



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Statistical modeling of isoform splicing dynamics from RNA-seq time series data

### Citation for published version:

Huang, Y & Sanguinetti, G 2016, 'Statistical modeling of isoform splicing dynamics from RNA-seq time series data', *Bioinformatics*, vol. 32, no. 19, pp. 2965-2972. <https://doi.org/10.1093/bioinformatics/btw364>

### Digital Object Identifier (DOI):

[10.1093/bioinformatics/btw364](https://doi.org/10.1093/bioinformatics/btw364)

### Link:

[Link to publication record in Edinburgh Research Explorer](#)

### Document Version:

Peer reviewed version

### Published In:

Bioinformatics

### General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



Gene expression

# Statistical modeling of isoform splicing dynamics from RNA-seq time series data

Yuanhua Huang<sup>1</sup> and Guido Sanguinetti<sup>1,2,\*</sup>

<sup>1</sup>School of Informatics, University of Edinburgh, Edinburgh, EH8 9AB, UK and

<sup>2</sup>Centre for Synthetic and Systems Biology (SynthSys), University of Edinburgh, Edinburgh, EH9 3BF, UK.

\*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

## Abstract

**Motivation:** Isoform quantification is an important goal of RNA-seq experiments, yet it remains problematic for genes with low expression or several isoforms. These difficulties may in principle be ameliorated by exploiting correlated experimental designs, such as time series or dosage response experiments. Time series RNA-seq experiments, in particular, are becoming increasingly popular, yet there are no methods that explicitly leverage the experimental design to improve isoform quantification.

**Results:** Here we present DICEseq, the first isoform quantification method tailored to correlated RNA-seq experiments. DICEseq explicitly models the correlations between different RNA-seq experiments to aid the quantification of isoforms across experiments. Numerical experiments on simulated data sets show that DICEseq yields more accurate results than state-of-the-art methods, an advantage that can become considerable at low coverage levels. On real data sets, our results show that DICEseq provides substantially more reproducible and robust quantifications, increasing the correlation of estimates from replicate data sets by up to 10% on genes with low or moderate expression levels (bottom third of all genes). Furthermore, DICEseq permits to quantify the trade-off between temporal sampling of RNA and depth of sequencing, frequently an important choice when planning experiments. Our results have strong implications for the design of RNA-seq experiments, and offer a novel tool for improved analysis of such data sets.

**Availability:** Python code is freely available at <http://diceseq.sf.net>.

**Contact:** G.Sanguinetti@ed.ac.uk

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

## 1 Introduction

In most eukaryotes, alternative splicing is an important post-transcriptional mechanism of regulation of gene expression, and largely increases the diversity of the proteome (Graveley, 2001). For example, over 90% of human genes have multiple isoforms (Wang *et al.*, 2008). Several lines of evidence indicate that alternative splicing plays a vital role in regulating biological processes (Blencowe, 2006), and its failure often causes serious diseases (Scotti and Swanson, 2016). The study of splicing has been revolutionised by the advent of high-throughput transcriptome sequencing (RNA-seq) techniques which enable unbiased sampling of the transcriptome and have greatly contributed to uncover novel biological

functions for alternative splicing (Wang *et al.*, 2009). More recently, RNA-seq technologies have been combined with biotin labelling treatment to provide kinetic measurements of RNA transcription and splicing with high temporal resolution (Windhager *et al.*, 2012; Barrass *et al.*, 2015; Eser *et al.*, 2016), providing invaluable mechanistic insights in the dynamics of splicing.

At the current stage of the technology, sequenced reads in RNA-seq experiments are much shorter than almost all eukaryotic transcripts. Thus, most reads from an RNA-seq experiment cannot be unambiguously aligned to a specific isoform. While in some cases a high level of coverage may obviate the problems, in many cases the number of reads that map to a single isoform is too low; when many isoforms are present, there may be no unambiguously assigned reads. To address this problem, several probabilistic methods were proposed to quantify the isoform proportion,

for example IsoEM (Nicolae *et al.*, 2011), Cufflinks (Trapnell *et al.*, 2010), MISO (Katz *et al.*, 2010), and BitSeq (Glaus *et al.*, 2012). All of these methods introduce latent variables to model the *identity* of a read, i.e. which isoform it came from, and then reconstruct isoform proportions by maximum likelihood or by computing a posterior distribution from the observed read distribution.

