

1 **Supplementary Methods: Effects of transcriptional noise on estimates of gene and transcript**
2 **expression in RNA sequencing experiments**

3 Ales Varabyou ^{1,3*}, Steven L. Salzberg ^{1,2,3,4}, and Mihaela Pertea ^{1,2,3*}

4 1 Center for Computational Biology, Johns Hopkins University, Baltimore, MD

5 2 Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD

6 3 Department of Computer Science, Johns Hopkins University, Baltimore, MD

7 4 Department of Biostatistics, Johns Hopkins University, Baltimore, MD

8 *Corresponding author. Email: ales.varabyou@jhu.edu, mperte@jhu.edu

9

1 **Methods**

2 **Reference data**

3 To avoid unnecessary complexity, we restricted our analysis to the primary chromosomes of the
4 GRCh38.12 human assembly (Graves-Lindsay et al. 2017), excluding all alternative scaffolds and
5 patches. As reference annotation we used the version of CHES 2.2 human gene catalog tailored for
6 transcriptome assembly, which does not include pseudogenes, V|J|D|C segments, snoRNAs,
7 snRNAs, telomerase RNAs, guide RNAs, or other small RNAs.

8 **Command:** [build_indices.py](#) --idx_dir <output_directory> --reference <reference_genome.fasta> --
9 annotation <reference_annotation.gff> --hisat <path_to_hisat2_executable> --salmon
10 <path_to_salmon_executable> --kallisto <path_to_kallisto_executable> --gffread
11 <path_to_gffread_executable>

12 **Identification and typing of noisy transcription**

13 To identify other types of transcription observed across the GTEx dataset (Lonsdale et al. 2013), we
14 relied on the set of ~30 million transcripts assembled as part of building the CHES 2.2 gene database
15 (Pertea et al. 2018). These 30 million transcripts were the unfiltered results of assembling the GTEx
16 data.

17 **Filtering and typing**

18 We first used gffcompare (Pertea and Pertea 2020) to identify overlaps between all ~30 million
19 assembled transcripts and known features in the CHES 2.2 annotation. We then used gffread
20 (Pertea and Pertea 2020) to extract the loci of all intron-compatible transcripts that overlapped any of
21 the annotated genes. Next, to reduce the number of confounding factors in our analysis, we removed
22 any noisy isoforms that overlapped annotated loci on the opposite strand, contained annotated loci
23 within their introns, or were in close proximity to known genes but did not overlap their exons.
24 Together, our filtering criteria reduced the initial set of transcripts from ~30 million down to ~20 million
25 assembled isoforms.

26 Transcripts that did not represent annotated isoforms in CHES 2.2 were further classified into:

- 1 1. Intronic noise: any assembled feature entirely contained within an intronic region of a gene,
2 including both coding and noncoding genes.
- 3 2. Splicing noise: transcripts that shared a splice site with any of the annotated transcripts in
4 CHES. These could either be un-annotated but nonetheless functional transcripts, or
5 transcripts that occur due to errors produced by the splicing machinery, incompletely
6 processed pre-mRNA, or other biological phenomena not accounted for by the annotation.
- 7 3. Intergenic noise: the set of all transcripts that do not overlap any of the annotated transcripts.

8 Multiple transcripts in the assembled dataset did not have a strand assigned to them. For consistency
9 in the downstream analysis, we randomly assigned a strand to all intergenic transcripts, making sure
10 that all transcripts in the same locus share the same strand. For all other transcripts, we assigned the
11 strand of the annotated gene containing or overlapping the transcript.

12 **Command:** [setup.py](#) --base_dir_data <directory_with_input_data> --base_dir_out <output_directory>
13 --gtex_stats <[gtex_stats](#)> --filters <[filters.sh](#)> --seed 101

14 **Extracting parameters of the simulation from GTEx**

15 Expression estimates and frequency of occurrence of each transcript and locus in the GTEx tissues
16 and samples were computed based on the tracking information which gffcompare builds between
17 merged transcripts and constituent assemblies. These summaries were computed for each type of
18 real and noise transcripts and loci.

