

# FreePSI: an alignment-free approach to estimating exon-inclusion ratios without a reference transcriptome

## Supplementary information

Jianyu Zhou<sup>1,2</sup>, Shining Ma<sup>3</sup>, Dongfang Wang<sup>1</sup>, Jianyang Zeng<sup>4</sup> and Tao Jiang<sup>5,1,2</sup>

1 MOE Key Laboratory of Bioinformatics; Bioinformatics Division and Center for Synthetic  
& Systems Biology, TNLIST, Tsinghua University

2 Department of Computer Science and Technology, Tsinghua University

3 Department of Statistics, Stanford University

4 Institute for Interdisciplinary Information Sciences, Tsinghua University

5 Department of Computer Science and Engineering, University of California, Riverside

## Contents

<b>S1 Abundance flow graph</b>	<b>2</b>
S1.1 Formal definitions of $\alpha$ and PSI . . . . .	2
S1.2 Properties of the abundance flow graph . . . . .	2
<b>S2 Probabilistic generative model</b>	<b>5</b>
S2.1 Subscript rearrangement . . . . .	5
S2.2 Random variables and parameters . . . . .	5
S2.3 Graphical model . . . . .	5
S2.4 Relationship between $\gamma$ , $\theta$ and $\alpha$ . . . . .	6
S2.5 Theoretical distribution of k-mers . . . . .	7
<b>S3 Algorithms</b>	<b>9</b>
S3.1 Maximum likelihood estimation . . . . .	9
S3.2 The expectation-maximization algorithm . . . . .	9
S3.3 The conjugate gradient projection descent algorithm . . . . .	10
<b>S4 Some implementation details</b>	<b>12</b>
S4.1 Linear indexing algorithm . . . . .	12
S4.2 Techniques for improving the efficiency of CGPD . . . . .	12
S4.2.1 Offline computation for part of $\nabla Q_g^{\Pi}(\theta_g)$ . . . . .	12
S4.2.2 Replacing outer product of vectors . . . . .	13
<b>S5 Supplementary results and discussion</b>	<b>14</b>
S5.1 Simulated data evaluation . . . . .	14
S5.2 Real data evaluation . . . . .	16
S5.3 Impact of the quality of transcriptome assembly . . . . .	17
S5.4 Impact of k-mer length on FreePSI . . . . .	18
S5.5 Supplementary tables . . . . .	18
<b>S6 Software configurations</b>	<b>20</b>

## S1 Abundance flow graph

### S1.1 Formal definitions of $\alpha$ and PSI

Let  $\alpha_h$  denote the relative abundance of isoform  $h$  of a gene. For convenience, denote the junction segment formed by exon segments  $i$  and  $j$  as a pair  $(i, j)$ , where  $i \neq j$ . For simplicity, we will also denote an exon segment  $i$  as the identity pair  $(i, i)$ . All such pairs will uniformly be referred to as segments. The indicator variable  $\mathcal{I}_{hij}$  is defined as

$$\mathcal{I}_{hij} = \begin{cases} 1 & \text{if isoform } h \text{ covers segment } (i, j) \\ 0 & \text{otherwise} \end{cases}$$

where the segment  $(i, j)$  represents an exon segment if  $i = j$  or a junction segment otherwise. The parameter  $\alpha_{ij}$  for segment  $(i, j)$  is formally defined as

$$\alpha_{ij} = \sum_h \mathcal{I}_{hij} \alpha_h \quad (\text{S1.1})$$

Note that isoforms will not be explicitly dealt with in our method (*i.e.*, they will be marginalized out), but they are used here to help explain the construction of our model, and will be referred to frequently below.

PSI is defined as the ratio of the total relative abundance of all isoforms containing a given exon over the total relative abundance of all isoforms of the gene containing the exon. Hence, the PSI value of exon segment  $i$  can be expressed by the following equation:

$$\psi_i = \frac{\sum_h \mathcal{I}_{hii} \alpha_h}{\sum_h \alpha_h} = \frac{\alpha_{ii}}{\sum_h \alpha_h} \quad (\text{S1.2})$$

### S1.2 Properties of the abundance flow graph

In order to consider the first and last exon segments of an isoform, we extend the above indicator variables as follows:

$$\begin{aligned} \mathcal{I}_{hsi} &= \begin{cases} 1 & \text{if exon segment } i \text{ is the first exon segment of isoform } h \\ 0 & \text{otherwise} \end{cases} \\ \mathcal{I}_{hit} &= \begin{cases} 1 & \text{if exon segment } i \text{ is the last exon segment of isoform } h \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Then, the edge weights  $\alpha_{si}$  and  $\alpha_{it}$  in an abundance flow graph can be defined formally as

$$\alpha_{si} = \sum_h \mathcal{I}_{hsi} \alpha_h, \quad \alpha_{it} = \sum_h \mathcal{I}_{hit} \alpha_h \quad (\text{S1.3})$$

**Theorem S1.** *The following equalities hold:*

$$\sum_i \alpha_{si} = \sum_h \alpha_h = \sum_i \alpha_{it} \quad (\text{S1.4})$$

*Proof.* Since every isoform has a unique first exon segment, we have the following equality:

$$\sum_i \mathcal{I}_{hsi} = 1$$

Hence,

$$\begin{aligned} \sum_i \alpha_{si} &= \sum_i \sum_h \mathcal{I}_{hsi} \alpha_h \\ &= \sum_h \alpha_h \sum_i \mathcal{I}_{hsi} \\ &= \sum_h \alpha_h \end{aligned}$$

The second equality can be derived similarly. □

The abundance flow graph satisfies the flow conservation property, if we consider the edge weights as flow amounts. Figure S1 illustrates an example of this flow conservation property.

