported by multiple submissions, and importantly, inter-submission variability within a tissue or developmental context was significantly smaller than inter-tissue/developmental stage variability (p-value =  $1.23 \times 10^{-16}$ , F-test; Figures S2a-b). We also compared the expression of *ELF3*, *LHY*, and *TOC1* -- early and late elements of the *A. thaliana* circadian clock -- and found strong correlation in their expression with a direction and magnitude that fit established expectations (Figure S2c)<sup>2</sup>.

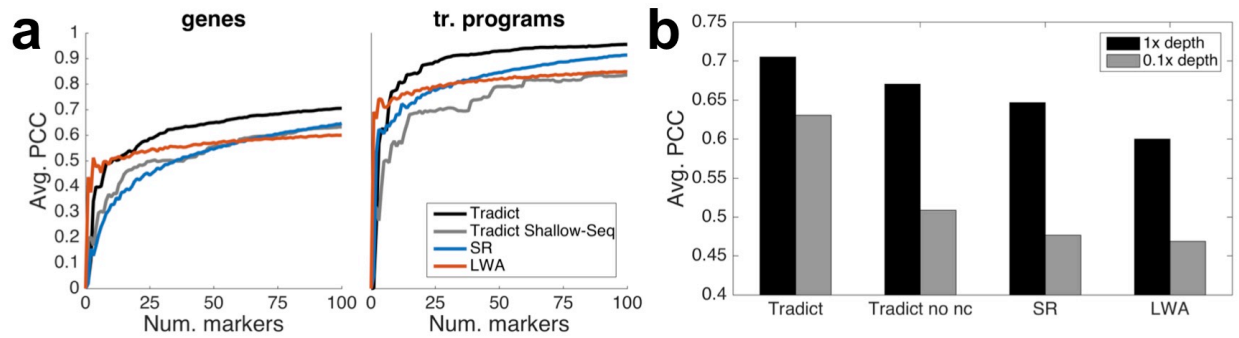
We next performed a temporal rarefaction analysis. We compared (measured by Pearson correlation) how past distributions of samples along each of the first 100 principal components compared to their present distribution. Figures S3a-b illustrate that the expression space stabilized 2-3 years ago, and that new transcriptome samples that are added to the SRA tend to fall within already established clusters. We further note that the amount of usable transcriptomic data deposited on the SRA, and hence the representativeness of our sample, is increasing exponentially (Figure S4).



**Figure S3. The expression space has stabilized.** For each of the first 100 principal components (PCs), depicted is the Pearson correlation between how samples were distributed along the PC at a select point in the past and how they are distributed currently. Each line, representing a PC, is shaded by the percent variance explained by that PC. a) *A. thaliana*. b) *M. musculus*.



**Figure S4. The number of high quality transcriptomes deposited in the SRA is growing exponentially.** SRA growth for a) *A. thaliana*, and b) *M. musculus*.



**Figure S5. Tradict outperforms leading methods and is robust to noise.** Tradict was trained on the first (historically speaking) 90% of SRA submissions and then tasked with predicting the remaining 10% of “test-set” submissions. a) Average intra-submission Pearson correlation coefficients between predicted and actual expression of genes (left) and transcriptional programs (tr. programs; right) in the test-set as a function of the number of markers used in the model. b) Intra-submission prediction accuracy of gene expression on the same test-set processed normally or rarefied to 0.1x depth. ‘Tradict no nc’ uses the same algorithm as Tradict, however, a diagonal covariance is used over markers, instead of a full one.

## Supplemental Note 2 - Tradict outperforms leading approaches and is robust to noise

As baselines for Tradict, we consider three alternative approaches. The first two, locally weighted averaging (LWA) and structured regression (SR) are the two best performing methods used in Donner *et al.* (2012)<sup>3</sup>. LWA, a non-parametric and non-linear approach, formulates predictions as weighted averages of the entire training set, where weights are determined by the distance between a query set of marker expressions and the expression of those markers in a training transcriptome. The exact weighting function is given by a Gaussian kernel, whose bandwidth we learn through cross-validation. In contrast, SR selects markers and predicts expression using regularized regression and the  $L_{0,\infty}$  objective. The appropriate level of regularization is again learned through cross-validation. Given these methods were built for use on microarray data and hence their dependence on normality, we applied them to a log-transformed version of our training collection ( $\log[\text{TPM} + 0.1]$ ).

In the third baseline (Tradict Shallow-Seq), we employ Tradict as usual; however, we restrict Tradict’s selected markers to be the 100 most abundant genes in the transcriptome. This provides a control for Tradict’s marker selection algorithm, and simulates a situation that would be typical of shallow sequencing, where only the most abundant genes are used to make conclusions about the rest of the transcriptome.

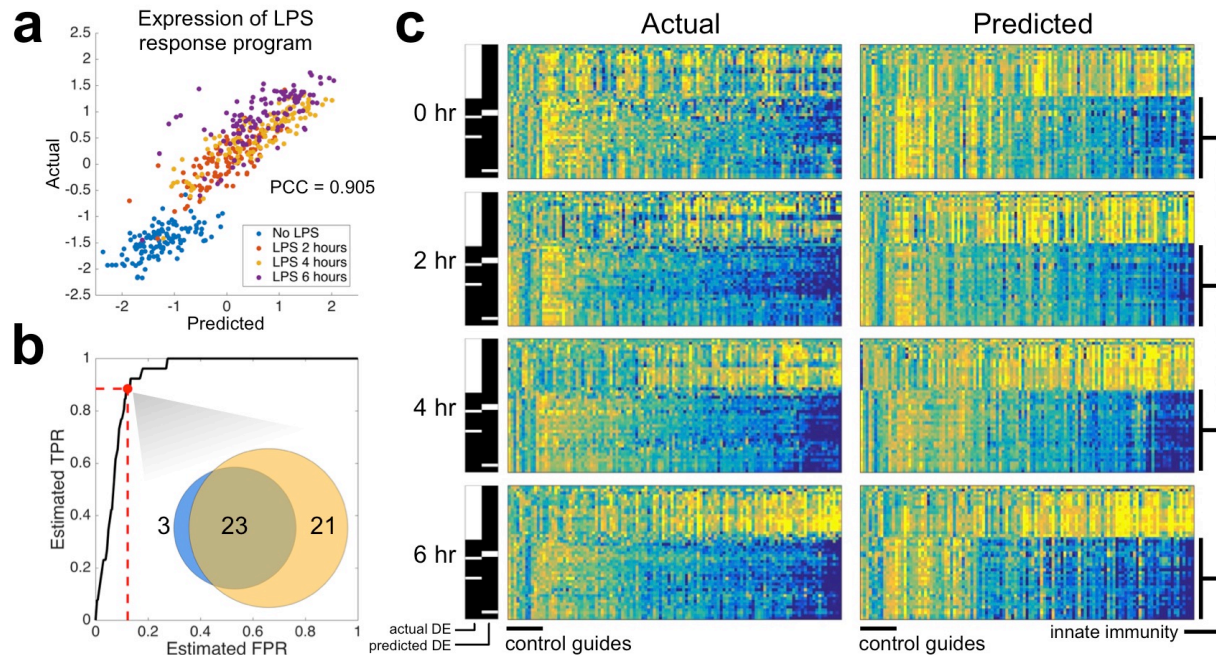
To assess Tradict’s prospective predictive performance and how it compares to the baseline models, we first partitioned our transcriptome collection for *A. thaliana* into a training set and test set by submission and historical date. In order to mimic Tradict’s use in practice as closely as possible, the training set contained the first 90% of submissions (208 submissions comprised of 2,389 samples) deposited on the SRA, and the test set contained the remaining 10% (17 submissions comprised of 208 samples). Tradict and the baseline models were each first trained on the training set. Their intra-submission predictive performance on the test set was then determined by providing only the expression values of selected markers as input, and subsequently examining the within-submission Pearson correlation coefficient (PCC) between the predicted and actual expression of transcriptional programs and the remaining >20,000 non-marker genes in the transcriptome.

