

Text S1. Supplementary Materials for IQSeq: Integrated Isoform Quantification Analysis Based on Next-generation Sequencing

Jiang Du¹, Jing Leng², Lukas Habegger², Andrea Sboner³, Drew McDermott¹, Mark Gerstein^{1,2,3,*}

1 Department of Computer Science, Yale University, New Haven, Connecticut, United States of America

2 Program in Computational Biology and Bioinformatics, Yale University, New Haven, Connecticut, United States of America

3 Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut, United States of America

* E-mail: Mark.Gerstein@yale.edu

Supplementary Materials

Maximum Likelihood Estimation

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \log(\Pr(S|\Theta)) \quad (1)$$

$$= \operatorname{argmax}_{\Theta} \log \prod_{m=1}^M \prod_{i=1}^{N_m} \Pr(s_{m,i}|\Theta, \text{Samp}_m) \quad (2)$$

$$= \operatorname{argmax}_{\Theta} \log \prod_{m=1}^M \prod_{s=s_{m,*}} \Pr(s|\Theta, \text{Samp}_m) \quad (3)$$

$$= \operatorname{argmax}_{\Theta} \log \prod_{m=1}^M \prod_{s=s_{m,*}} \sum_{k=1}^K \delta_{s,k} \theta_k \Pr(s|\Theta, \text{Samp}_m, I_k) \quad (4)$$

$$= \operatorname{argmax}_{\Theta} \log \prod_{m=1}^M \prod_{s=s_{m,*}} \sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)} \quad (5)$$

$$= \operatorname{argmax}_{\Theta} \sum_{m=1}^M \sum_{s=s_{m,*}} \log \sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)} \quad (6)$$

Expectation Maximization

$$\zeta_{s,k}^{(n)} = \mathbf{E}_{Z|S, \Theta^{(n)}} [Z_{s,k}] \quad (7)$$

$$= \mathbf{E} [Z_{s,k} | S, \Theta^{(n)}] \quad (8)$$

$$= \Pr(Z_{s,k} = 1 | S, \Theta^{(n)}) \quad (9)$$

$$= \frac{\Pr(Z_{s,k} = 1, s | \Theta^{(n)})}{\Pr(s | \Theta^{(n)})} \quad (10)$$

$$= \frac{\delta_{s,k} \theta_k^{(n)} G_{s,k}^{(m)}}{\sum_{k'=1}^K \delta_{s,k'} \theta_{k'}^{(n)} G_{s,k'}^{(m)}} \quad (11)$$

E step:

$$Q^{(n)}(\Theta) = \mathbf{E}_{Z|S, \Theta^{(n)}} [\log(\text{Pr}(Z, S|\Theta))] \quad (12)$$

$$= \mathbf{E}_{Z|S, \Theta^{(n)}} \left[\sum_{m=1}^M \sum_{s=s_{m,*}} \log \sum_{k=1}^K Z_{s,k} \theta_k G_{s,k}^{(m)} \right] \quad (13)$$

$$= \mathbf{E}_{Z|S, \Theta^{(n)}} \left[\sum_{m=1}^M \sum_{s=s_{m,*}} \sum_{k=1}^K Z_{s,k} \log \theta_k G_{s,k}^{(m)} \right] \quad (14)$$

$$\text{(for all } Z_{s,*}, \text{ one and only one can have a value of 1)} \quad (15)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} \sum_{k=1}^K \mathbf{E}_{Z|S, \Theta^{(n)}} (Z_{s,k}) \log \theta_k G_{s,k}^{(m)} \quad (16)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} \sum_{k=1}^K \zeta_{s,k}^{(n)} (\log \theta_k + \log G_{s,k}^{(m)}) \quad (17)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} \sum_{k=1}^K \zeta_{s,k}^{(n)} \log \theta_k + C \quad (18)$$

M step:

$$\theta_k^{(n+1)} = \frac{\sum_{m=1}^M \sum_{s=s_{m,*}} \zeta_{s,k}^{(n)}}{N} \quad (19)$$

$$= \frac{\sum_{m=1}^M \sum_{s=s_{m,*}} \frac{\delta_{s,k} \theta_k^{(n)} G_{s,k}^{(m)}}{\sum_{k'=1}^K \delta_{s,k'} \theta_{k'}^{(n)} G_{s,k'}^{(m)}}}{N} \quad (20)$$