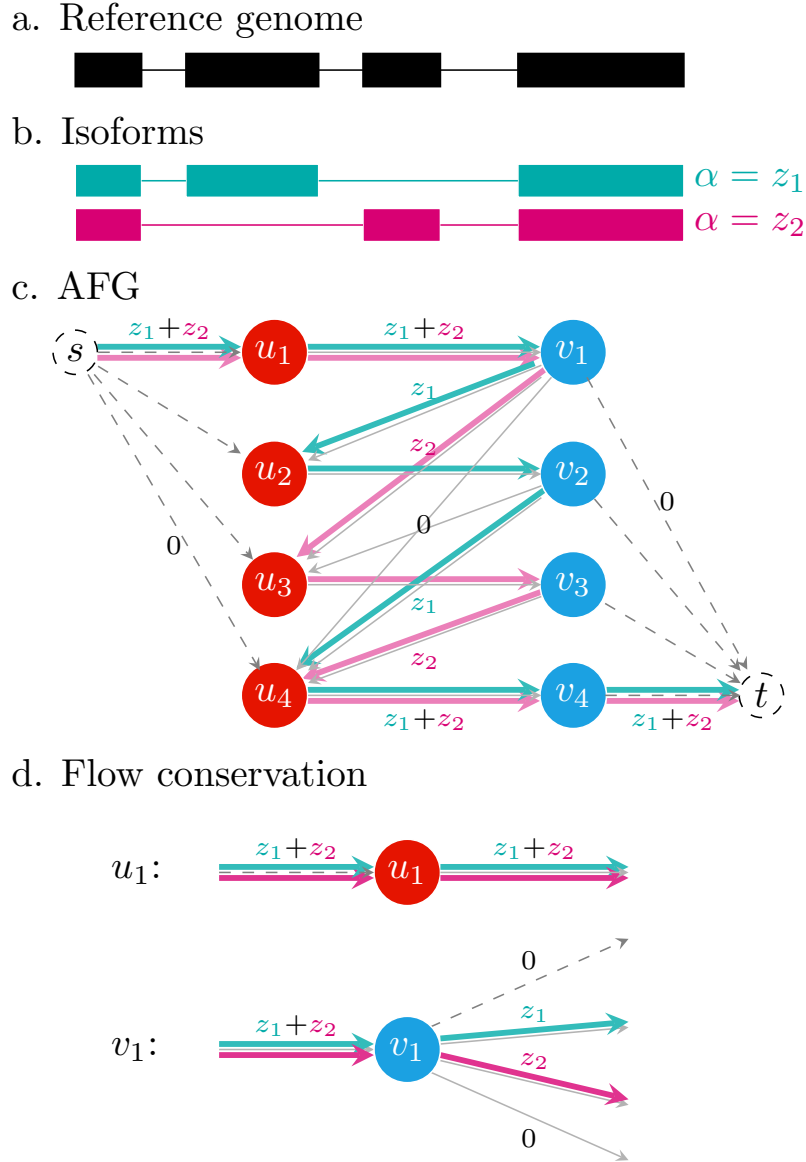


Figure S1: This figure illustrates the flow conservation property of an example AFG for a gene consisting of four exon segments and two isoforms. The exon boundaries are shown in (a) and the relative abundance of the two isoforms ( $z_1$  and  $z_2$ ) are shown in (b). The resulting AFG is shown in (c). The two isoforms are highlighted as two paths with colors cyan and pink from  $s$  to  $t$ , respectively. Considering the edge weights as flow amounts, the flow conservation property clearly holds for all vertices in  $U$  and  $V$ . In particular, Figure (d) shows that the flow conservation property holds for vertices  $u_1$  and  $v_1$ .

The following theorem provides a formal derivation of the flow conservation property.

**Theorem S2.**

$$\alpha_{ii} = \alpha_{si} + \sum_{j < i} \alpha_{ji}, \quad \text{for all } i \quad (\text{S1.5})$$

$$\alpha_{ii} = \alpha_{it} + \sum_{j > i} \alpha_{ij}, \quad \text{for all } i \quad (\text{S1.6})$$

*Proof.* We first prove that equality S1.5 holds. For any given isoform, an exon segment  $i$  of the isoform is either its first exon segment or next to exon segment  $j$  for some  $j < i$ . Therefore, for each isoform  $h$ ,

$$\mathcal{I}_{hii} = \mathcal{I}_{hsi} + \sum_{j < i} \mathcal{I}_{hji} \quad (\text{S1.7})$$

Multiplying both sides of Equation S1.7 by  $\alpha_h$  and summing up the equations for all isoforms  $h$ , we get

$$\begin{aligned} \sum_h \mathcal{I}_{hii} \alpha_h &= \sum_h \left( \mathcal{I}_{hsi} + \sum_{j < i} \mathcal{I}_{hji} \right) \alpha_h \\ \Rightarrow \sum_h \mathcal{I}_{hii} \alpha_h &= \sum_h \mathcal{I}_{hsi} \alpha_h + \sum_{j < i} \sum_h \mathcal{I}_{hji} \alpha_h \\ \Rightarrow \alpha_{ii} &= \alpha_{si} + \sum_{j < i} \alpha_{ji} \quad (\text{by Equations S1.1 and S1.3}) \end{aligned}$$

Hence, equality S1.5 holds. Equality S1.6 can be proven similarly.  $\square$

**Corollary S1.** *The PSI values can be rewritten as*

$$\psi_i = \frac{\alpha_{ii}}{\sum_i \alpha_{ii} - \sum_i \sum_{j > i} \alpha_{ij}} \quad (\text{S1.8})$$

*Proof.* By Equation S1.2, we only need to focus on the denominator.

$$\begin{aligned} \sum_h \alpha_h &= \sum_i \alpha_{si} \quad (\text{by Equation S1.4}) \\ &= \sum_i \left( \alpha_{ii} - \sum_{j < i} \alpha_{ji} \right) \quad (\text{by Equation S1.5}) \\ &= \sum_i \alpha_{ii} - \sum_i \sum_{j < i} \alpha_{ji} \\ &= \sum_i \alpha_{ii} - \sum_i \sum_{j > i} \alpha_{ij} \end{aligned}$$

$\square$

Since  $\alpha_{si} \geq 0$  and  $\alpha_{it} \geq 0$  for all  $i$ , we have the following constraints for the parameters  $\alpha_{ij}$ :

**Corollary S2.**

$$\alpha_{ii} \geq \sum_{j < i} \alpha_{ji}, \quad \text{for all } i \quad (\text{S1.9})$$

$$\alpha_{ii} \geq \sum_{j > i} \alpha_{ij}, \quad \text{for all } i \quad (\text{S1.10})$$

These constraints are crucial for an accurate estimation of our probabilistic model.

## S2 Probabilistic generative model

### S2.1 Subscript rearrangement

In order to perform a genome-wide analysis, we use a single index  $s$  to indicate a segment (exon or junction) of any gene  $g$ , and  $\alpha_{gs}$  to denote the total relative abundance of all isoforms containing segment  $s$  of gene  $g$ . Hence, the PSI value of exon segment  $i$  of gene  $g$  can be rewritten as

$$\psi_{gi} = \frac{\alpha_{gi}}{\sum_{\substack{s \in g \\ s \text{ is an exon segment}}} \alpha_{gs} - \sum_{\substack{s \in g \\ s \text{ is a junction segment}}} \alpha_{gs}} \quad (\text{S2.1})$$

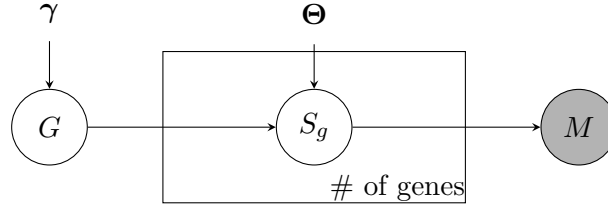
### S2.2 Random variables and parameters

Three random variables are defined for the probabilistic generative model.

1.  $G$  represents a random gene. Thus,  $P(G = g) = \gamma_g$  indicates that the probability that a read is generated from gene  $g$  is  $\gamma_g$ .
2.  $S_g$  represents a random segment of gene  $g$ , and  $P(S_g = s|G = g) = \theta_{gs}$  indicates that each read from gene  $g$  is generated from segment  $s$  with conditional probability  $\theta_{gs}$ .
3.  $M$  represents a random k-mer, and  $P(M = m)$  indicates the probability of k-mer  $m$  being observed.

### S2.3 Graphical model

The probabilistic generative model can be represented by the following graphical model:



The probability of observing a k-mer  $m$  can be calculated as

$$\begin{aligned} P(M = m) &= \sum_g \sum_{s \in g} P(M = m, S_g = s, G = g) \\ &= \sum_g \sum_{s \in g} P(M = m|S_g = s, G = g) P(S_g = s|G = g) P(G = g) \\ &= \sum_g P(G = g) \sum_{s \in g} P(S_g = s|G = g) P(M = m|S_g = s, G = g) \\ &= \sum_g \gamma_g \sum_{s \in g} \theta_{gs} P(M = m|S_g = s, G = g) \end{aligned}$$

where  $P(M = m|S_g = s, G = g)$  denotes the theoretical distribution of k-mers on segment  $s$ .