Most of these computational methods can quantify the isoform proportions accurately in many cases (Kanitz *et al.*, 2015), however for all methods isoform quantification at low coverages remains challenging. A natural approach in these cases is to exploit additional information, for example exploiting correlations across different experiments arising out of structured experimental designs such as time series or dosage response experiments. Time series RNA-seq designs, in particular, are becoming increasingly popular as an effective tool to investigate the dynamics of gene expression in a range of systems (Bar-Joseph *et al.*, 2012; Tuomela *et al.*, 2012; Zhang *et al.*, 2014; Honkela *et al.*, 2015). To our knowledge, no methods have been proposed that can exploit structured experimental designs in order to improve isoform estimation. This methodological gap also negatively affects the ability to design effectively experiments: for example, it is difficult to understand whether resources should be invested in gathering more time points, or in sequencing at a deeper level on a more limited number of samples.

In this article, we present a new methodology, DICEseq (Dynamic Isoform splicing Estimator via sequencing data) to jointly estimate the dynamics of isoform proportions from RNA-seq experiments with structured experimental designs. DICEseq is a Bayesian method based on a mixture model whose mixing proportions represent isoform ratios, as in (Katz *et al.*, 2010; Glaus *et al.*, 2012); however, DICEseq incorporates the correlations induced by the structured design by coupling the isoform proportions in different samples through a latent Gaussian process (GP). By doing so, DICEseq effectively transfers information between samples, borrowing strength which can aid to identify the isoform proportions. Our results show that DICEseq consistently improves in accuracy and reproducibility over the state of the art. This improvement can be very significant for a large fraction of genes: on one real data set, the correlation between estimates from replicate data sets increased by over 10% across one third of the genes as a result of taking temporal information into account. Furthermore, simulation studies indicate that DICEseq can be an important tool in experimental design, enabling an effective trade-off of resources between sequencing depth and sample numbers. DICEseq therefore offers an effective way to maximise information extraction from complex high-throughput data sets.

## 2 Methods

### 2.1 Mixture modelling of RNA-seq data

We briefly review here the mixture modelling framework for isoform identification (MISO), as described in (Katz *et al.*, 2010). We will describe the model on a per gene basis; the output of an RNA-seq experiment is therefore  $N$  reads  $R_{1:N}$  aligned to a gene with  $C$  isoforms. Each read  $R_n$  has its *identity*  $I_n \in \{1, \dots, C\}$ , i.e. which specific isoform it originated from, but, unless the read is aligned to isoform specific region, e.g., a junction, we will not know its identity. The proportion of each specific isoform within the pool of total mRNA is defined by the vector  $\Psi$ , whose entries must be positive and sum to 1. We can then define the likelihood of isoform proportions  $\Psi$  as mixture model as follows

$$P(R_{1:N}|\Psi) = \prod_{n=1}^N \sum_{I_n=1}^C P(R_n|I_n)P(I_n|\Psi). \quad (1)$$

The conditional distribution of  $I_n|\Psi$  is assumed to be Multinomial,  $(I_n|\Psi) \sim \text{Multinomial}(\Psi * w)$  where  $w$  is a weight vector adjusting the isoform proportion by the effective length of each isoform. The term  $P(R_n|I_n)$  encodes the probability of observing a certain read coming from a specific isoform  $I_n$ . This term automatically adjusts for the different informativeness of different reads: for example, junction reads will generally have a reduced number of possible isoforms (in extreme cases, only one), and as such will carry considerably more information through a reduced-entropy term  $P(R_n|I_n)$ . In this way, while the approach uses all sequenced reads for inference, the architectural information of the various transcript is still retained and automatically used. The model is completed by specifying a prior distribution over the isoform proportion vector  $\Psi$ , which in (Katz *et al.*, 2010) was chosen to be a Dirichlet on  $\Psi$ . Extending the MISO model to time series RNA-seq experiments involves a choice on how to model temporal correlations between the values of  $\Psi$  at different time points; we will use a flexible non-parametric prior in the form of a Gaussian process for this.