Figure S5a illustrates the intra-submission performance of each method as a function of the number of markers entered into the model. LWA demonstrates the quickest performance gain, but then saturates after 10 markers. This is likely because a non-linear kernel based approach makes the cleverest use of a few markers, but is plagued by the curse of

dimensionality as more markers are added. The parametric methods (Tradict, SR) navigate this dimensionality increase more efficiently and ultimately realize better performance for still reasonable numbers of markers. Tradict outperforms SR and Tradict Shallow-Seq, ultimately obtaining a PCC between predicted and actual expression of 0.71 for genes and strikingly 0.96 for transcriptional programs. This suggests Tradict's probabilistic framework is more reasonable than SR's and that Tradict's marker selection is more optimal than picking the most abundant genes.

We noticed that though Tradict iteratively selects markers to maximize explanatory power, these markers are not orthogonal. Consequently, during inference of the marker latent abundances, on which all expression predictions are based, the internal covariance among the markers will be used during estimation. In increasing data (larger sequencing depth, higher *a priori* abundance) the latent abundance inference will place less emphasis on this internal covariance; however, in situations of measurement error or inadequacy, the internal covariance will help to learn the correct latent abundances, which in turn, should stabilize predictions in noisy situations. To test this hypothesis, we considered a version of Tradict, 'Tradict no nc' (noise correction), in which only the diagonal of the internal marker covariance was used, effectively decoupling marker abundances in Tradict's underlying model. We re-evaluated intra-submission prediction accuracy for all of the methods, excluding Tradict Shallow-Seq, on the same training and test set above using 100 markers. However this time, in order to simulate situations of high measurement error, we rarefied samples in the test set to 0.1x depth and evaluated each method's predicted (depth-normalized) expression accuracy; the original 1x depth values formed the basis of comparison. The 10<sup>th</sup>, 25<sup>th</sup>, 50<sup>th</sup>, 75<sup>th</sup>, and 90<sup>th</sup> percentiles of read depths in the 0.1x scenario were 0.65, 1.1, 2.1, 3.1, and 4.4 million reads, respectively -- all below the recommended depths for *A. thaliana*. 30-40% of the markers had zero abundance in nearly half of the samples. Figure 3b illustrates that though all methods perform worse at 0.1x depth, Tradict is least affected. Importantly, we notice that Tradict no nc's performance is substantially reduced at lower depth, confirming our hypothesis that the internal marker covariance provides a valuable source of noise correction.





**Figure S6. Tradict accurately predicts temporal transcriptional responses to lipopolysaccharide treatment in a dendritic cell line CRISPR library.** a) Actual vs. predicted z-score standardized expression of the “response to lipopolysaccharide” transcriptional program. Samples are colored by time point. b) Receiver operator characteristic (ROC) curve illustrating Tradict’s accuracy for identifying differentially expressed (DE) transcriptional programs. Here the “truth set” was considered to be all DE programs with FDR < 0.01 based on actually measured expression values. The marked point along the ROC curve and the inset venn diagram depict the concordance between the predicted and actual set of DE transcriptional programs when an FDR threshold of 0.01 for predicted DE programs was also used. c) Predicted vs actual heatmaps of DE transcriptional programs (rows) across time for different CRISPR lines (columns). Here, DE programs included those found either in actuality or by prediction and are accordingly marked by the black and white indicator bars on the left of each sub-block. Columns of these heat maps represent different profiled lines. The first 12 correspond to negative control guides, whereas the remaining columns correspond to positive regulators of Tnf expression. The expression of programs in each sub-block is z-score normalized to their expression in the negative control guide lines. The bottom 26 programs are all of those directly related to innate immunity among the 368 programs we’ve defined for *M. musculus*. All heatmaps are clustered in the same order across time, genotype, and between predicted and actual.

### Supplemental Note 3 - Tradict accurately predicts temporal dynamics of innate immune signaling in CRISPRed in primary immune cells

To further dissect Tradict’s capabilities, we examined a *M. musculus* dataset from Parnas *et al.* (2015) in which one of the first CRISPR screens was performed on primary immune cells to look for regulators of tumor necrosis factor (Tnf) expression<sup>4</sup>. They found many positive regulators of Tnf expression and created clonal bone-marrow derived dendritic cell (BDMC) lines where each positive regulator was disrupted using CRISPR. They used shallow RNA-sequencing (2.75 +/- 1.2 million reads) to profile the transcriptomes of these lines for 6 hours after lipopolysaccharide (LPS) treatment.

We asked whether Tradict’s predictions could quantitatively recapitulate actuality, despite the challengingly noisy marker measurements due to the low sequencing depth. To be specific, approximately 30% of the markers had zero measured expression in greater than 40% of samples. After performing the batch correction described in Parnas *et al.* (2015), we examined the expression of the “response to lipopolysaccharide” transcriptional program. Figure S6a illustrates that despite the limitation on marker measurement accuracy, Tradict predicts response to LPS with a PCC accuracy of 0.905. Differential transcriptional program expression analysis revealed that DE programs based on Tradict’s predictions were highly concordant with those based on actual measurements (Figure S6b). Strikingly, programs found

DE based on Tradict predictions included 92% of those directly related to innate immune signaling in mice.

We next examined the quantitative quality of Tradict predictions by observing how the DE programs found by either analysis of actual measurements or predictions behave across time. Figure S6c illustrates that despite the high marker measurement error, Tradict's predictions are quantitatively concordant with actuality. As expected most lines of CRISPRed positive regulators demonstrate loss of innate immune signaling.

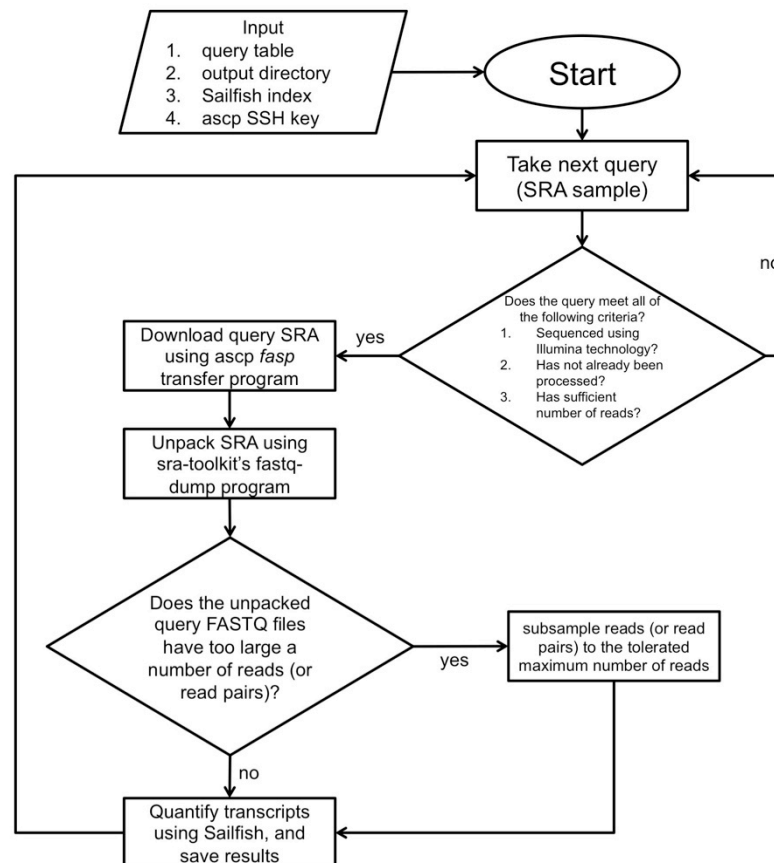
#### **Supplemental Note 4 - Potential improvements to existing targeted RNA sequencing strategies**