## S2.4 Relationship between $\gamma$ , $\theta$ and $\alpha$

Let  $L_{gs}$  denote the length of segment  $s$  of gene  $g$  and  $L_{gh}$  denote the length of isoform  $h$  of gene  $g$ . The number of possible starting sites for a read in the segment or isoform are

$$\begin{aligned}\tilde{L}_{gs} &= L_{gs} - L_{\text{read}} + 1 \\ \tilde{L}_{gh} &= L_{gh} - L_{\text{read}} + 1\end{aligned}$$

respectively, where  $L_{\text{read}}$  denotes the read length. If the relative abundance is measured by TPM,  $\alpha_{gh}$  can be calculated as

$$\alpha_{gh} = \frac{\frac{f_{gh}}{\tilde{L}_{gh}}}{\sum_g \sum_{h' \in g} \frac{f_{gh'}}{\tilde{L}_{gh'}}} \times 10^{-6} \quad (\text{S2.2})$$

where  $f_{gh}$  represents the number of reads from isoform  $h$  of gene  $g$ . Similarly, let  $f_{gs}$  denote the number of reads from segment  $s$  of gene  $g$ , and the indicator variable  $\mathcal{I}_{ghs}$  to indicate whether isoform  $h$  contains segment  $s$  of gene  $g$ . Assuming that the reads are uniformly distributed in each segment/isoform, we get

$$f_{gs} = \sum_{h \in g} \mathcal{I}_{ghs} f_{gh} \frac{\tilde{L}_{gs}}{\tilde{L}_{gh}} \quad (\text{S2.3})$$

Note that the reads are generated by an implicit two-level process in the above probabilistic model: from genes to isoforms and then to segments. Equation S2.3 marginalizes out isoforms to associate the reads with segments directly.

Since the reads are assumed to be generated from the probabilistic model, we have

$$\gamma_g \theta_{gs} \approx \frac{f_{gs}}{\sum_g \sum_{h \in g} f_{gh}} \quad (\text{S2.4})$$

The following proposition shows a relationship among  $\gamma$ ,  $\theta$  and  $\alpha$ .

**Proposition S1.**

$$\alpha_{gs} \approx \frac{Z_2}{Z_1} \frac{\gamma_g \theta_{gs}}{\tilde{L}_{gs}} \quad (\text{S2.5})$$

where  $Z_1 = \sum_g \sum_{h \in g} \frac{f_{gh}}{\tilde{L}_{gh}} \times 10^{-6}$  and  $Z_2 = \sum_g \sum_{h \in g} f_{gh}$ .

*Proof.*

$$\begin{aligned}\gamma_g \theta_{gs} &\approx \frac{f_{gs}}{Z_2} && \text{(by Equation S2.4)} \\ &= \frac{1}{Z_2} \sum_{h \in g} \mathcal{I}_{ghs} f_{gh} \frac{\tilde{L}_{gs}}{\tilde{L}_{gh}} && \text{(by Equation S2.3)} \\ &= \frac{\tilde{L}_{gs}}{Z_2} \sum_{h \in g} \mathcal{I}_{ghs} \frac{f_{gh}}{\tilde{L}_{gh}} \\ &= \frac{\tilde{L}_{gs}}{Z_2} \sum_{h \in g} \mathcal{I}_{ghs} \alpha_{gh} Z_1 && \text{(by Equation S2.2)} \\ &= \frac{Z_1}{Z_2} \tilde{L}_{gs} \sum_{h \in g} \mathcal{I}_{ghs} \alpha_{gh} \\ &= \frac{Z_1}{Z_2} \tilde{L}_{gs} \alpha_{gs} && \text{(by Equation S1.1)}\end{aligned}$$

□

Therefore, the definition of PSI given in Equation S2.1 can be rewritten in terms of  $\theta$  as follows:

$$\psi_{gi} \approx \frac{\frac{\theta_{gi}}{\bar{L}_{gi}}}{\sum_{\substack{s \in g \\ s \text{ is an exon segment}}} \frac{\theta_{gs}}{\bar{L}_{gs}} - \sum_{\substack{s \in g \\ s \text{ is a junction segment}}} \frac{\theta_{gs}}{\bar{L}_{gs}}} \quad (\text{S2.6})$$

The linear inequalities S1.9 and S1.10 can also be transformed into the constraints on  $\theta$ . For gene  $g$ , the parameter  $\theta$  should satisfy  $\mathbf{A}_g \boldsymbol{\theta}_g \geq 0$ , which is the matrix-form of the constraints. More specifically,  $\mathbf{A}_g$  is composed of three parts:

$$\mathbf{A}_g = \begin{pmatrix} \mathbf{D}_g \\ \mathbf{U}_g \\ \mathbf{B}_g \end{pmatrix} \quad (\text{S2.7})$$

Each row of  $\mathbf{D}_g \in \mathbb{R}^{N_e(g) \times N_s(g)}$  denotes the coefficients of the linear constraints between each exon segment and its downstream junction segments (*i.e.*, the inequality in Equation S1.10).  $N_e(g)$  denotes the number of exon segments in gene  $g$  and  $N_s(g)$  the number of (exon and junction) segments in gene  $g$ . The element at row  $i$  and column  $s$  is defined as

$$\mathbf{D}_g(i, s) = \begin{cases} \frac{1}{\bar{L}_{gs}} & \text{if } s \text{ is an exon segment } i \\ -\frac{1}{\bar{L}_{gs}} & \text{if } s \text{ is a junction segment beginning at exon segment } i \\ 0 & \text{otherwise} \end{cases} \quad (\text{S2.8})$$

Similarly, each row of  $\mathbf{U}_g \in \mathbb{R}^{N_e(g) \times N_s(g)}$  denotes the coefficients of the linear constraints between each exon segment and its upstream junction segments (*i.e.*, the inequality in Equation S1.9). Its elements are defined as

$$\mathbf{U}_g(i, s) = \begin{cases} \frac{1}{\bar{L}_{gs}} & \text{if } s \text{ is exon segment } i \\ -\frac{1}{\bar{L}_{gs}} & \text{if } s \text{ is a junction segment ending at exon segment } i \\ 0 & \text{otherwise} \end{cases} \quad (\text{S2.9})$$

$\mathbf{B}_g \in \mathbb{R}^{(N_s(g)-N_e(g)) \times N_s(g)}$  represents the non-negativity constraints on the relative abundance of junction segments, which is defined as

$$\begin{aligned} \mathbf{B}_g &= (\mathbf{0}, \mathbf{I}) \\ \text{where } \mathbf{0} &\in \mathbb{R}^{(N_s(g)-N_e(g)) \times N_e(g)} \\ \mathbf{I} &\in \mathbb{R}^{(N_s(g)-N_e(g)) \times (N_s(g)-N_e(g))} \end{aligned} \quad (\text{S2.10})$$

The non-negativity constraints on the relative abundance of exon segments are omitted here since they are already implied by the constraints in  $\mathbf{D}_g, \mathbf{U}_g$  and  $\mathbf{B}_g$ .

## S2.5 Theoretical distribution of k-mers

Let  $r$  represent a read and  $m_r$  represent a k-mer in  $r$ . Define the following indicator function:

$$\mathcal{I}(m, m_r) = \begin{cases} 1 & \text{if } m = m_r \\ 0 & \text{if } m \neq m_r \end{cases}$$

Recall  $c_{gsm} = P(M = m | S_g = s, G = g)$  denotes the theoretical distribution of k-mers on segment  $s$  of gene  $g$ , under the assumption that the reads are uniformly distributed on each segment. The following proposition is easy to prove.