Observed Fisher information matrix

$$\mathfrak{J}(\Theta)_{p,q} = -\frac{\partial^2 \log(\text{Pr}(S|\Theta))}{\partial \theta_p \partial \theta_q}, \text{ where } p, q = 1, \dots, K-1 \quad (21)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} -\frac{\partial^2 \log \sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)}}{\partial \theta_p \partial \theta_q} \quad (22)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} -\frac{\partial^2 \left[\log \left(\sum_{k=1}^{K-1} \delta_{s,k} \theta_k G_{s,k}^{(m)} + \delta_{s,K} G_{s,K}^{(m)} (1 - \sum_{l=1}^{K-1} \theta_l) \right) \right]}{\partial \theta_p \partial \theta_q} \quad (23)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} -\frac{\partial \left[\frac{1}{\sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)}} \left(\delta_{s,p} G_{s,p}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right) \right]}{\partial \theta_q} \quad (24)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} \frac{\left(\delta_{s,p} G_{s,p}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right) \left(\delta_{s,q} G_{s,q}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right)}{\left[\sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)} \right]^2} \quad (25)$$

Covariance matrix of the maximum likelihood estimator

Let $T(S) = (\hat{\theta}_1, \dots, \hat{\theta}_{K-1}, 1 - \sum_{k=1}^{K-1} \hat{\theta}_k)^T$, and $\psi(\Theta) = \mathbf{E}[T(S)]$. The Cramér-Rao bound states that:

$$\text{cov}_{\Theta}(T(S)) \geq \frac{\partial \psi(\Theta)}{\partial \Theta} [\mathcal{I}(\Theta)]^{-1} \left(\frac{\partial \psi(\Theta)}{\partial \Theta} \right)^T \quad (26)$$

, where $[\partial \psi(\Theta)/\partial \Theta]_{u,v} = \partial \psi_u(\Theta)/\partial \theta_v$, $u = 1, \dots, K$; $v = 1, \dots, K-1$.

We then estimate $\psi(\Theta)$ by Θ , and use the bound above to estimate the covariance matrix:

$$\frac{\partial \psi_u(\Theta)}{\partial \theta_v} \approx \frac{\partial \theta_u}{\partial \theta_v} \quad (27)$$

$$= \begin{cases} 1 & \text{if } u = v \text{ and } u, v < K; \\ -1 & \text{if } u = K; \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

$$\text{cov}_{\Theta}(T(S)) \approx \frac{\partial \psi(\Theta)}{\partial \Theta} [\mathcal{I}(\Theta)]^{-1} \left(\frac{\partial \psi(\Theta)}{\partial \Theta} \right)^T \quad (29)$$

$$= \begin{bmatrix} \mathbf{I}_{(K-1) \times (K-1)} & & \\ -1 & \dots & -1 \end{bmatrix}_{K \times (K-1)} \times \left([\mathcal{I}(\Theta)]^{-1} \right)_{(K-1) \times (K-1)} \quad (30)$$

$$\times \begin{bmatrix} & -1 \\ \mathbf{I}_{(K-1) \times (K-1)} & \vdots \\ & -1 \end{bmatrix}_{(K-1) \times K}, \quad \mathbf{I} \text{ is the identity matrix} \quad (31)$$

$$= \begin{bmatrix} \mathcal{I}_{(K-1) \times (K-1)}^{-1} & & \\ -\sum_{k=1}^{K-1} \mathcal{I}_{k,1}^{-1} & \dots & -\sum_{k=1}^{K-1} \mathcal{I}_{k,K-1}^{-1} \end{bmatrix}_{K \times (K-1)} \quad (32)$$

$$\times \begin{bmatrix} & -1 \\ \mathbf{I}_{(K-1) \times (K-1)} & \vdots \\ & -1 \end{bmatrix}_{(K-1) \times K} \quad (33)$$