We take as an example the workflow of Illumina's TruSeq® Targeted RNA Expression kit, which includes cDNA synthesis, target-probe hybridization, extension-ligation, and PCR with dual indexed primers. By directly probing a total RNA extract with the same probes, but forcing them to hybridize adjacently and using a DNA ligase capable of working on DNA-RNA hybrids<sup>5,6</sup>, it may be possible to skip cDNA synthesis and extension, and proceed directly to ligation and PCR. One could also design more (custom) dual index primers, thereby quadratically increasing the number of multiplexable samples. Given that ~10,000 reads should be sufficient for measuring Tradict's suggested marker panel, one could multiplex ~1,500 samples on a MiSeq (2x39 dual index primers required), or 20,000-30,000 samples on a HiSeq 2500 (2x174 dual index primers required). Coupled with even more reagent-focused improvements (e.g. custom DNA bead purification<sup>7</sup>, custom-ordered enzymes), for example those presented in Yang *et al.* (2016)<sup>8</sup> and Lundberg *et al.* (2013)<sup>9</sup>, the above innovations and Tradict could bring the cost of obtaining actionable transcriptome-wide information close to \$1 per sample.

**Materials and Methods**

**Data acquisition and transcript quantification**

Data acquisition and transcript quantification were managed using a custom script, `srafish.pl`. The `srafish.pl` algorithm and its dependencies are described below. Complete instructions for installing (including all dependencies) and using `srafish.pl` are available on our GitHub page: [https://github.com/surgebiswas/transcriptome\\_compression/tree/master/data\\_download](https://github.com/surgebiswas/transcriptome_compression/tree/master/data_download).



**Figure SM1.** Algorithmic workflow of data acquisition and quantification as implemented by `srafish.pl`.

Figure SM1 illustrates the workflow of `srafish.pl`. Briefly, after checking it meets certain quality requirements, `srafish.pl` uses the `ascp fasp` transfer program to download the raw sequence read archive (.sra file) for an SRA RNA-Seq sample. Transfers made using `ascp` are substantially faster than traditional FTP. The .sra file is then unpacked to FASTQ format using the `fastq-dump` program provided with the SRA Toolkit (NCBI)<sup>10</sup>. The raw FASTQ read data is then passed to `Sailfish`<sup>11</sup>, which uses a fast alignment-free algorithm to quantify transcript abundances. To preserve memory, files with more than 40 million reads for *A. thaliana* and 70 million reads for *M. musculus* are downsampled prior to running `Sailfish`. Samples with fewer than 4 million reads are not downloaded at all. This workflow is then iterated for each SRA RNA-Seq sample available for the organism of interest.



The main inputs into `srafish.pl` are a query table, output directory, Sailfish index, and ascp SSH key, which comes with each download of the aspera ascp client. `srafish.pl` depends on Perl (v5.8.9 for Linux x86-64), the aspera ascp client (v3.5.4 for Linux x86-64), SRA Toolkit (v2.5.0 for CentOS Linux x86-64), and Sailfish (v0.6.3 for Linux x86-64).

### Query table construction

For each organism, using the following (Unix) commands, we first prepared a “query table” that contained all SRA sample ID's as well as various metadata required for the download:

```
qt_name=<query_table_file_name>
sra_url=http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?save=efetch&db=sra&rettype=ru
ninfo&term=
organism=<organism_name>
wget -O $qt_name '$url($organism[Organism]) AND "strategy rna seq"[Properties]'
```

Where fields in between <> indicate input arguments. As an example,

```
qt_name=Athaliana_query_table.csv
sra_url=http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?save=efetch&db=sra&rettype=ru
ninfo&term=
organism="Arabidopsis thaliana"
wget -O $qt_name '$url($organism[Organism]) AND "strategy rna seq"[Properties]'
```

### Reference transcriptomes and index construction

Sailfish requires a reference transcriptome -- a FASTA file of cDNA sequences -- from which it builds an index it can query during transcript quantification. For the *A. thaliana* transcriptome reference we used cDNA sequences of all isoforms from the TAIR10 reference. For the *M. musculus* transcriptome reference we used all protein-coding and long non-coding RNA transcript sequences from the Gencode vM5 reference.

Sailfish indices were created using the following command:

```
sailfish index -t <ref_transcriptome.fasta> -k 20 -p 6 -o .
```

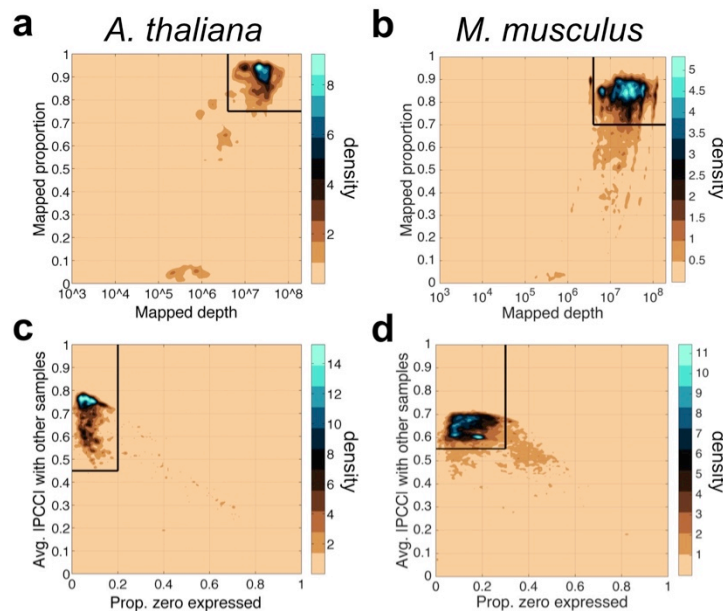
Here, <ref\_transcriptome.fasta> refers to the reference transcriptome FASTA file. Copies of the reference transcriptome FASTA files used in this study are available upon request.

### Quality and expression filtering

Upon completion of download and transcript quantification of all samples, we assembled an n-samples x p-isoforms matrix of transcripts per million (TPM) values as calculated by Sailfish. We then proceeded to quality and expression filter the data as follows:

1. We first removed samples with a read depth and mapping rate below 4 million reads and 0.75 respectively for *A. thaliana* and 4 million reads and 0.70 for *M. musculus* (Figure SM2a-b). We used a slightly lower mapping rate threshold for *M. musculus* because the average mapping rate for *M. musculus* was lower than that of *A. thaliana*. We reasoned this was due to the fact that the Gencode vM5 reference is likely less complete than the TAIR10 reference for *M. musculus*. Though these read count thresholds may be considered slightly lower than what is ideal for both organisms, raising them much higher removed a large number of samples from analysis. Importantly, low read count samples should only add to the noise in the dataset, and so

- the performance results presented in the main text are, if anything, artificially lower than they should be.
- Subsequently, we collapsed the isoform expression table into a gene expression table by setting a gene's expression to be the sum of the expression values of all isoforms of that gene.
  - We next removed all non-protein coding transcripts except for long non-coding RNAs, and removed samples with large amounts (>30%) of non-protein coding contamination (e.g. rRNA).
  - The dataset was then expression filtered by only keeping genes with expression greater than 1 TPM in at least 5% of all samples. The latter requirement ensures that outlier or extreme expression in just a few samples is not enough to keep the gene for analysis.
  - We then removed samples with an abnormally large number of genes with expression values of zero. To do this we calculated the mean and standard deviation of the number of genes with zero expression across all samples. Samples with the number of zero expressed genes greater than the mean plus two times the standard deviation were removed.
  - Finally, we removed outlier samples by first examining the proportion of zeros contained in each sample and by computing the pairwise Pearson correlation coefficient between the gene expression profiles of all samples. To improve heteroskedasticity, raw TPM values for each gene were converted to a log-scale ( $\log_{10}[\text{TPM} + 0.1]$ ) prior to calculating correlations. For *A. thaliana*, the majority of samples had an average correlation with other samples of greater than 0.45 and fewer than 20% percent zero values. Samples with lower correlation or a greater percentage of zeros were removed (Figure SM2c). By similar arguments, samples with less average correlation than 0.55 with other samples and greater than 30% zeros were removed for *M. musculus* (Figure SM2d). Manual inspection of ~100 of these samples revealed they were highly enriched for non-polyA selected samples and samples made from low-input RNA (e.g. single-cells).