**Proposition S2.**

$$c_{gsm} := P(M = m | S_g = s, G = g) = \frac{\mathcal{F}_{gsm}}{\tilde{L}_{gs}(L_{read} - K + 1)} \quad (\text{S2.11})$$

where

$$\mathcal{F}_{gsm} = \sum_{r \in s} \sum_{m_r} \mathcal{I}(m_r, m) \quad (\text{S2.12})$$

Finally, the probability of observing a k-mer  $m$  is

$$P(M = m) = \sum_g \gamma_g \sum_{s \in g} \theta_{gs} c_{gsm} \quad (\text{S2.13})$$

where  $\gamma_g$  and  $\theta_{gs}$  are model parameters to be estimated.



## S3 Algorithms

### S3.1 Maximum likelihood estimation

Let  $n_m$  denote the number of occurrences of k-mer  $m$  in the input RNA-seq reads. The likelihood of observing all k-mers in the input is

$$\begin{aligned}\mathcal{L}(\gamma, \Theta) &= \prod_m P(M = m)^{n_m} \\ \log \mathcal{L}(\gamma, \Theta) &= \sum_m n_m \log P(M = m) \\ &= \sum_m n_m \log \left( \sum_g \gamma_g \sum_{s \in g} \theta_{gs} c_{gsm} \right)\end{aligned}$$

The maximum likelihood estimation is to solve the following nonlinear constrained optimization:

$$\begin{aligned}\max \quad & \log \mathcal{L}(\gamma, \Theta) \\ \text{s.t.} \quad & \mathbf{A}_g \boldsymbol{\theta}_g \geq 0, \quad \text{for all gene } g \\ & \sum_{s \in g} \theta_{gs} = 1, \quad \text{for all gene } g \\ & \sum_g \gamma_g = 1, \quad \forall \gamma_g \geq 0, \quad \forall \theta_{gs} \geq 0\end{aligned}$$

### S3.2 The expectation-maximization algorithm

An initial feasible solution is obtained via the following algorithm.

---

<b>Algorithm S1:</b> Initial feasible solution construction	
<b>Input:</b> $c_{gsm}$ , $n_m$ and $\mathbf{A}_g$ (for all $g, s$ and $m$ )	
<b>Output:</b> $\boldsymbol{\theta}^{(0)}$ and $\gamma^{(0)}$	
1	<b>begin</b>
2	$Z_m \leftarrow \sum_g \sum_{s \in g} c_{gsm}, \quad \theta_{gs} \leftarrow \sum_m n_m \frac{c_{gsm}}{Z_m}, \quad \gamma_g \leftarrow \sum_{s \in g} \theta_{gs} \quad \triangleright$ distribute $n_m$ to segments
3	$Y \leftarrow \sum_g \gamma_g, \quad \gamma_g \leftarrow \frac{\gamma_g}{Y} \quad \triangleright$ normalize $\gamma$
4	<b>forall</b> $g$ <b>do</b>
5	<b>forall</b> row vector $\mathbf{a}_r$ in $\mathbf{A}_g$ with $\mathbf{a}_r \boldsymbol{\theta}_g < 0$ <b>do</b>
6	<b>if</b> $\theta_{gs} = 0$ for some exon segment $s$ <b>then</b>
7	set all $\theta_{gs'} \leftarrow 0$ for all adjacent junction segments $s'$ <span style="float: right;"><math>\triangleright</math> make <math>\theta</math> feasible</span>
8	<b>end</b>
9	<b>while</b> $\mathbf{a}_r \boldsymbol{\theta}_g < 0$ <b>do</b>
10	$\theta_{gs} \leftarrow \theta_{gs} \times 10$ for all exon segments $s$ <span style="float: right;"><math>\triangleright</math> make <math>\theta</math> feasible</span>
11	<b>end</b>
12	$X_g \leftarrow \sum_{s \in g} \theta_{gs}, \quad \theta_{gs} \leftarrow \frac{\theta_{gs}}{X_g} \quad \triangleright$ normalize $\theta$
13	<b>end</b>
14	<b>end</b>
15	<b>return</b> $\boldsymbol{\theta}^{(0)}$ and $\gamma^{(0)}$
16	<b>end</b>

---

In the E-step of the EM algorithm, we derive the expected log-likelihood using the current

estimation of  $\gamma_g^{(t)}$  and  $\theta_{gs}^{(t)}$  as follows:

$$\mathcal{Q}(\gamma, \Theta) = \sum_m n_m \sum_g \mu_{gm}^{(t)} \log \left( \gamma_g \sum_{s \in g} c_{gsm} \theta_{gs} \right) \quad (\text{S3.1})$$

where

$$\mu_{gm}^{(t)} = \frac{\gamma_g^{(t)} \sum_{s \in g} \theta_{gs}^{(t)} c_{gsm}}{\sum_g \gamma_g^{(t)} \sum_{s \in g} \theta_{gs}^{(t)} c_{gsm}}$$

By expanding the product in the logarithmic term of Equation S3.1,  $\mathcal{Q}(\gamma, \Theta)$  can be decomposed as the summation of two independent parts:

$$\mathcal{Q}(\gamma, \Theta) = \mathcal{Q}^I(\gamma) + \sum_g \mathcal{Q}_g^{\text{II}}(\theta_g) \quad (\text{S3.2})$$

where

$$\begin{aligned} \mathcal{Q}^I(\gamma) &= \sum_m \sum_g \mu_{gm}^{(t)} \log(\gamma_g) \\ \mathcal{Q}_g^{\text{II}}(\theta_g) &= \sum_m \mu_{gm}^{(t)} \log \left( \sum_{s \in g} \theta_{gs} c_{gsm} \right) \end{aligned}$$

The M-step of the algorithm is to maximize the expectation of the log-likelihood given in Equation S3.2. This is divided into two independent parts. The first part is to solve

$$\begin{aligned} \max \quad & \mathcal{Q}^I(\gamma) \\ \text{s.t.} \quad & \sum_g \gamma_g = 1, \quad \forall \gamma_g \geq 0 \end{aligned}$$

By using the Lagrangian multiplier method, a closed-form solution for this part can be derived:

$$\gamma_g^{(t+1)} = \frac{\sum_m \mu_{gm}^{(t)}}{\sum_m \sum_g \mu_{gm}^{(t)}}$$

The second part consists of a similar optimization problem for each gene  $g$ :

$$\begin{aligned} \max \quad & \mathcal{Q}_g^{\text{II}}(\theta_g) \\ \text{s.t.} \quad & \mathbf{A}_g \theta_g \geq 0, \quad \sum_{s \in g} \theta_{gs} = 1, \quad \forall \theta_{gs} \geq 0 \end{aligned} \quad (\text{S3.3})$$

Since a closed-form solution for this problem is unavailable due to the linear inequality constraints, the conjugate gradient projection descent (CGPD) algorithm is applied to solve the problem for all genes concurrently.

### S3.3 The conjugate gradient projection descent algorithm

The CGPD algorithm is an extension of the well-known gradient projection descent (GPD) algorithm [1]. For completeness, a pseudocode of the CGPD algorithm is given below.