$$= \begin{bmatrix} & & & -\sum_{k=1}^{K-1} \mathcal{I}_{1,k}^{-1} \\ & \mathcal{I}_{(K-1) \times (K-1)}^{-1} & & \vdots \\ -\sum_{k=1}^{K-1} \mathcal{I}_{k,1}^{-1} & \dots & -\sum_{k=1}^{K-1} \mathcal{I}_{k,K-1}^{-1} & -\sum_{k=1}^{K-1} \mathcal{I}_{K-1,k}^{-1} \\ & & & \sum_{i=1}^{K-1} \sum_{j=1}^{K-1} \mathcal{I}_{i,j}^{-1} \end{bmatrix}_{K \times K} \quad (34)$$

We also have:

$$\sum_{k=1}^K \text{var}(\hat{\theta}_k) = \text{tr}(\text{cov}_{\Theta}(\hat{\Theta})) \quad (35)$$

$$= \text{tr}(\text{cov}_{\Theta}(T(S))) \quad (36)$$

$$\approx \text{tr}([\mathcal{I}(\Theta)]^{-1}) + \sum_{i=1}^{K-1} \sum_{j=1}^{K-1} [\mathcal{I}(\Theta)]_{i,j}^{-1} \quad (37)$$

This means that we only need $\mathcal{I}(\Theta)$ in order to estimate the performance of our MLE with different sampling method combinations.

Efficient Computation of $\mathcal{I}(\Theta)$

$$\mathcal{I}(\Theta)_{p,q} = \mathbf{E} [\mathfrak{J}(\Theta)_{p,q}] \quad (38)$$

$$= \sum_{m=1}^M \sum_{s=s_{m,*}} \mathbf{E} \left[-\frac{\partial^2 \log \sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)}}{\partial \theta_p \partial \theta_q} \right] \quad (39)$$

$$= \sum_{m=1}^M N_m \mathcal{I}^{(m)}(\Theta)_{p,q} \quad (40)$$

where

$$\mathcal{I}^{(m)}(\Theta)_{p,q} = \mathbf{E}_{s \sim \text{Samp}_m} \left[-\frac{\partial^2 \log \sum_{k=1}^K \delta_{s,k} \theta_k G_{s,k}^{(m)}}{\partial \theta_p \partial \theta_q} \right] \quad (41)$$

is the expected Fisher information matrix of a single partial sample based on Samp_m . Thus we need to be able to compute $\mathcal{I}^{(m)}(\Theta)$ in order to obtain $\mathcal{I}(\Theta)$.

$$\mathcal{I}^{(m)}(\Theta)_{p,q} = \sum_{k=1}^K \theta_k \left\{ \sum_{s=s_{[a,b]}^{(k)}; \forall [a,b] \in I_k} -G_{s,k}^{(m)} \frac{\partial^2 \log \sum_{k'=1}^K \delta_{s,k'} \theta_{k'} G_{s,k'}^{(m)}}{\partial \theta_p \partial \theta_q} \right\} \quad (42)$$

$$= \sum_{k=1}^K \theta_k \sum_{s=s_{[a,b]}^{(k)}; \forall [a,b] \in I_k} G_{s,k}^{(m)} \mathfrak{J}_{s=s_{[a,b]}^{(k)}}^{(m)}(\Theta)_{p,q} \quad (43)$$

where

$$\mathfrak{J}_{s=s_{[a,b]}^{(k)}}^{(m)}(\Theta)_{p,q} = -\frac{\partial^2 \log \sum_{k'=1}^K \delta_{s,k'} \theta_{k'} G_{s,k'}^{(m)}}{\partial \theta_p \partial \theta_q} \quad (44)$$

$$= \frac{\left(\delta_{s,p} G_{s,p}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right) \left(\delta_{s,q} G_{s,q}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right)}{\left[\sum_{k'=1}^K \delta_{s,k'} \theta_{k'} G_{s,k'}^{(m)} \right]^2} \quad (45)$$

is the Fisher information matrix of a partial sample s from Samp_m at $[a, b]$ in I_k .