19 **Command:** [gtex_stats](#) -t <path_to_tissues_list> -a ALL.merged.filtered.tracking -r real.gtf -g
20 ALL.merged.filtered.gtf -s splicing.gtf -i intronic.gtf -p RNAPol.gtf -o <output_basename>

21 **Simulation**

22 **Selecting the loci to be expressed in a simulated tissue**

23 We assumed a Gaussian distribution of the number of real and intergenic loci per each tissue, where
24 we determined the mean and the standard deviation based on the observations extracted from the
25 GTEx analysis described in section 2.2. To generate the number of loci expressed in a tissue, we
26 sampled from these two distributions independently. Furthermore, we assumed that if a locus is

1 chosen to be expressed in a tissue, all transcripts at that locus can also be expressed in the given
2 tissue.

3 **Command:** [step1_sim_tissues.py](#) --out_dir <output_directory> --base_dir_data
4 <input_directory_from_setup.py> --base_dir_out <output_directory_from_setup.py> --num_tissues 3 -
5 -seed 101

6 **Selecting expression levels for the transcripts in a simulated tissue**

7 During our analysis of the GTEx dataset, we noticed that the relationship between expression of real
8 and noisy transcripts in a given locus is non-linear and strongly correlated, and annotated transcripts
9 tend to dominate the expression of a locus (Fig. 1D). We also observed that the expression values of
10 a transcript within a tissue are more correlated than the expression values of a transcript between
11 different tissues. To preserve these and any other possible relationships existing in the data, we
12 sampled the number of transcripts per each locus as well as corresponding expression values as they
13 appeared in GTEx. To randomize our sampling algorithm without changing any potential relationships
14 we relied on the following procedure:

- 15 1. For each locus selected for a simulated tissue, we identified all loci with identical numbers of
16 real, splicing noise and intronic noise isoforms.
- 17 2. We then randomly selected one of the matching loci and obtained a list of expression values
18 as computed in section 2.2. Each such list contained expression values of each transcript from
19 the selected locus observed in a single tissue, with absent transcripts marked by a zero.
20 Values within the list were grouped by sample and each group was sorted by unique transcript
21 identifier, thus guaranteeing identical ordering between groups. Each group in the list
22 represented a possible set of expressions for a single simulated sample.
- 23 3. Keeping the groups unchanged, we shuffled a list of transcript identifiers for the selected locus
24 to further randomize associations between isoforms and expression values, without changing
25 any relationships in expression between values. This step produced a matrix in which rows
26 denote individual transcripts of a locus, columns denote possible samples and values contain

1 expression values of transcripts in a sample, with zeros indicating that a transcript is not
2 expressed in a sample.

3 The current approach is limited by using observed values only. A possible addition to the method
4 would be to replace observed values with a multivariate distribution of expression values of each type,
5 which can be estimated from the observed values. This would allow creating more random samples.
6 However, further work needs to be done to determine the types of distributions suitable for modeling
7 expression values.

8 **Command:** [step2_sim_attach_exp.py](#) --out_dir <output_directory_from_sim_tissues.py> --
9 base_dir_data <input_directory_from_setup.py> --base_dir_out <output_directory_from_setup.py> --
10 num_tissues 3 --seed 101

11 **Selecting transcripts for the simulated samples**

12 Given a set of loci to be expressed in a simulated tissue and a set of possible expression values for
13 each real and noisy isoform at each locus, we randomly drew individual samples. Not all tissue loci
14 are typically expressed in a single sample of that tissue. To simulate this process, we took a random
15 subset of the tissue loci based on a random sampling from the distribution of the number of loci per
16 sample as observed in GTEx.

17 Most of the tissues in GTEx contain many more samples than are simulated in our analysis and not all
18 loci that are detected in a tissue are expressed in all samples of that tissue, which may be caused by
19 multiple intrinsic and extrinsic factors contributing to the stochasticity of gene expression in living
20 organisms (Swain et al. 2002). As such, simulating a smaller number of samples per tissue often
21 results, as evident from our simulation, in a lower total number of loci within simulated tissues than
22 observed in GTEx.

23 Next, for each locus we selected a random column from the simulated matrix of possible expression
24 values as described in step 3 from section 3.2. The number of transcripts expressed for a given locus
25 in a sample is then dictated by the number of non-zero values in the sampled column.

1 At the end of this step the simulation produces 4 groups of files for each type of transcription being
2 simulated. Each group includes:

- 3 1. an annotation in a GTF format, listing transcript and exon features only.
- 4 2. a list of expression values in TPM format.