---

**Algorithm S2:** The CGPD algorithm
 

---

**Input:**  $\theta_g, \mathbf{A}_g$ 
**Output:**  $\arg \max \mathcal{Q}_g^{\Pi}(\theta_g)$  s.t.  $\mathbf{A}_g \theta_g \geq 0, \sum_{s \in g} \theta_{gs} = 1, \forall \theta_{gs} \geq 0$ 

```

1 begin
2    $i \leftarrow 0, \mathbf{x}^{(i)} \leftarrow \theta_g, \mathbf{g}^{(i)} \leftarrow \nabla \mathcal{Q}_g^{\Pi}(\mathbf{x}^{(i)})$ 
3    $q \leftarrow 1, \mathbf{N}_q \leftarrow (1, 1, \dots, 1)^{\top}$ 
4   forall row vector  $\mathbf{a}_j$  in  $\mathbf{A}_g$  do
5     if  $\mathbf{a}_j \mathbf{x}^{(i)} = 0$  then
6        $\mathbf{N}_{q+1} \leftarrow (\mathbf{N}_q, \mathbf{a}_j^{\top})$ 
7        $q \leftarrow q + 1$ 
8     end
9   end
10   $\mathbf{H}_q^{(i)} \leftarrow \mathbf{I}^{q \times q}$ 
11  repeat
12     $\mathbf{s}^{(i)} = \mathbf{H}_q^{(i)} \mathbf{g}^{(i)}$  ▷ construct conjugate gradient direction
13     $\boldsymbol{\alpha} \leftarrow (\mathbf{N}_q^{\top} \mathbf{N}_q)^{-1} \mathbf{N}_q^{\top} \mathbf{g}^{(i)}$  ▷ construct projected gradient direction
14     $b_{jj} \leftarrow (\mathbf{N}_q^{\top} \mathbf{N}_q)^{-1}_{jj}$ 
15    if  $\mathbf{s}^{(i)} = \mathbf{0}$  and  $\forall \alpha_j \leq 0$  then
16      return  $\mathbf{x}^{(i)}$ 
17    else if  $\|\mathbf{s}^{(i)}\| \leq \max \left\{ \frac{1}{2} \alpha_j b_{jj}^{-1/2} \right\}$  then
18      update  $\mathbf{N}_q$  as  $\mathbf{N}_{q-1}$  ▷ deactivate a constraint
19      update  $\mathbf{H}_q^{(i)}$  as  $\mathbf{H}_{q-1}^{(i)}$ 
20       $q \leftarrow q - 1$ 
21    else
22       $\lambda^{(i)} \leftarrow \arg \max \mathcal{Q}_g^{\Pi}(\mathbf{x}^{(i)} + \lambda \mathbf{s}^{(i)})$  s.t.  $0 \leq \lambda \leq \lambda_{\text{bound}}$  ▷ perform line search
23      if  $\lambda^{(i)} = 0$  then
24        return  $\mathbf{x}^{(i)}$ 
25      else if  $\lambda^{(i)} = \lambda_{\text{bound}}$  then
26        update  $\mathbf{N}_q$  as  $\mathbf{N}_{q+1}$  ▷ activate a constraint
27        update  $\mathbf{H}_q^{(i)}$  as  $\mathbf{H}_{q+1}^{(i+1)}$ 
28         $q \leftarrow q + 1$ 
29      else
30        update  $\mathbf{H}_q^{(i)}$  as  $\mathbf{H}_q^{(i+1)}$  ▷ keep the constraints
31      end
32       $\mathbf{x}^{(i+1)} \leftarrow \mathbf{x}^{(i)} + \lambda^{(i)} \mathbf{s}^{(i)}, \mathbf{g}^{(i+1)} \leftarrow \nabla \mathcal{Q}_g^{\Pi}(\mathbf{x}^{(i+1)})$ 
33       $i \leftarrow i + 1$ 
34    end
35  end
36 end

```

---

## S4 Some implementation details

### S4.1 Linear indexing algorithm

---

**Algorithm S3:** Linear indexing algorithm

---

**Input:** sequence  $S$  with length  $L$   
**Output:** indices of all k-mers in  $S$

```

1 begin
2    $a \leftarrow 0$ 
3   for  $i = 1, \dots, K - 1$  do
4      $a \leftarrow (a \text{ lsh } 2) \text{ or } H(S_i)$ 
5   end
6    $mask \leftarrow (1 \text{ lsh } K) - 1$ 
7    $A \leftarrow \{\}$ 
8   for  $i = K, \dots, L$  do
9      $a \leftarrow (a \text{ lsh } 2) \text{ or } H(S_i)$ 
10     $a \leftarrow a \text{ and } mask$ 
11     $A \leftarrow A \cup \{a\}$ 
12  end
13  return  $A$ 
14 end

Input: base pair  $s$ 
Output: index of  $s$ 
15 function  $H$ 
16   switch  $s$  do
17     case 'A' do return 0
18     case 'C' do return 1
19     case 'G' do return 2
20     case 'T' do return 3
21   end
22 end

```

---

### S4.2 Techniques for improving the efficiency of CGPD

#### S4.2.1 Offline computation for part of $\nabla \mathcal{Q}_g^{\text{II}}(\theta_g)$

$\nabla \mathcal{Q}_g^{\text{II}}(\theta_g)$  can be represented by

$$\nabla \mathcal{Q}_g^{\text{II}}(\theta_g) = \text{diag}^{-1} \left( \mathbf{C}_g \mathbf{C}_g^\top \theta_g \right) (\mathbf{C}_g \boldsymbol{\mu}_g) \quad (\text{S4.1})$$

where  $\mathbf{C}_g \in \mathbb{R}^{n_s(g) \times n_k}$  is the matrix form of  $c_{gsm}$  and  $\boldsymbol{\mu}_g \in \mathbb{R}^{n_k \times 1}$  is the vector form of  $\mu_{gm}$ , with  $n_k$  denoting the number of k-mers and  $n_s(g)$  the number of segments in gene  $g$ . Assuming that the CGPD algorithm converges in  $T$  iterations, its time complexity is then  $O(T n_k n_s(g)^2)$ , if  $\nabla \mathcal{Q}_g^{\text{II}}(\theta_g)$  is computed directly according to Equation S4.1. Since only  $\theta_g$  is changed during the iterations,  $\mathbf{C}_g \mathbf{C}_g^\top$  and  $\mathbf{C}_g \boldsymbol{\mu}_g$  (the iteration-invariant parts) can be computed in advance. This way, the time complexity is reduced into  $O(n_k n_s(g)^2 + T n_s(g)^2)$ .

### S4.2.2 Replacing outer product of vectors

The CGPD algorithm performs many vector outer product operation in the following form:

$$\mathbf{H}^{(t+1)} = \mathbf{H}^{(t)} \pm \mathbf{q}\mathbf{q}^\top$$

where  $\mathbf{H} \in \mathbb{R}^{n \times n}$  and  $\mathbf{q} \in \mathbb{R}^{n \times 1}$ . A direct computation requires allocating  $n \times n$  new memory to store the matrix  $\mathbf{q}\mathbf{q}^\top$ , which is redundant and becomes an efficiency bottleneck of FreePSI. To speed up this frequent operation, we update  $\mathbf{H}$  by in-space column-wise operations as follows:

$$\mathbf{H}_{:,i}^{(t+1)} \leftarrow \mathbf{H}_{:,i}^{(t)} \pm q_i \mathbf{q} \quad \text{for } i = 1, \dots, n$$

which does not require temporary memory allocation.