**Figure SM2.** Quality filtering thresholds for mapping depth and proportion (a,b), and for average correlation to other samples and proportion of zeros (c,d).

## Metadata annotation

RNA-Seq samples are submitted to the SRA with non-standardized metadata annotations. For example, for some samples tissue and developmental stage are clearly noted as separate fields, whereas in others such information can only be found the associated paper's abstract or sometimes only in its main text. In order to ensure the maximum accuracy when performing metadata annotations, we annotated samples manually until the structure of the gene expression space represented by the first three principal components was clear. Annotation was accomplished by first finding those few submissions with samples in multiple clusters. These submissions revealed that the likely separating variables of interest were issue and developmental context. For each major cluster in the PCA (determined visually) we then annotated samples by size of their submission until the tissue or developmental context of that cluster became qualitatively clear.

## Tradict algorithm

Tradict's usage can be broken down into two parts: 1) Encoding, and 2) Decoding. Encoding is the process of learning, from training data, the marker panel and its predictive relationship to the expression of transcriptional programs and to the remaining genes in the transcriptome. In essence, during encoding we begin with full transcriptome data and collapse its information into a subset of marker genes. Decoding is the reverse process of predicting the expression of transcriptional programs and non-marker genes from the expression measurements of just the selected markers.

Our encoding algorithm can be broken down into several steps: 1) Computing the latent logarithm of the training transcriptome collection, 2) Defining transcriptional programs, 3) marker selection via Simultaneous Orthogonal Matching Pursuit, 4) building a predictive Multivariate Normal Continuous-Poisson hierarchical model.

**1) Computing the latent logarithm of the training transcriptome collection** - Expression values in our training dataset are stored as transcripts per million (TPM), which are non-negative, variably scaled, and strongly heteroscedastic, similar to read counts. For subsequent steps in our algorithm and analysis it will be important transform this data to improve its scaling and heteroscedasticity.

Often, one log transforms such data. However, to avoid undefined values where the data are zeros, one also adds a pseudocount (e.g. 1). This pseudocount considers neither the gene's *a priori* abundance nor the confidence with which the measurement was made, making this practice convenient but statistically unfounded. In previous work, we introduced the latent logarithm, or "lag"<sup>12</sup>. lag assumes that each observed expression value is actually a noisy realization of an unmeasured *latent abundance*. By taking the logarithm of this latent abundance, which considers both sampling depth and the gene's *a priori* abundance, lag provides a more nuanced and statistically principled alternative to the conventional "log(x + pseudocount)". In increasing data, lag quickly converges to log, but in the absence of it, lag relies on both sampling depth and the gene's *a priori* abundance to make a non-zero estimate of the gene's latent abundance.

With these intuitions in mind, we apply the lag transformation to our entire training dataset. The lag-transformed expression matrix demonstrated a Pearson correlation of 0.98 to the log(TPM + 0.1) transformed expression matrix for both *A. thaliana* and *M. musculus*, but again, especially for samples with 0 expression, lag is able to make better estimates of their true abundance in the log-domain.

Availability: [https://github.com/surgebiswas/latent\\_log.git](https://github.com/surgebiswas/latent_log.git)

2) **Defining transcriptional programs** - We define a transcriptional program to be the first principal component of the z-score standardized lag expression of the set of genes involved in a certain response or pathway<sup>13,14</sup>. This virtual program marker maximally captures (in one dimension) the information contained in the transcriptional program. We considered three criteria for defining a globally comprehensive, but interpretable list of transcriptional programs for *A. thaliana* and *M. musculus*:

- a) In order to capture as much information about the transcriptome as possible, we wanted to maximize the number of genes covered by the transcriptional programs.
- b) In order to improve interpretability, we wanted to minimize the total number of transcriptional programs.
- c) The number of genes in a transcriptional program should not be too large or too small -- genes in a transcriptional program should be in the same pathway.

Rather than defining these transcriptional programs *de novo*, we took a knowledge-based approach and defined them using Gene Ontology (GO). We also tried using KEGG pathways, but found these were less complete and nuanced than GO annotations. Gene Ontology is made of three sub-ontologies or aspects: Molecular Function, Biological Process, and Cellular Component. Each of these ontologies contains terms that are arranged as a directed acyclic graph with the above three terms as roots. Terms higher in the graph are less specific than those near the leaves<sup>15,16</sup>. Thus, with respect to the three criteria above, we wanted to find GO terms with low-to-moderate height in the graph such that they were neither too specific nor too general. Given we were interested in monitoring the status of different processes in the organism, we focused on the Biological Process ontology.

We downloaded gene association files for *A. thaliana* and *M. musculus* from the Gene Ontology Consortium (<http://geneontology.org/page/download-annotations>). We then examined for each of several minimum and maximum GO term sizes (defined by the number of genes annotated with that GO term) the number of GO terms that fit this size criterion and the number of genes covered by these GO terms.

Supplemental Tables 1 and 2 contain the results of this analysis for *A. thaliana* and *M. musculus*, respectively. *A. thaliana* has 3333 GO annotations for 27671 genes. We noticed that when the minimum GO term size was as small as it could be (1) and we moved from a maximum GO term size of 5000 to 10000, we jumped from covering 18432 genes (67% of the transcriptome) to covering the full transcriptome (black bolded two rows of Supplemental Table 1). This is due to the addition of one GO term, which was the most general, "Biological Process," term. Thus, we concluded that 33% of the genes in the transcriptome have only "Biological Process" as a GO annotation, and therefore that we do not need to capture these genes in our GO term derived gene sets. We hereafter refer to the set of genes annotated with more than just the "Biological Process" term as *informatively annotated*.

We reasoned that a minimum GO term size of 50 and a maximum size of 2000, best met our aforementioned criteria for defining globally representative GO term derived gene sets. These size thresholds defined 150 GO terms, which in total covered 15124 genes (82.1% of the informatively annotated, and 54.7% of the full transcriptome). These 150 GO-term derived, globally comprehensive transcriptional programs covered the major pathways related to growth, development, and response to the environment.

We performed a similar GO term size analysis for *M. musculus*. *M. musculus* has 10990 GO annotations for 23566 genes. Of these genes, 6832 (29.0%) had only the "Biological Process" term annotation and were considered not informatively annotated. As we did for *A. thaliana*, we selected a GO term size minimum of 50 and a maximum size of 2000. These

size thresholds defined 368 GO terms, which in total covered 14873 genes (88.9% of the informatively annotated, 63% of the full transcriptome). As we found for *A. thaliana*, these 368 GO-term derived, globally comprehensive transcriptional programs covered the major pathways related to growth, development, and response to the environment.

Supplemental Tables 3 and 4 contain the lists of the globally comprehensive transcriptional programs as defined by the criteria above. For each of these programs, we then computed its first principal component over all constituent genes.