Proofs on Equivalent partial samples

$$\mathcal{I}^{(m)}(\Theta)_{p,q} = \sum_{k=1}^K \theta_k \left\{ \sum_{s=s_{[a,b]}^{(k)}; \forall [a,b] \in I_k} -G_{s,k}^{(m)} \frac{\partial^2 \log \sum_{k'=1}^K \delta_{s,k'} \theta_{k'} G_{s,k'}^{(m)}}{\partial \theta_p \partial \theta_q} \right\} \quad (46)$$

$$= \sum_{k=1}^K \theta_k \sum_{s=s_{[a,b]}^{(k)}; \forall [a,b] \in I_k} G_{s,k}^{(m)} \mathfrak{J}_{s=s_{[a,b]}^{(k)}}^{(m)}(\Theta)_{p,q} \quad (47)$$

where

$$\mathfrak{J}_{s=s_{[a,b]}^{(k)}}^{(m)}(\Theta)_{p,q} = -\frac{\partial^2 \log \sum_{k'=1}^K \delta_{s,k'} \theta_{k'} G_{s,k'}^{(m)}}{\partial \theta_p \partial \theta_q} \quad (48)$$

$$= \frac{\left(\delta_{s,p} G_{s,p}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right) \left(\delta_{s,q} G_{s,q}^{(m)} - \delta_{s,K} G_{s,K}^{(m)} \right)}{\left[\sum_{k'=1}^K \delta_{s,k'} \theta_{k'} G_{s,k'}^{(m)} \right]^2} \quad (49)$$

is the Fisher information matrix of a partial sample s from $Samp_m$ at $[a, b]$ in I_k .

Definition 1. Two partial samples s_1 and s_2 are equivalent w.r.t. $Samp_m$ if and only if $\mathfrak{J}_{s_1}^{(m)}(\Theta) = \mathfrak{J}_{s_2}^{(m)}(\Theta)$.

Lemma 1. If $\forall I_k \in I$, $\delta_{s_1,k} G_{s_1,k}^{(m)} = \delta_{s_2,k} G_{s_2,k}^{(m)}$, then s_1 and s_2 are equivalent w.r.t. $Samp_m$.

Proof. According to Equation 45, we have: $\mathfrak{J}_{s_1}^{(m)}(\Theta) = \mathfrak{J}_{s_2}^{(m)}(\Theta)$. Thus the two partial samples are equivalent w.r.t. $Samp_m$ according to Definition 1. \square

Definition 2. A set of partial samples S is an equivalent sample set w.r.t. $Samp_m$ if and only if $\forall s_1, s_2 \in S$, s_1 and s_2 are equivalent w.r.t. $Samp_m$.

Lemma 2. Given an isoform I_k and a sampling method $Samp_m$, if we divide all its possible partial samples into n non-overlapping equivalent sample sets S_1, S_2, \dots, S_n , then:

$$\mathcal{I}^{(m)}(\Theta)_{p,q} = \sum_{k=1}^K \theta_k \sum_{i=1}^n |S_i| G_{s_i,k}^{(m)} \mathfrak{J}_{s_i}^{(m)}(\Theta)_{p,q}, \text{ for any } s_i \in S_i \quad (50)$$

Proof. We can rewrite the $\sum_{s=s_{[a,b]}^{(k)}; \forall [a,b] \in I_k} G_{s,k}^{(m)} \mathfrak{J}_{s=s_{[a,b]}^{(k)}}^{(m)}(\Theta)_{p,q}$ part in Equation 43 by dividing all the possible $s_{a,b}^{(k,m)}$ s into equivalent sample sets S_1, S_2, \dots, S_n , and then obtain the equation above. \square

Definition 3. A simple shotgun sampling method $Samp_m$ generates samples with fixed read length r_m . When sampling from an isoform I_k with length l_k , there are in total $l_k - r_m + 1$ different samples $s_{[a,b]}^{(k)}$, where $a = 0, 1, 2, \dots, (l_k - r_m)$; and $b = a + r_m$. Each of these samples has equal probability of being generated from I_k : $G_{s,k}^{(m)} = 1/(l_k - r_m + 1)$.