## S5 Supplementary results and discussion

### S5.1 Simulated data evaluation

#### Genome-wide correlation on simulated data

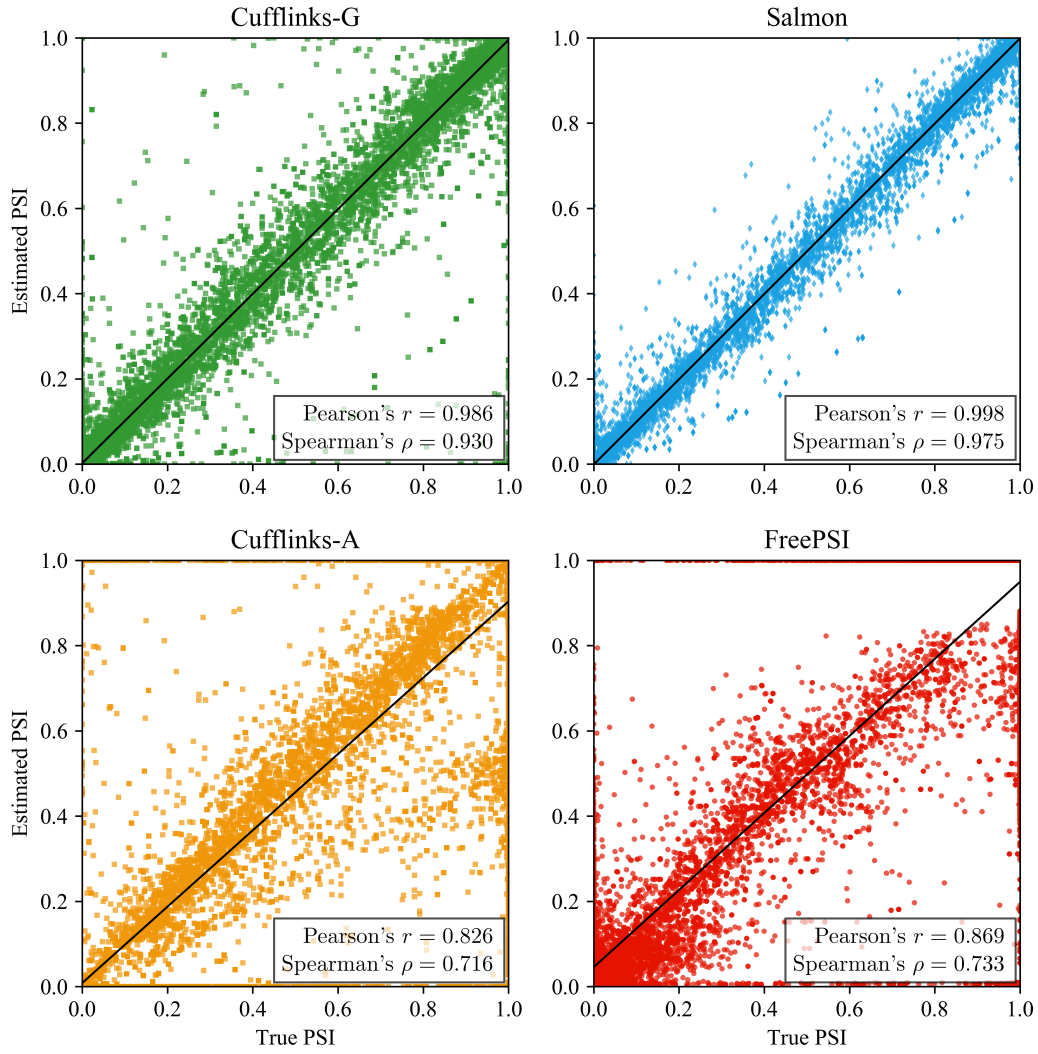


Figure S2: The scatter plot for genome-wide evaluation of different methods on the simulated data. The X-axis shows the true PSI values in the simulation and the Y-axis the PSI values estimated by different methods.

### Exon-centric correlation on simulated data

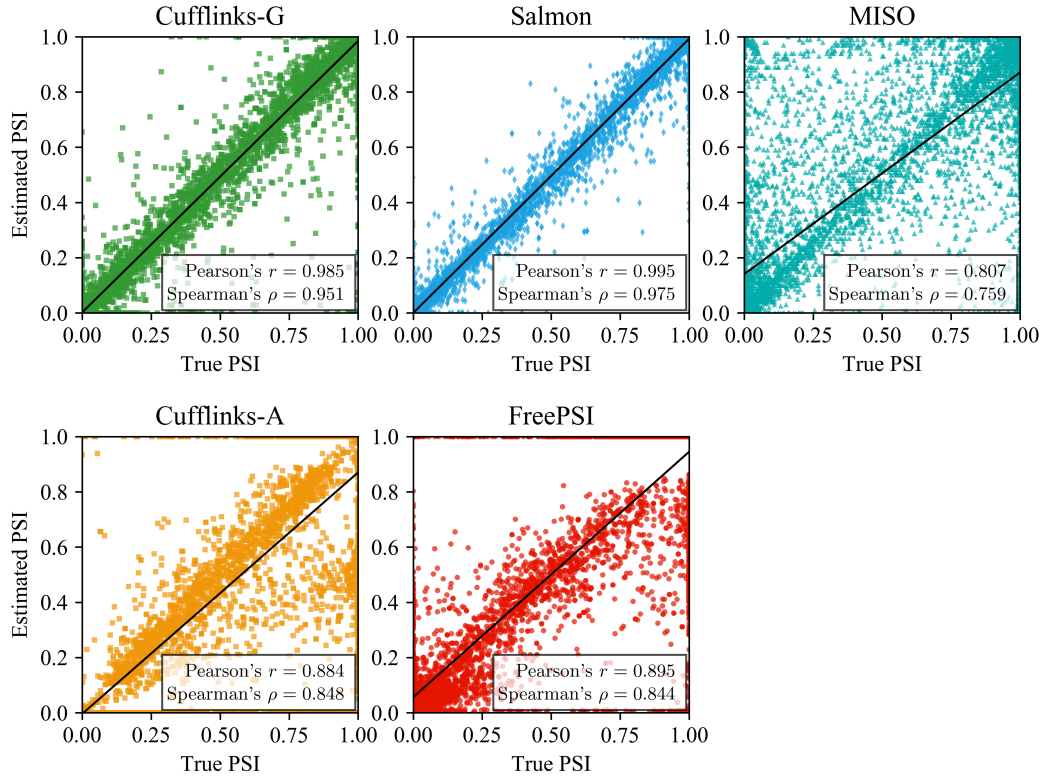


Figure S3: The scatter plot for exon-centric evaluation of different methods on the simulated data. The X-axis shows the true PSI values in the simulation and the Y-axis the PSI values estimated by different methods.