**3) Marker selection via Simultaneous Orthogonal Matching Pursuit** - After defining transcriptional programs we are left with a #-training-samples x #-transcriptional-programs table of expression values. We decompose this matrix using an adapted version of the Simultaneous Orthogonal Matching Pursuit, using the #-training-samples x #-transcriptional-programs table as a dictionary<sup>17,18</sup>. Because transcriptional programs are often correlated with other programs, we first cluster them using consensus clustering<sup>19,20</sup>, which produces a robust and stable clustering by taking the consensus of many clusterings performed by a base clustering algorithm. 100 independent iterations of K-means are used as the base-clusterings, and the number of clusters is determined using the Davies-Bouldin criterion<sup>21</sup>. The decomposition is greedy, in which foreach iteration the algorithm first find the transcriptional program cluster with the largest unexplained variance. It then finds the gene contained within this cluster of transcriptional programs with the maximum average absolute correlation to the expression of all transcriptional programs. This gene is then added to an “active set,” onto which the transcriptional program expression matrix is orthogonally projected. This fit is subtracted to produce a residual, on which the above steps are repeated until a predefined number of genes have been added to the active set or the residual variance of the transcriptional program expression matrix falls below some predefined threshold.

#### **4) Building a predictive Multivariate Normal Continuous-Poisson hierarchical model**

Here we describe conceptually how we fit a predictive model that allows us to predict gene and transcriptional program expression from expression measurements of our selected markers. Readers interested in the full mathematical details of the Multivariate Normal Continuous-Poisson hierarchical model are referred to the attached “Tradict - mathematical details” document.

The Multivariate Normal Continuous Poisson distribution offers us a way of modeling statistically coupled count based or, more generally, non-negative random variables, such as the TPM or count-based expression values of genes<sup>22-26</sup>. Here it is assumed the TPM expression of each gene in a given sample is a noisy, Poisson realization of some unmeasured latent abundance, the logarithm of which comes from Multivariate-Normal distribution over the log-latent abundances of all genes in the transcriptome.

Given the marginalization properties of the multivariate normal distribution, we are only interested in learning relationships between the selected markers and non-marker genes. For the purposes of decoding, we need to estimate 1) the mean vector and 2) covariance matrix over the log-latent TPMs of the markers, 3) the mean vector of the log-latent TPMs of the non-markers, and 4) cross-covariance matrix between the log-latent TPMs of markers and non-markers.

Note that before we can estimate these parameters, we must learn the log-latent TPMs of all genes. To do this we first lag-transform the entire training dataset. We then learn the marker log-latent TPMs, and their associated mean vector and covariance matrix using an iterative conditional modes algorithm. Specifically, we initialize our estimate of the marker log-latent TPMs to be the lag-transformed expression values, which by virtue of the lag’s probabilistic assumptions are also derived from a Normal Continuous-Poisson hierarchical



model. We then iterate 1) estimation of the mean vector and the covariance matrix given the current estimate of log-latent TPMs, and 2) maximum *a posteriori* estimation of log-latent TPMs given the estimated mean vector, covariance matrix, and the measured TPM values of the selected markers. A small regularization is added during estimation of the covariance matrix in order to ensure stability and to avoid infinite-data-likelihood singularities that arise from singular covariance matrices. This is most often happens when a gene's TPM abundance is mostly zero (i.e. there is little data for the gene), giving the multivariate normal layer an opportunity to tightly couple this gene's latent abundance to that of another gene, thereby producing a nearly singular covariance matrix.

Learning the mean vector of the non-marker genes and the marker x non-marker cross-covariance matrix is considerably easier. For the mean vector, we simply take the sample mean of the lag-transformed TPM values. For the cross-covariance matrix we compute sample cross-covariance between the learned log-latent marker TPMs and the log-latent non-marker TPMs obtained from the lag transformation. We find that these simple sample estimates are highly stable given that our training collection includes thousands to tens of thousands of transcriptomes.

Using similar ideas, we can also encode the expression of the transcriptional programs. Recall that a principal component output by PCA is a linear combination of input features. Thus by central limit theorem, the expression of these transcriptional programs should behave like normal random variables. Indeed, after regressing out the first 3 principal components computed on the entire training samples x genes expression matrix from the expression values of the transcriptional programs (in order to remove the large effects of tissue and developmental stage), 85-90% of the transcriptional programs had expression that was consistent with a normal distribution (average p-value = 0.43, Pearson's chi-squared test). Consequently, as was done for non-marker genes and as will be needed for decoding, we compute the mean vector of the transcriptional programs and the markers x transcriptional programs cross covariance matrix. These are given by the standard sample mean of the training transcriptional program expression values and sample cross-covariance between the learned log-latent TPMs of the markers and the transcriptional program expression values.

To perform decoding, we must translate newly obtained TPM measurements of our marker genes into expression predictions for transcriptional programs and the remaining non-marker genes. To do this, we first learn the latent abundances of our markers using the measured TPMs, and the 1 x markers mean vector and markers x markers covariance matrix previously learned from the training data. This is done using maximum *a posteriori* estimation. Using these latent abundances and the previously estimated mean vectors and cross-covariance matrices, we can use standard Gaussian conditioning to estimate the log-latent expression of the transcriptional programs and the remaining genes in the transcriptome.

## References

1. Crowley, J. J. *et al.* Analyses of allele-specific gene expression in highly divergent mouse crosses identifies pervasive allelic imbalance. *Nat. Genet.* **47**, (2015).
2. Greenham, K. & McClung, C. R. Integrating circadian dynamics with physiological processes in plants. *Nat Rev Genet* **16**, 598–610 (2015).
3. Donner, Y., Feng, T., Benoist, C. & Koller, D. Imputing gene expression from selectively reduced probe sets. *Nat. Methods* **9**, (2012).
4. Parnas, O., Jovanovic, M., Eisenhaure, M. & Zhang, F. A Genome-wide CRISPR Screen in Primary Immune Cells to Dissect Regulatory Networks. *Cell* **162**, 1–12 (2015).
5. New England BioLabs Inc. SplintR Ligase. at <<https://www.neb.com/products/m0375-splintr-ligase>>

6. Lohman, G. J. S., Zhang, Y., Zhelkovsky, A. M., Cantor, E. J. & Jr, T. C. E. Efficient DNA ligation in DNA – RNA hybrid helices by Chlorella virus DNA ligase. *Nucleic Acids Res.* 1–14 (2013). doi:10.1093/nar/gkt1032
7. Rohland, N. & Reich, D. Cost-effective, high-throughput DNA sequencing libraries for multiplexed target capture. *Genome Res.* **22**, 939–946 (2012).
8. Yang, L. *et al.* The *Pseudomonas syringae* type III effector HopBB1 fine tunes pathogen virulence by gluing together host transcriptional regulators for degradation. *Submitted* (2016).
9. Lundberg, D. S., Yourstone, S., Mieczkowski, P., Jones, C. D. & Dangl, J. L. Practical innovations for high-throughput amplicon sequencing. *Nat. Methods* **10**, 999–1002 (2013).
10. Leinonen, R., Sugawara, H. & Shumway, M. The Sequence Read Archive. **39**, 2010–2012 (2011).
11. Patro, R., Mount, S. M. & Kingsford, C. Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nat. Biotechnol.* **32**, 462–4 (2014).
12. Biswas, S. The latent logarithm. *arXiv* 1–11 (2016).
13. Ma, S. & Kosorok, M. R. Identification of differential gene pathways with principal component analysis. *Bioinformatics* **25**, 882–889 (2009).
14. Fan, J. *et al.* Characterizing transcriptional heterogeneity through pathway and gene set overdispersion analysis. *Nat. Methods* **13**, 241–244 (2016).
15. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nat Genet* **25**, 25–29 (2000).
16. The Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res.* **43**, D1049–D1056 (2015).
17. Tropp, J. a & Gilbert, A. C. Signal Recovery From Random Measurements Via Orthogonal Matching Pursuit. *IEEE Trans. Inf. Theory* **53**, 4655–4666 (2007).
18. Tropp, J. a., Gilbert, A. C. & Strauss, M. J. Algorithms for simultaneous sparse approximation. Part I: Greedy pursuit. *Signal Processing* **86**, 572–588 (2006).
19. Monti, S., Tamayo, P., Mesirov, J. & Golub, T. Consensus Clustering : A Resampling-Based Method for Class Discovery and Visualization of Gene Expression Microarray Data. *Mach. Learn.* **52**, 91–118 (2003).
20. Yu, Z., Wong, H.-S. & Wang, H. Graph-based consensus clustering for class discovery from gene expression data. *Bioinforma.* **23** , 2888–2896 (2007).
21. Davies, D. L. & Bouldin, D. W. A Cluster Separation Measure. *IEEE Trans. Pattern Anal. Mach. Intell.* **2**, 224–227 (1979).
22. Aitchison, J. & Shen, S. M. Logistic-Normal Distributions: Some Properties and Uses. *Biometrika* **67**, 261 (1980).
23. Aitchison, J. & Ho, C. H. The multivariate Poisson-log normal distribution. *Biometrika* **76**, 643–653 (1989).
24. Biswas, S., Mcdonald, M., Lundberg, D. S., Dangl, J. L. & Jojic, V. Learning Microbial Interaction Networks from Metagenomic Count Data. in *Res. Comput. Mol. Biol.* **1**, 32–43 (2015).
25. Ho, C. H. & Kong, H. The multivariate Poisson-log normal distribution. **2**, (1989).
26. Madsen, L. & Dalthorp, D. Simulating correlated count data. *Environ. Ecol. Stat.* **14**, 129–148 (2007).