Lemma 3. Given the sample generation model $Samp_m$ above, if two samples s_1 and s_2 generated by this method are compatible with the same set of isoforms, i.e. $\delta_{s_1,k} = \delta_{s_2,k}, \forall I_k \in I$, then s_1 and s_2 are equivalent w.r.t. $Samp_m$.

Proof. If $\delta_{s_1,k} = \delta_{s_2,k} = 0$, then obviously $\delta_{s_1,k} G_{s_1,k}^{(m)} = \delta_{s_2,k} G_{s_2,k}^{(m)} = 0$. Otherwise, if $\delta_{s_1,k} = \delta_{s_2,k} = 1$, then both s_1 and s_2 are partial samples that may be generated by $Samp_m$ from I_k . In this case, according to Definition 3, $G_{s_1,k}^{(m)} = G_{s_2,k}^{(m)} = 1/(l_k - r_m + 1)$. Thus we always have: $\delta_{s_1,k} G_{s_1,k}^{(m)} = \delta_{s_2,k} G_{s_2,k}^{(m)}, \forall I_k \in I$. According to Lemma 1, s_1 and s_2 are equivalent w.r.t. $Samp_m$. \square

Theorem 1. Given the sample generation model $Samp_m$ above, if two samples s_1 and s_2 generated by this method overlap with all the junctions in the same set of connected exons $e_{k_1} \rightarrow e_{k_2} \rightarrow \dots \rightarrow e_{k_n}$, then s_1 and s_2 are equivalent w.r.t. $Samp_m$.

Proof. We first prove by contradiction that $\forall I_k \in I, \delta_{s_1,k} = \delta_{s_2,k}$:

If $\delta_{s_1,k} \neq \delta_{s_2,k}$, without loss of generality, we assume that $\delta_{s_1,k} = 1$ and $\delta_{s_2,k} = 0$. Then there must exist an exon junction $e_{k_i} \rightarrow e_{k_{i+1}}, i \in \{1, 2, \dots, n-1\}$, such that $e_{k_i} \rightarrow e_{k_{i+1}}$ is not compatible with I_k (otherwise, as a part of $e_{k_1} \rightarrow e_{k_2} \rightarrow \dots \rightarrow e_{k_n}, s_2$ will be compatible with I_k). Since s_1 overlaps with the junction of $e_{k_i} \rightarrow e_{k_{i+1}}, s_1$ is not compatible with I_k either, which will lead to a contradiction to the previous assumption that $\delta_{s_1,k} = 1$. Thus the original statement $\delta_{s_1,k} = \delta_{s_2,k}$ must be true.

Then according to Lemma 3, s_1 and s_2 are equivalent w.r.t. $Samp_m$. \square

Use of empirical G function

To illustrate how non-uniform G function works, we modeled the bias of RNA-Seq data by aggregating signal of mapped reads along annotated transcripts. A signal map of the first base of mapped reads was generated. The signal was subsequently mapped onto the transcript and aggregated for all genes with signal isoform. The aggregation plot for Young Adult is shown in Figure S1. Each transcript is divided into 100 bins, accounting for their different lengths. The signals are normalized by the sum of signals in all the bins. The normalized signal at each bin represents the probability that a read is generated at certain position of the transcript. These non-uniform probabilities gave more realistic estimation of how the reads are generated, and were plugged into the EM calculations.

Results with different read generation assumptions

We have also carried out additional analysis comparing the result of IQSeq with uniform read generation assumption and with an empirical read generation model derived from the MAQC-3 dataset [1]. Figure S2 shows that while the overall results are still correlated, there are a significant number of isoforms with different estimated quantity under such different read generation assumptions.

Replicate variance vs. FIM based variance estimation

With the same MAQC-3 data, we also compared the isoform quantification variances between replicates with the FIM based variance estimations, and their logarithmic values have a correlation of 0.59 (Figure S3).

References

1. Bullard JH, Purdom E, Hansen KD, Dudoit S (2010) Evaluation of statistical methods for normalization and differential expression in mrna-seq experiments. BMC Bioinformatics 11: 94.