5 **Command:** [step3_sim_samples.py](#) --out_dir <output_directory_from_sim_tissues.py> --
6 base_dir_data <input_directory_from_setup.py> --base_dir_out <output_directory_from_setup.py> --
7 num_tissues 3 --num_samples 10 --seed 101

8 **Normalization of transcript abundances in the simulated samples**

9 In the next step we randomly sample the number of reads to be simulated for a given sample from the
10 observations in GTEx. To compute the target number of reads to be simulated for each transcript in a
11 sample given the TPM of each molecule and the total number of reads, we reversed the TPM
12 calculation as follows:

$$13 \quad C_i = \left(\frac{E_i \times L_i}{\sum_{j=0}^N (E_j \times L_j)} \times N \times l \right) \div L_i$$

14 where C_i is the coverage of transcript i , E_i is its expression measured in TPM, L_i is length, and N is the
15 number of reads in the sample and l is the read length (Li and Dewey 2011).

16 **Command:** [step4_sim_normalize.py](#) --out_dir <output_directory_from_sim_tissues.py> --
17 base_dir_data <input_directory_from_setup.py> --base_dir_out <output_directory_from_setup.py> --
18 num_tissues 3 --num_samples 10 --seed 101 --readlen 101

19 **Visualization of summarized statistics for the GTEx and simulated datasets**

20 All scripts for re-producing visualizations from our analysis of GTEx and simulated datasets are
21 provided in a public Python **notebook** at: [evaluation_stage1_results.ipynb](#)

1 **Simulating reads**

2 We used the simulated transcripts selected for each sample and their corresponding coverages to
3 simulate reads for each experiment. Reads were simulated using the Polyester package (Frazee et al.
4 2015), and the input FASTA file was prepared using gffread.

5 In our preliminary analysis we found that Polyester is unable to accurately model the sequencing
6 protocol when simulating paired-end reads. Specifically, for transcripts shorter than the fragment
7 length, Polyester left gaps in coverage unless the fragment length was reduced, or the standard
8 deviation was set high enough to cause frequent overlaps between pairs. Setting fragment length and
9 standard deviation according to the minimum transcript length instead of realistic observations from
10 sequencing data would inaccurately represent sequencing protocols. Therefore, to avoid introducing
11 artificial coverage gaps we chose to simulate single-end reads instead.

12 We also found that when simulating single-end reads Polyester was often unable to extend read
13 coverage to the end of the last exon in the transcript. We assume that internally Polyester always
14 simulates paired-end reads and only reports the first mate in the single-end mode. We found that
15 setting the fragment length to be the same as the read length and standard deviation to zero fixed the
16 issue.

17 Each sample was simulated with 101-bp reads, as was done in the majority of the GTEx samples.
18 Error rate was set at 0.4 percent.

19 Next, we shuffled reads within each simulated experiment using a custom script. A non-random order
20 could present issues for many alignment and quantification methods, which perform a quick scan
21 through a subset of the reads for calibration.

22 Lastly, we combined reads generated from real transcripts, splicing noise, intronic noise and
23 intergenic noise transcripts together for each sample.

24 **Notebook:** [evaluation_stage2_genReads.ipynb](#)

1 **Ground truth alignment**

2 While not used in our final analysis, we found it useful to have a ground-truth alignment of the
3 simulated reads, since any alignment algorithm inevitably introduces errors due to sequencing errors,
4 splice-site detection, low-complexity or homologous regions. We designed a separate method to
5 convert standard Polyester and RSEM output into SAM formatted alignments.

6 **Notebook:** [evaluation_stage3_process.ipynb](#)

7 **Alignment of simulated reads and assembly-based quantification**

8 Reads were aligned with HISAT2 version 2.2 (Kim et al. 2019) against the indices constructed for this
9 analysis. We used a copy of HISAT2 from the rna_sensitive branch which introduces the "--rna-
10 sensitive" flag. This flag guarantees that for all multimapping reads at least one mapping will be
11 reported, unlike the default behavior of HISAT2, which skips reporting mappings for reads which align
12 to more than "-k" locations on the genome.

13 Alignments were sorted with samtools (Li et al. 2009) and assembled with StringTie2 (Kovaka et al.
14 2019) providing chess2.2_assembly.gff as guides.

15 **Notebook:** [evaluation_stage4_align_assemble.ipynb](#)

16 **Abundance quantification via pseudo-mapping**

17 Abundance of transcripts in each sample was estimated using Salmon version 0.14.0 (Patro et al.
18 2017) and Kallisto version 0.46.1 (Bray et al. 2016). Salmon was run with "--validateMappings" flag in
19 compliance with recommended parameters. For Kallisto, we set the read length to 101 using "-l"
20 option and the standard deviation was set to 0.0001 with the "-s" option.