# Tradict - mathematical details

Surojit Biswas, Konstantin Kerner, Paulo José Pereira Lima Texeira,  
Jeffery L. Dangl, Vladimir Jojic, Philip A. Wigge

## 3 Contents

4	1 Preliminaries	1
5	2 Model	1
6	3 Encoding	2
7	3.1 Inference of $z_m$ given $\mu^{(m)}$ and $\Sigma^{(m)}$	3
8	3.2 Complete inference of $\mu^{(m)}$ , $\Sigma^{(m)}$ , and $z_m$	4
9	4 Decoding	5
10	5 References	6

This document describes the full mathematical details for the concepts presented in the “Tradict algorithm” section, “Building a predictive Multivariate Normal Continuous-Poisson hierarchical model” subsection of the Materials and Methods in the Supplemental Information. Specifically, we present exactly how Tradict uses a selected set of markers to 1) complete the encoding, and 2) to perform decoding.

## 15 1 Preliminaries

For a matrix  $A$ ,  $A_{:i}$  and  $A_{i:}$  index the  $i^{th}$  column and row, respectively. For a set of indices,  $q$ , we use  $-q$  to refer to all indices not specified by  $q$ .

## 18 2 Model

Tradict uses a Continuous-Poisson Multivariate Normal (CP-MVN) hierarchical model to model the expression of transcriptional programs and all genes in the transcriptome. Multivariate Normal hierarchies have been explored in the past as a means of modeling correlation structure among count based random variables [REF]. However, given we will be working with abundances as transcripts per million (TPM), which are non-negative (can equal zero) and fractional, we relax the integral assumption of the Poisson so it is continuous on  $[0, \infty)$ . Specifically, we define the continuous relaxation of the Poisson distribution (hereafter, Continuous-Poisson) to have the following density function:

$$f(x|\lambda) = C_\lambda \frac{e^{-\lambda} \lambda^x}{\Gamma(x+1)}$$

where  $C_\lambda$  is a normalization constant [REF]. The mean of this distribution is given by  $\lambda$ , just as the Poisson. We begin by building a predictive model of gene expression, and thereafter discuss a predictive model for the expression of transcriptional programs. Let  $z_j$  denote the log-latent abundance of gene  $j$ , such that

29  $\exp(z_j)$  is the *latent abundance* of that gene (in TPM) whose measured abundance is given by  $t_j$ . Let  
 30  $T_j = t_j o$  be the measured total number of transcripts of gene  $j$ . Here  $o$  is the sequencing depth in millions  
 31 of reads of the sample under consideration. We assume then,

$$z \sim \mathcal{N}(\mu, \Sigma)$$

$$T_j \sim \text{Continuous-Poisson}(\exp(z_j)o)$$

32 where  $\mu$  and  $\Sigma$  are of dimension  $1 \times \#\text{-genes}$  and  $\#\text{-genes} \times \#\text{-genes}$ , respectively. In effect, we are assuming  
 33 that the measured number of transcripts for gene  $j$  is a noisy realization of a latent abundance  $\exp(z_j)$  times  
 34 the sequencing depth,  $o$ . The dependencies between log-latent abundances (the  $z_j$ 's) are then encoded by  
 35 the covariance matrix of the Multivariate Normal layer of the model.

36 Note that we could model the TPM measurements directly in the second layer by assuming  $t_j \sim$   
 37  $\text{Continuous-Poisson}(\exp(z_j))$ ; however, this formulation does not consider sequencing depth, which can be a  
 38 valuable source of information when inferring latent abundances for rare/poorly sampled genes [1].

39 During decoding, we are interested in building a predictive model between markers and all genes in the  
 40 transcriptome. Therefore, we need to consider a conditional model of the transcriptome given the log-latent  
 41 abundances of the markers. Let  $m$  be the set of indices for the given panel of selected markers, which are the  
 42 subset of genes Tradict selects as representative of the transcriptome. To perform prediction we therefore  
 43 need  $p(z_{-m}|z_m)$ , and given this we would like to ultimately compute as our estimate of the abundance of all  
 44 genes in the transcriptome  $\hat{T} = \text{argmax}_T p(T|z_m)$ . We have,

$$z_m \sim \mathcal{N}(\mu^{(m)}, \Sigma^{(m)})$$

$$z_{-m}|z_m \sim \mathcal{N}(\mu_{z_{-m}|z_m}, \Sigma_{z_{-m}|z_m})$$

$$T_j \sim \text{Continuous-Poisson}(\exp(z_j)o)$$

45 Here,  $\mu^{(m)}$  and  $\Sigma^{(m)}$  refer to mean vector and covariance matrix of  $z_m$ . Given these, the conditional  
 46 mean of the log-latent abundances for all non-marker genes can be obtained through Gaussian conditioning.  
 47 Specifically, for two normally distributed row-vector variables  $a$  and  $b$  the conditional mean of  $b$  given  $a$  is  
 48 given by  $\mu_{b|a} = \mu_b + (a - \mu_a)\Sigma_a^{-1}\sigma_{ab}$  and  $\Sigma_{b|a} = \Sigma_b - \sigma_{ab}^T\Sigma_a^{-1}\sigma_{ab}$ , where  $\sigma_{ab}$  is the cross-covariance between  
 49  $a$  and  $b$ , and  $\Sigma_a$  and  $\Sigma_b$  are the covariance matrices of  $a$  and  $b$ , respectively.

50 Given the expression of a transcriptional program is a linear combination of the latent abundances  
 51 of its constituent genes, they will be normally distributed given 1) Central Limit Theorem, and 2) the  
 52 latent abundances themselves are normally distributed (convolutions of normals are normals). Let  $s$  be the  
 53 expression of all transcriptional programs. We posit the following model,

$$z_m \sim \mathcal{N}(\mu^{(m)}, \Sigma^{(m)})$$

$$s|z_m \sim \mathcal{N}(\mu_{s|z_m}, \Sigma_{s|z_m})$$

54 To use these models for prediction, we must learn their parameters from training data. This would complete  
 55 the process of encoding described in the Supplemental Information. Specifically, we need to learn  $\mu^{(m)}$ ,  $\Sigma^{(m)}$ ,  
 56  $\mu_s$ ,  $\mu_{z_{-m}}$ ,  $\sigma_{z_m, s}$  and  $\sigma_{z_m, z_{-m}}$ .