## S5.2 Real data evaluation

### Exon-centric correlation on real data

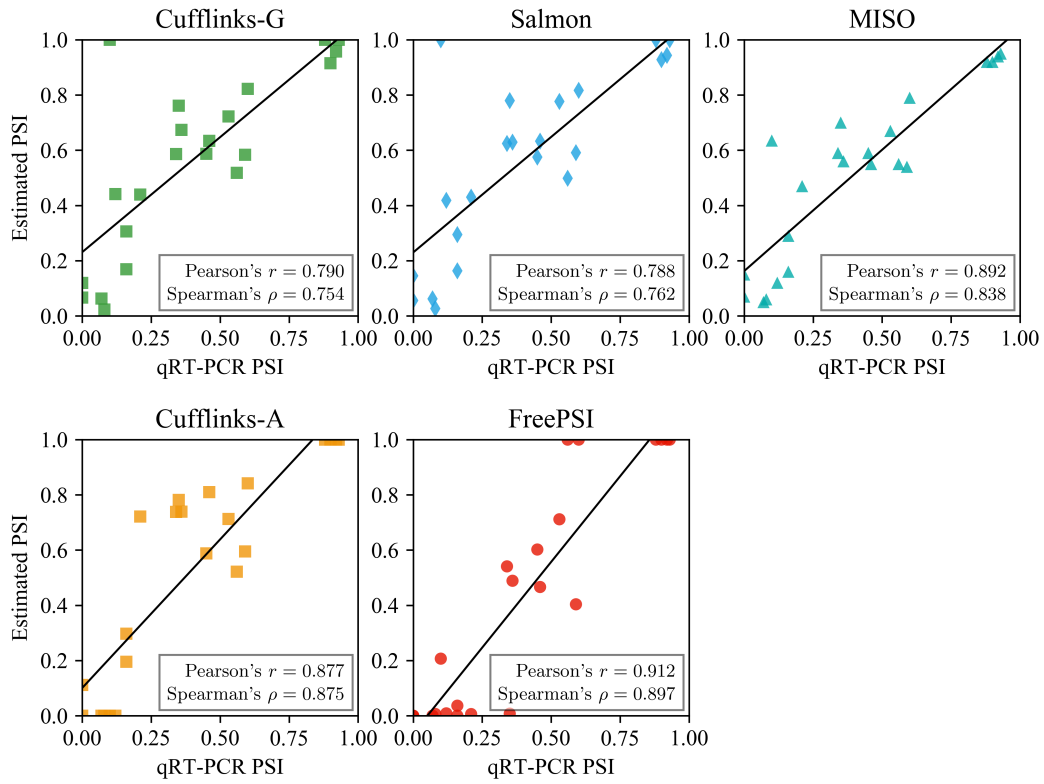


Figure S4: The scatter plot for exon-centric evaluation of different methods on the real data. The X-axis shows the true PSI values calculated from the qRT-PCR PSI results and the Y-axis the PSI values estimated by different methods.



### S5.3 Impact of the quality of transcriptome assembly

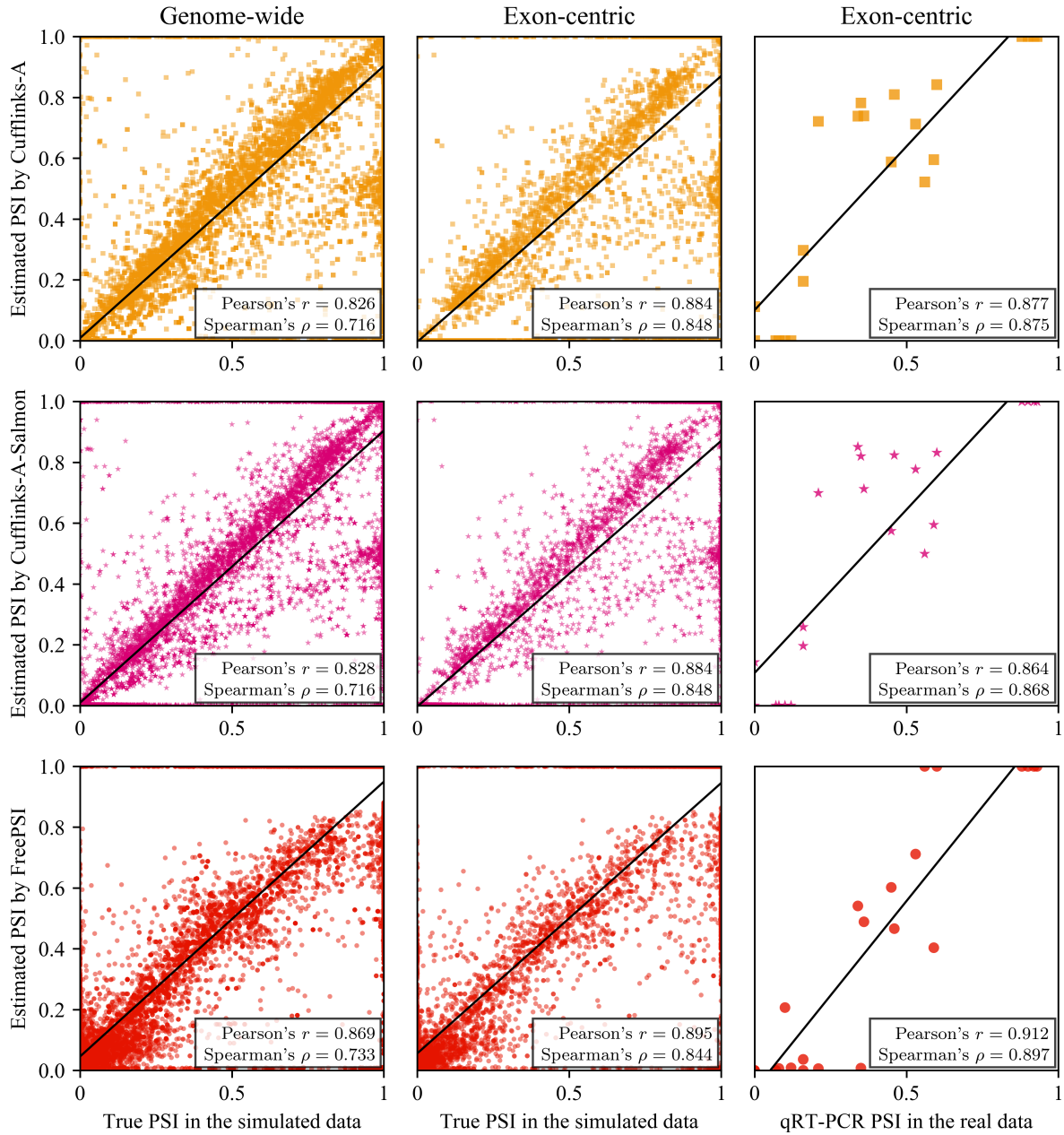


Figure S5: The scatter plots for Cufflinks-A (row 1), Cufflinks-A-Salmon (row 2) and FreePSI (row 3). Note that the plots for Cufflinks-A and FreePSI also appear in Figure 3 and Supplementary Figures S2, S3 and S4. We include them here again for the reader's convenience.

## S5.4 Impact of k-mer length on FreePSI

The parameter  $K$  representing the length k-mers considered in FreePSI is critical to the performance of FreePSI. We use the simulated dataset with 100 million reads to study the impact of  $K$ .

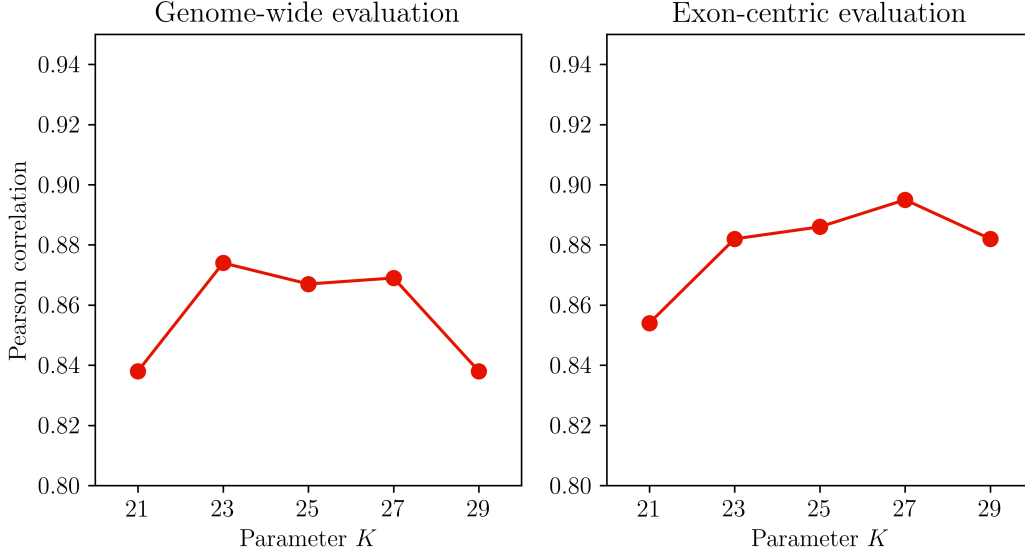


Figure S6: The performance of FreePSI under several choices of  $K$ .