Figures

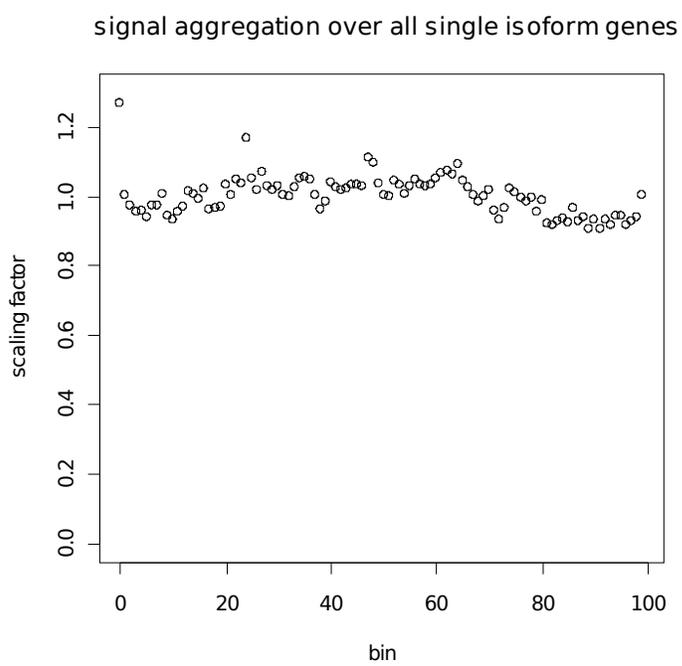
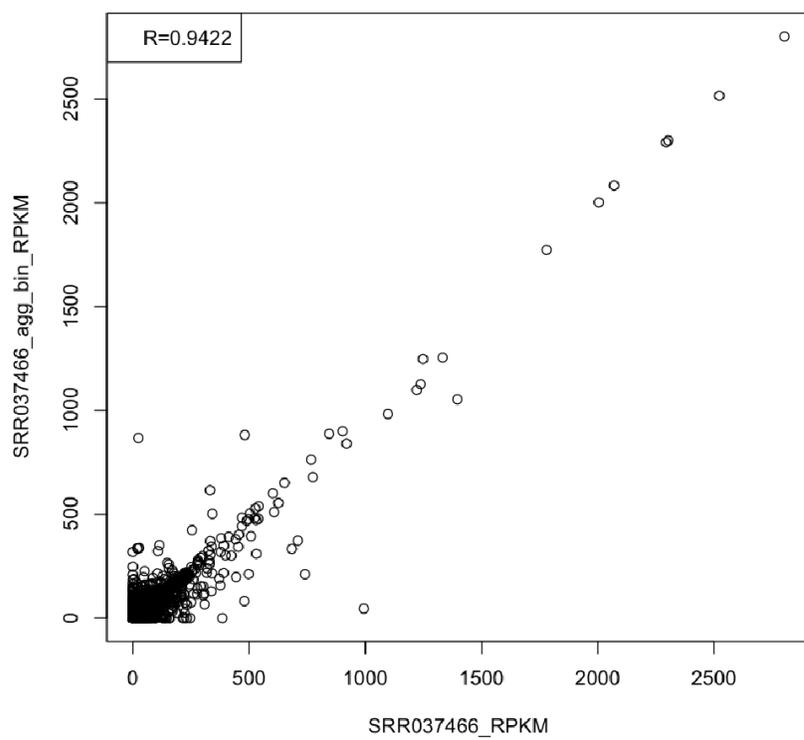
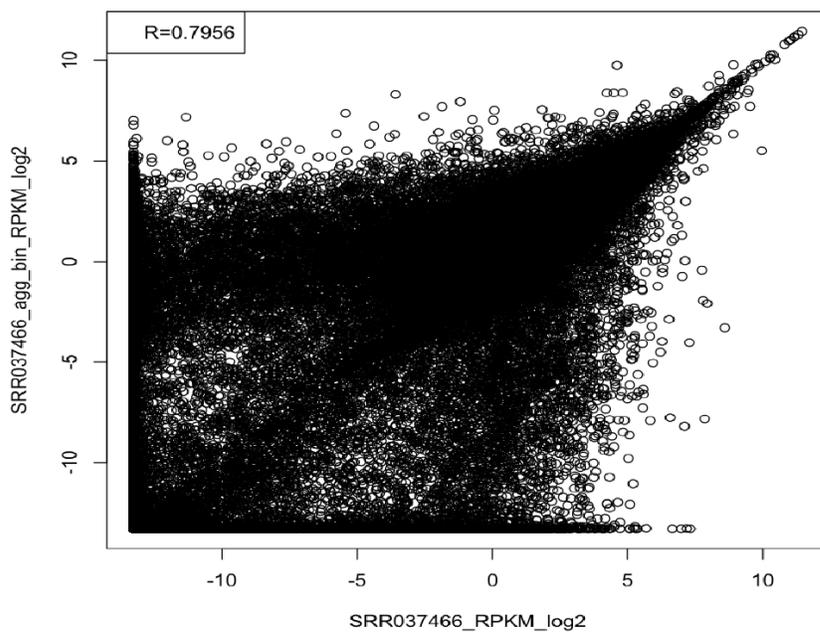