### 57 3 Encoding

58 As described in the Supplemental Information, given an estimate of  $z_m$ ,  $\hat{z}_m$ , inference of  $\mu_s$ ,  $\mu_{z_{-m}}$ ,  $\sigma_{z_m, s}$  and  
 59  $\sigma_{z_m, z_{-m}}$  is straightforward. In lag transforming the entire training TPM expression matrix,  $t \in \mathbb{R}^{\text{samples} \times \text{genes}}$ ,  
 60 we have an estimate of  $z$ ,  $\hat{z} = \text{lag}(t)$  [1]. Thus, an estimate of  $\mu_{z_{-m}}$  is given by the usual column-wise sample  
 61 mean of  $\hat{z}_{-m}$ .

62 Let  $\Lambda \in \mathbb{R}^{\text{genes} \times \text{transcriptional programs}}$  be a matrix of principal component 1 coefficients over genes for each  
 63 transcriptional program. Note, that  $\Lambda_{ij} = 0$  if gene  $i$  is not in transcriptional program  $j$ . An estimate of  $s$   
 64 is given by  $\hat{s} = \hat{z}\Lambda$ , and so an estimate for  $\mu_s$ ,  $\hat{\mu}_s$ , is given by the usual column-wise mean of  $\hat{s}$ .

65 Given  $\hat{z}_m$  the cross-covariances,  $\sigma_{z_m, s}$  and  $\sigma_{z_m, z_{-m}}$ , are given by the usual sample cross-covariance between  
 66  $\hat{z}_m$  and  $\hat{s}$  and between  $\hat{z}_m$  and  $\hat{z}_{-m}$ , respectively.

67 Now, though we could use the lag-transformed values of  $t_m$  as our estimate for  $z_m$ , we have an opportunity  
 68 to improve this estimate by virtue of having to estimate  $\mu^{(m)}$  and  $\Sigma^{(m)}$ . More specifically, given  $z_m$ , estimates  
 69 of  $\mu^{(m)}$  and  $\Sigma^{(m)}$  are given by – up to some regularization – the usual sample mean and covariance of  $z_m$ .  
 70 Furthermore, given  $\mu^{(m)}$  and  $\Sigma^{(m)}$ , we can update our estimate of  $z_m$  to the maximum of its posterior  
 71 distribution. This suggests an alternating iterative procedure in which we iterate 1) estimation of  $\mu^{(m)}$  and  
 72  $\Sigma^{(m)}$ , and 2) maximum *a posteriori* inference of  $z_m$  until convergence of their joint likelihood. It is the  $\hat{z}_m$   
 73 that we obtain from this procedure that we use in the cross-covariance calculations above. The following  
 74 section details this procedure.

### 75 3.1 Inference of $z_m$ given $\mu^{(m)}$ and $\Sigma^{(m)}$

76 Suppose Tradict has estimates of  $\mu^{(m)}$  and  $\Sigma^{(m)}$  given by  $\hat{\mu}^{(m)}$  and  $\hat{\Sigma}^{(m)}$ , and let  $T_m = t_m(o \times \mathbf{1}_{1 \times \text{markers}})$   
 77 be a matrix of the total measured number of transcripts for each marker. Here  $o \in \mathbb{R}^{\text{samples} \times 1}$  is a vector  
 78 of sample sequencing depths in millions of reads. Given these, we would like to calculate the maximum *a*  
 79 *posteriori* (MAP) estimate of  $\hat{z}_m = \text{argmax}_{z_m} p(z_m | o, T_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$ .

80 The posterior distribution over  $z_m$  is given by

$$\begin{aligned} p(z_m | o, T_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) &= \frac{p(T_m | o, z_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(z_m | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})}{\int_k p(T_m | o, k, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(k | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) dk} \\ &\propto \prod_{i=1}^n p(T_{im} | o, z_{im}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(z_{im} | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\ &= \prod_{i=1}^n \left[ \prod_{j=1}^{|m|} C_{[\exp(z_{ij}) o_i]} [\exp(z_{ij}) o_i]^{T_{ij}} e^{-[\exp(z_{ij}) o_i]} / \Gamma(T_{ij} + 1) \right] \\ &\quad \times \frac{1}{\sqrt{2\pi |\hat{\Sigma}^{(m)}|}^{|m|}} \exp \left( -\frac{1}{2} (z_{i:} - \hat{\mu}^{(m)}) \text{inv} \left( \hat{\Sigma}^{(m)} \right) (z_{i:} - \hat{\mu}^{(m)})^T \right) \end{aligned}$$

81 where for notational clarity we have used  $\text{inv}(\cdot)$  to represent matrix inverse.

82 Given  $z$  is a matrix parameter, this may be difficult to solve directly. However, note that given  $z_{ij}$ ,  $T_{ij}$   
 83 is conditionally independent of  $T_{i,-j}$ . Additionally, given  $z_{i,-j}$ ,  $z_{ij}$  is normally distributed with mean and  
 84 covariance

$$\begin{aligned} a_{ij} &= \mu_j^{(m)} + \left( z_{i,-j} - \mu_{-j}^{(m)} \right) \text{inv} \left( \Sigma_{-j,-j}^{(m)} \right) \Sigma_{-j,j}^{(m)} \\ \sigma_{m(j)} &= \Sigma_{j,j}^{(m)} - \Sigma_{j,-j}^{(m)} \text{inv} \left( \Sigma_{-j,-j}^{(m)} \right) \Sigma_{-j,j}^{(m)} \end{aligned}$$

85 respectively. Taken together, this suggests an iterative conditional modes algorithm [2] in which we maximize  
 86 the posterior one column of  $z$  at a time, while conditioning on all others.

87 Let  $\hat{z}_m$  denote our current estimate of  $z_m$ . Let  $m(j)$  denote the index of the  $j^{\text{th}}$  marker and let  $m(-j)$



88 denote the indices of all markers but the  $j^{th}$  one. The above sub-objective is given by,

$$\begin{aligned}
\hat{z}_{im(j)} &= \operatorname{argmax}_{z_{im(j)} | z_{im(-j)}} \log p(z_{im(j)} | T_{im(j)}, o_i, \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\
&= \operatorname{argmax}_{z_{im(j)} | z_{im(-j)}} \log p(T_{im(j)} | z_{im(j)}, o_i, \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(z_{im(j)} | \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\
&= \operatorname{argmax}_{z_{im(j)} | z_{im(-j)}} \log p(T_{im(j)} | z_{im(j)}, o_i) p(z_{im(j)} | \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\
&= \operatorname{argmax}_{z_{im(j)} | z_{im(-j)}} \log \left[ [\exp(z_{im(j)}) o_i]^{T_{im(j)}} e^{-[\exp(z_{im(j)}) o_i]} \exp \left( -\frac{1}{2\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})^2 \right) \right] \\
&= \operatorname{argmax}_{z_{im(j)} | z_{im(-j)}} T_{im(j)} \exp(z_{im(j)}) o_i - \exp(z_{im(j)}) o_i - \frac{1}{2\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})^2
\end{aligned}$$

89 Differentiating we get,

$$\begin{aligned}
\frac{\partial}{\partial z_{im(j)}} T_{im(j)} z_{im(j)} o_i - \exp(z_{im(j)}) o_i - \frac{1}{2\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})^2 \\
= T_{im(j)} o_i - \exp(z_{im(j)}) o_i - \frac{1}{\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})
\end{aligned}$$

90 Because  $z_{im(j)}$  appears as a linear and exponential term, we cannot solve this gradient analytically. We  
91 therefore utilize Newton-Raphson optimization. For this we also require the Hessian, which is given by,

$$\begin{aligned}
\frac{\partial}{\partial z_{im(j)}} T_{im(j)} o_i - \exp(z_{im(j)}) o_i - \frac{1}{\sigma_{m(j)}} (z_{im(j)} - a_{im(j)}) \\
= -\exp(z_{im(j)}) o_i - \frac{1}{\sigma_{m(j)}} < 0
\end{aligned}$$

92 Notice the Hessian is always negative-definite, which implies each update has a single, unique optimum.

93 In practice, the Newton-Raphson updates can be performed in vectorized fashion iteratively for each  
94 column of  $z$ . We generally find that this optimization takes 5-15 iterations (full passes over all columns  
95 of  $z$ ) and less than a minute to converge. We refer to the program that performs these calculations as  
96  $\hat{z}_m = \text{MAP\_Z}(t, o, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$ .