The above figure shows the performance of FreePSI on the simulated dataset when different values of  $K$  was applied. In both evaluations, the performance peaked under a moderate  $K$ , although the optimal  $K$  values were different. The reason is that a smaller  $K$  induces more k-mers shared by different segments, which increases the difficulty of the estimation. On the other hand, a larger  $K$  results in fewer k-mers representing a segment, which makes the estimation more sensitive to sequencing errors. Hence, a moderate  $K$  in the range of 23 and 27 seems to work well for FreePSI generally, and we set the default  $K$  as 27.

## S5.5 Supplementary tables

Table S1: Impact of sequencing depth on the performance of MISO, Salmon, Cufflinks-A, and FreePSI on simulated data. The numbers of genes and exons selected for the genome-wide and exon-centric evaluations, respectively, are also shown in the table.

	# reads	20M	50M	100M
	# genes	6907	7025	7032
Pearson correlation for genome-wide evaluation	Salmon	0.997	0.998	0.998
	Cufflinks-A	0.783	0.836	0.826
	FreePSI	0.783	0.837	0.869
	# exons	10930	10958	10919
Pearson correlation for exon-centric evaluation	MISO	0.678	0.750	0.807
	Salmon	0.978	0.994	0.995
	Cufflinks-A	0.773	0.845	0.884
	FreePSI	0.817	0.862	0.895

Table S2: Performance of Salmon and Cufflinks-G on incomplete reference transcriptomes with different sampling rates. Here, a sampling rate represents what percentage of the true reference transcriptome would be covered in the provided reference transcriptome.

	Sampling rate	100%	90%	80%	70%
Pearson correlation for genome-wide evaluation	Salmon	0.998	0.913	0.804	0.725
	Cufflinks-G	0.986	0.904	0.796	0.718
Pearson correlation for exon-centric evaluation	Salmon	0.995	0.934	0.848	0.780
	Cufflinks-G	0.985	0.926	0.843	0.777

Table S3: Performance of Salmon, Cufflink-G, Cufflinks-A, and FreePSI on genes under different TPM thresholds.

TPM	0	1	2	5	10
Salmon	0.967	0.997	0.997	0.998	0.998
Cufflinks-G	0.964	0.982	0.983	0.985	0.986
Cufflinks-A	0.742	0.684	0.751	0.812	0.826
FreePSI	0.856	0.802	0.828	0.856	0.869

Table S4: Performance of Salmon, Cufflink-G, Cufflinks-A, and FreePSI on 14 gene families with large proportions of multi-mapped reads in the simulation. As a comparison, the corresponding numbers for the whole genome are given in the last row of the table.

Gene family	Multi-read proportion	FreePSI	Cufflinks-A	Salmon	Cufflinks-G	# genes	# isoforms	# expressed genes
PARI	38.16 %	0.968	0.230	1.000	0.434	21	103	8
GST	36.61 %	0.760	0.419	1.000	1.000	23	57	8
CDK	27.12 %	0.911	0.615	0.999	0.980	26	88	13
MAGE	25.64 %	0.978	0.541	0.999	0.997	39	103	7
NLR	15.90 %	0.805	0.524	1.000	0.998	23	60	9
TTC	13.37 %	0.709	0.698	0.999	0.984	115	373	46
TUB	10.78 %	0.907	0.755	0.999	0.998	24	59	11
SDR	9.95 %	0.907	0.825	1.000	0.999	75	179	23
NUP	7.88 %	0.536	0.884	1.000	0.994	32	89	15
DDX	5.63 %	0.943	0.867	0.994	0.974	42	108	17
CLEC	5.34 %	0.950	0.891	0.998	0.993	46	130	21
TRIM	5.23 %	0.931	0.849	1.000	0.997	80	189	26
SCAR	4.49 %	0.915	0.942	0.999	0.996	27	79	9
AKAP	3.04 %	0.947	0.849	1.000	0.882	29	73	10
Whole genome	2.48 %	0.869	0.826	0.998	0.986	23983	57822	7032

## S6 Software configurations

- Jellyfish on simulated data  
`jellyfish count -m 27 -s 100M -t 16 -Q 5 ${READS} -o ${OUTPUT}`
- Jellyfish on real data  
`jellyfish count -m 27 -s 100M -t 16 -Q A -L 10 ${READS} -o ${OUTPUT}`
- HISAT  
`hisat2 --fr --dta-cufflinks -p 16 -x ${GENOME_INDEX} -1 ${READS-1} -2 ${READS-2}  
| samtools view --threads 16 -Sbo ${BAM_FILE}`
- FreePSI  
`freePSI build -k 27 -p 16 -g ${REF_GENOME} -1 ${READS-1} -2 ${READS-2}  
-a ${EXON_BND} -o ${HASHTABLE}  
freePSI quant -k 27 -p 16 -i ${HASHTABLE} -o ${OUTPUT}`
- Salmon  
`salmon index -t ${REF_TRANSCRIPTOME} -i ${INDEX}  
salmon quant -p 16 -l ISF -i ${INDEX} -1 ${READS-1} -2 ${READS-2} -o ${OUTPUT}`
- Cufflinks-A  
`cufflinks -u -b ${GENOME_INDEX} -p 16 --library-type fr-secondstrand ${BAM}  
-o ${OUTPUT}`
- Cufflinks-G  
`cufflinks -u -b ${GENOME_INDEX} -p 16 --library-type fr-secondstrand ${BAM}  
-G ${REF_TRANSCRIPTOME} -o ${OUTPUT}`
- MISO  
`miso --run ${REF_SPLICING} ${BAM_FILE} --settings-filename=${DEFAULT_MISO_SETTING}  
-p 16 --read-len ${READ_LEN} output-dir ${OUTPUT}`
- Flux Simulator (.par file)  
`REF_FILE_NAME ${REF_TRANSCRIPTOME}  
GEN_DIR ${REF_GENOME}  
NB_MOLECULES 5000000  
TSS_MEAN 25  
POLYA_SCALE NaN  
POLYA_SHAPE NaN  
RTRANSCRIPTION YES  
RT_PRIMER RH  
FRAG_SUBSTRATE RNA  
FRAG_METHOD UR  
FRAG_UR_ETA NaN  
FRAG_UR_DO 1  
READ_NUMBER 100000000  
READ_LENGTH 76  
PAIRED_END YES  
ERR_FILE 76  
FASTA YES  
UNIQUE_IDS YES`

More details can be found in the source code.

## References

- [1] Jo Bo Rosen. The gradient projection method for nonlinear programming. part i. linear constraints. *Journal of the Society for Industrial and Applied Mathematics*, 8(1):181–217, 1960.