Figure S1. Aggregated signal over all single isoform genes



(a) Gene Level RPKM



(b) Isoform Level RPKM (log scale)

Figure S2. Results with uniform and non-uniform read generation models

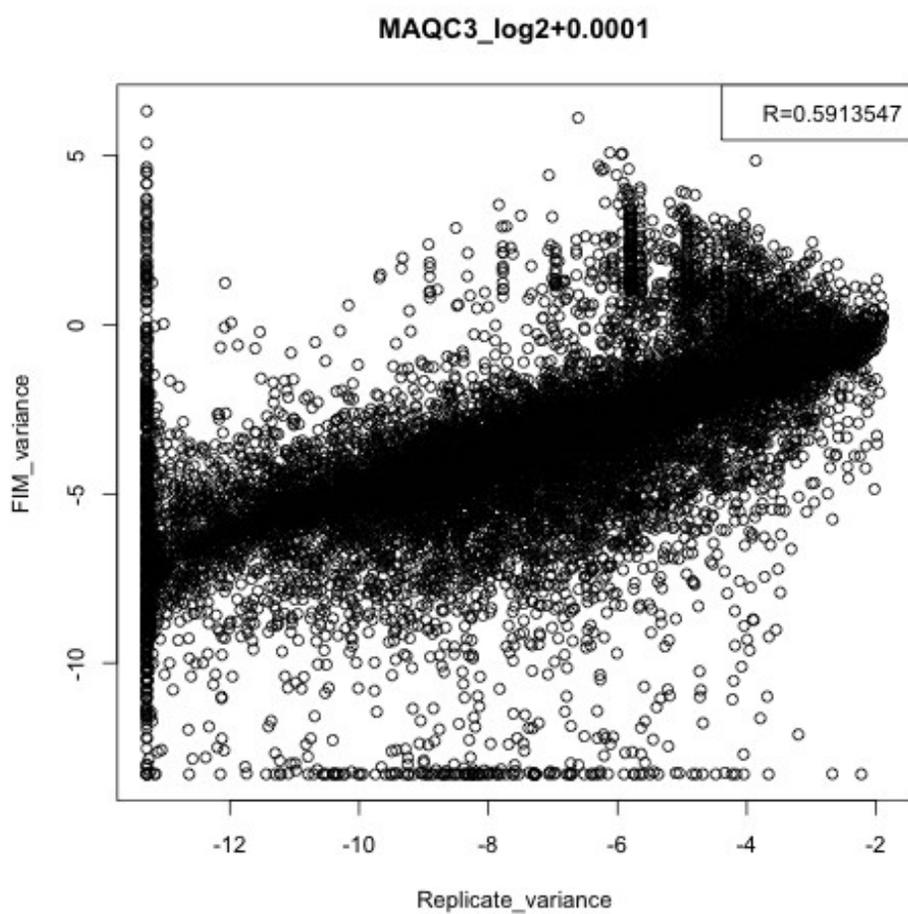


Figure S3. Replicate variance vs. FIM based variance estimation