21 **Notebook:** [evaluation_stage5_quantify.ipynb](#)

22 **Processing of transcript-level results**

23 We counted the number of simulated reads for each transcript in a sample using the metadata
24 provided by Polyester. Read counts were used to compute the TPM of all real and noisy transcripts in
25 each dataset. It must be noted that, since simulated transcripts came from independently assembled

1 data, the positions of the 3' and 5' ends may have been different from positions reported in the
2 annotation provided to all methods. In our analysis expression was measured based on the true
3 structure of the simulated molecules, to provide a more realistic simulated dataset.

4 Additionally, TPM values were computed separately for each of the combinations of noisy and real
5 transcripts, since fluctuations in the number of molecules of different effective lengths changes the
6 TPM values of all transcripts. However, in our analysis we refer to simulated TPM values only in the
7 context of abundance of false negatives reported by each method, while keeping other comparisons
8 strictly based on the values reported by each method.

9 Lastly, for StringTie2 we computed the raw number of reads assigned to each transcript based on the
10 reported coverage with respect to the effective length of the assembled molecule.

11 **Notebook:** [evaluation_stage6_get_results.ipynb](#)

12 **Processing of gene-level results**

13 To avoid differences in normalization and to measure raw effects of transcriptional noise, we
14 computed gene-level abundances using only raw read counts from real transcripts for each gene in
15 each sample. True read counts for each transcript in each sample were computed from the simulated
16 data and labeled by whether the transcript was simulated as part of a real or noisy transcription. To
17 obtain expected read counts per gene, transcript-level read counts were summed considering real
18 transcripts only. For each gene, we computed the number of reads which came from noisy transcripts.

19 **Notebook:** [evaluation_stage7_get_gene_expression.ipynb](#)

20 **Visualization of results**

21 All scripts for re-producing transcript-level visualizations from our analysis are provided in a Python

22 **notebook** at: [evaluation_stage8_results_tx.ipynb](#)

23 All scripts for re-producing gene-level visualizations from our analysis are provided in a Python

24 **notebook** at: [evaluation_stage9_results_gene.ipynb](#)

1 References

- 2 Bray NL, Pimentel H, Melsted P, Pachter L. 2016. Near-optimal probabilistic RNA-seq quantification. *Nature*
3 *Biotechnology* **34**: 525–527.
- 4 Frazee AC, Jaffe AE, Langmead B, Leek JT. 2015. Polyester: simulating RNA-seq datasets with differential
5 transcript expression. *Bioinformatics* **31**: 2778–2784.
- 6 Graves-Lindsay T, Albracht D, Fulton RS, Kremitzki M, Magrini V, Meltz K, et al. 2017. Evaluation of GRCh38 and
7 de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly.
8 *Genome Research* **27**: 849–864.
- 9 Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. 2019. Graph-based genome alignment and genotyping with
10 HISAT2 and HISAT-genotype. *Nature Biotechnology* **37**: 907–915.
- 11 Kovaka S, Zimin A v., Pertea GM, Razaghi R, Salzberg SL, Pertea M. 2019. Transcriptome assembly from long-
12 read RNA-seq alignments with StringTie2. *Genome Biology* **20**.
- 13 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The sequence
14 alignment/map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- 15 Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F, Young N, et al. 2013. The
16 Genotype-Tissue Expression (GTEx) project. *Nature Genetics* **45**: 580–585.
- 17 Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. 2017. Salmon provides fast and bias-aware quantification
18 of transcript expression. *Nature Methods* **14**: 417–419.
- 19 Pertea G, Pertea M. 2020. GFF Utilities: GffRead and GffCompare. *F1000Research* **9**.
- 20 Pertea M, Shumate A, Pertea G, Varabyou A, Breitwieser FP, Chang YC, Madugundu AK, Pandey A, Salzberg SL.
21 2018. CHES: A new human gene catalog curated from thousands of large-scale RNA sequencing
22 experiments reveals extensive transcriptional noise. *Genome Biology* **19**.
- 23 Swain PS, Elowitz MB, Siggia ED. 2002. Intrinsic and extrinsic contributions to stochasticity in gene expression.
24 *Proceedings of the National Academy of Sciences* **99**: 12795–12800.
- 25