## 97 3.2 Complete inference of $\mu^{(m)}$ , $\Sigma^{(m)}$ , and $z_m$

98 For complete inference we use the following iterative conditional modes algorithm [2]:

- 99 • Initialize  $T_m = t_m(o \times \mathbf{1}_{1 \times \text{markers}})$ ,  $\hat{z}_m = \text{lag}(t_m)$ .  
100 • Until convergence of  $\log p(T_m | o, \hat{z}_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) + \log p(\hat{z}_m | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$ , iterate:  
101 – Update  $\hat{\mu}^{(m)}$  and  $\hat{\Sigma}^{(m)}$ :

$$\begin{aligned}
\hat{\mu}^{(m)} &= \frac{1}{\# \text{samples}} \sum_i \hat{z}_{im} \\
\hat{\Sigma}^{(m)} &= \frac{1}{\# \text{samples} - 1} \sum_i (\hat{z}_{im} - \hat{\mu}^{(m)})^T (\hat{z}_{im} - \hat{\mu}^{(m)}) + \lambda \text{diag} \left[ \text{cov} \left( \hat{z}_m^{(\text{init})} \right) \right]
\end{aligned}$$

- 102 – Update  $\hat{z}_m = \text{MAP\_Z}(t, o, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$ .

Here  $\text{diag}(x)$  of the square matrix  $x$  returns an equivalently sized matrix with only the diagonal of  $x$  preserved and 0's for the off-diagonal terms.  $\text{cov}(\cdot)$  denotes the usual sample covariance matrix.

Note that in this algorithm we have added a regularization to the estimate of the covariance matrix. This is done in order to ensure stability and to avoid infinite-data-likelihood singularities that arise from singular covariance matrices. This is most often happens when a genes TPM abundance is mostly zero (i.e. there is little data for the gene), giving the multivariate normal layer an opportunity to increase the data likelihood (via the determinant of the covariance matrix) by tightly coupling this genes latent abundance to that of another gene, thereby producing a singularity. This regularization is probabilistically equivalent to adding an Inverse-Wishart prior over  $\Sigma^{(m)}$ . The parameter  $\lambda$  controls the strength of the regularization. In practice, we find  $\lambda = 0.1$  leads to good predictive performance, stable (non-singular) covariance matrices, and reasonably quick convergence.

## 4 Decoding

During decoding we are given new measured TPM measurements for our markers,  $t_m^* \in \mathbb{R}^{\text{query samples} \times |m|}$ , and we must make predictions about the expression of all transcriptional programs and the remaining non-marker genes. To do this we first need an estimate of the log-latent abundances  $\hat{z}_m^*$  associated with  $t_m^*$ . Given the estimates  $\hat{\mu}^{(m)}$  and  $\hat{\Sigma}^{(m)}$  obtained from the training data, we obtain these estimates as

$$\hat{z}_m^* = \text{MAP\_Z} \left( t_m^*, \mathbf{1}_{\text{query samples} \times 1}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)} \right)$$

Given the inferred marker latent abundances, we let our estimates of  $s^*$  and  $t_m^*$  be the maximizers of their probability distribution. In other words,  $\hat{s}^* = \text{argmax}_{s^*} p(s^* | \hat{z}_m^*)$  and  $\hat{t}_m^* = \text{argmax}_{t_m^*} p(t_m^* | \hat{z}_m^*)$ .

Our estimate for the expression of all transcriptional programs is given by

$$\text{argmax}_{s^*} p(s^* | \hat{z}_m^*) = \mathbb{E}[s^* | \hat{z}_m^*] = \mu_{s^* | \hat{z}_m^*} = \hat{\mu}_s + \left( \hat{z}_m^* - \hat{\mu}^{(m)} \right) \text{inv} \left( \hat{\Sigma}^{(m)} \right) \hat{\sigma}_{z_m, s}.$$

Here,  $\hat{\mu}_s$  and  $\hat{\sigma}_{z_m, s}$  represent estimates of the unconditional mean of  $s$  and the cross-covariance matrix between  $z_m$  and  $s$  previously learned during encoding.

Similarly, for the entire transcriptome we have,

$$\hat{t}_{ij}^* = \text{argmax}_t p(t | \hat{z}_{im}^*) = \exp \left( \mu_{z_{ij} | \hat{z}_{im}^*} \right).$$

where,

$$\mu_{z_{ij} | \hat{z}_{im}^*} = \hat{\mu}_j + \left( \hat{z}_{im}^* - \hat{\mu}^{(m)} \right) \text{inv} \left( \hat{\Sigma}^{(m)} \right) \hat{\sigma}_{z_m, z_j}$$

We could also use the expected value of  $t$  as our estimate.

$$\begin{aligned} \mathbb{E}[t_{ij}^* | \hat{z}_{im}^*] &= \int_{-\infty}^{\infty} \mathbb{E}[t_{ij}^* | z_{ij}] p(z_{ij} | \hat{z}_{im}^*) dz_{ij} \\ &= \int_{-\infty}^{\infty} \exp(z_{ij}) \mathcal{N}(z_{ij} | \mu_{z_{ij} | \hat{z}_{im}^*}, \Sigma_{z_{ij} | \hat{z}_{im}^*}) dz_{ij} \\ &= \mathbb{E}_{\mathcal{N}}[\exp(z_{ij}) | \hat{z}_{im}^*] \end{aligned}$$

The Moment Generating Function of a Normal random variable  $X$  with mean  $\mu$  and variance  $\sigma^2$  is given by  $M(t) = \mathbb{E}[\exp(tX)] = \exp(\mu t + \sigma^2 t^2 / 2)$ . Therefore we have,

$$\mathbb{E}[t_{ij}^* | \hat{z}_{im}^*] = \mathbb{E}_{\mathcal{N}}[\exp(z_{ij}) | \hat{z}_{im}^*] = M(1) = \exp \left( \mu_{z_{ij} | \hat{z}_{im}^*} + \frac{1}{2} \Sigma_{z_{ij} | \hat{z}_{im}^*} \right)$$

129 where,

$$\begin{aligned}\mu_{z_{ij}|\hat{z}_{im}^*} &= \hat{\mu}_j + \left(\hat{z}_{im}^* - \hat{\mu}^{(m)}\right) \text{inv} \left(\hat{\Sigma}^{(m)}\right) \hat{\sigma}_{z_m, z_j} \\ \Sigma_{z_{ij}|\hat{z}_{im}^*} &= \hat{\sigma}_{jj} - \hat{\sigma}_{z_m, z_j}^T \text{inv} \left(\hat{\Sigma}^{(m)}\right) \hat{\sigma}_{z_m, z_j}\end{aligned}$$

130 Here,  $\hat{\mu}_j$  and  $\hat{\sigma}_{z_m, z_j}$  represent estimates of the unconditional mean of  $z_j$  and the cross-covariance matrix  
131 between  $z_m$  and  $z_j$ . These were learned from the training data during encoding.

132 Though this predictor is unbiased, it does not produce a good prediction for most samples. This is due  
133 to the right-skew of the Poisson, which drags its mean away from the most likely values.

## 134 5 References

135 [1] Surojit Biswas. The latent logarithm. *arXiv*, pages 1–11, 2016.

136 [2] Julian Besag. On the Statistical Analysis of Dirty Pictures. *Journal of the Royal Statistical Society*,  
137 48(3):259–302, 1986.