### 2.2 Gaussian processes

Gaussian processes (GPs) are a generalisation of the multivariate normal distribution to infinite-dimensional random functions. The key property of a GP is that all of its finite dimensional marginals are multivariate normals; in other words, evaluating a random function drawn from a GP at a finite set of points yields a normally distributed random vector. A GP over a suitable input space  $\mathcal{T}$  is uniquely specified by a mean function  $m: \mathcal{T} \rightarrow \mathbb{R}$  and a covariance function  $k: \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}$ , which models how correlations between function outputs depend on the inputs. In this paper, we will identify the input space  $\mathcal{T}$  with the time axis, and use as a covariance function the *squared exponential* (or RBF) covariance

$$k(t_1, t_2) = \theta_1 \exp\left(-\frac{1}{2\theta_2}(t_1 - t_2)^2\right). \quad (2)$$

The covariance function depends on two hyper-parameters, the prior variance  $\theta_1$  and the (squared) correlation lengthscale  $\theta_2$ .

The fundamental property of GPs relates the abstract function space view of GPs reported above with the explicit parametric form of their finite dimensional marginals. Let  $f$  denote a random function sampled from a GP,  $(t_1, \dots, t_N)$  denote a set of input (time) points and  $\mathbf{f} = (f(t_1), \dots, f(t_N))$  the vector obtained by evaluating the function  $f$  over the input points. Then, we have that

$$f \sim \mathcal{GP}(m, k) \leftrightarrow \mathbf{f} \sim \mathcal{N}(\mathbf{m}, K) \quad (3)$$

where  $\mathbf{m}$  and  $K$  are obtained by evaluating the mean and covariance functions over the set of points  $(t_1, \dots, t_N)$  (and pairs thereof). The fundamental property (3) is key to the success of GPs as a practical tool for Bayesian inference: given observations of the function values  $\mathbf{y}$ , it is in principle straightforward to obtain posterior predictions of the function values everywhere by applying Bayes' theorem

$$p(f(t_{new}|\mathbf{y})) \propto \int d\mathbf{f} p(\mathbf{f}, f(t_{new}))p(\mathbf{y}|\mathbf{f}) \quad (4)$$

If the observation noise model  $p(\mathbf{y}|\mathbf{f})$  is Gaussian, then the integral in (4) is analytically computable. Notice that equation (4) provides a way of predicting the latent function *at all time points*, not just the observation points. In the following, we describe an algorithm to approximate the computation of (4) for multinomial observations. For a thorough review of GPs and their use in modern machine learning, we refer the reader to the excellent book (Rasmussen and Williams, 2006).

### 2.3 Posterior of splicing dynamics with GP prior

Given a set of RNA-seq reads  $\mathbf{R} = [R_{1:N_1}^{(1)}, \dots, R_{1:N_T}^{(T)}]$  for  $T$  time points that are aligned to a gene with  $C$  isoforms, the posterior of the splicing dynamics for the isoform proportions  $\Psi = [\Psi_{1:C}^{(1)}, \dots, \Psi_{1:C}^{(T)}]$  is as follows,

$$\begin{aligned} P(\Psi|\Theta, \mathbf{R}) &\propto P(\Theta)P(\Psi|\Theta) \times \prod_{t=1}^T P(R_{1:N_t}^{(t)}|\Psi^{(t)}) \\ &\propto P(\Theta)P(\Psi|\Theta) \times \prod_{t=1}^T \prod_{n=1}^{N_t} \sum_{I_n^{(t)}=1}^C P(R_n^{(t)}|I_n^{(t)})P(I_n^{(t)}|\Psi^{(t)}) \end{aligned} \quad (5)$$

where  $\Psi$  is assumed as a **Softmax** function of latent variable  $Y$ , i.e.,  $\psi_c = e^{y_c} / \sum_{i=1}^C e^{y_i}$ , and  $y_c = 0$  to make the correspondence. Also  $Y_c = [y_c^{(1)}, \dots, y_c^{(T)}]$  follows a Gaussian process with its isoform specific hyperparameters  $\theta_c$  and mean  $m_c$ . By introducing the GP prior here, the joint analysis of time series RNA-seq data becomes possible, as shown in a cartoon in Figure 1.

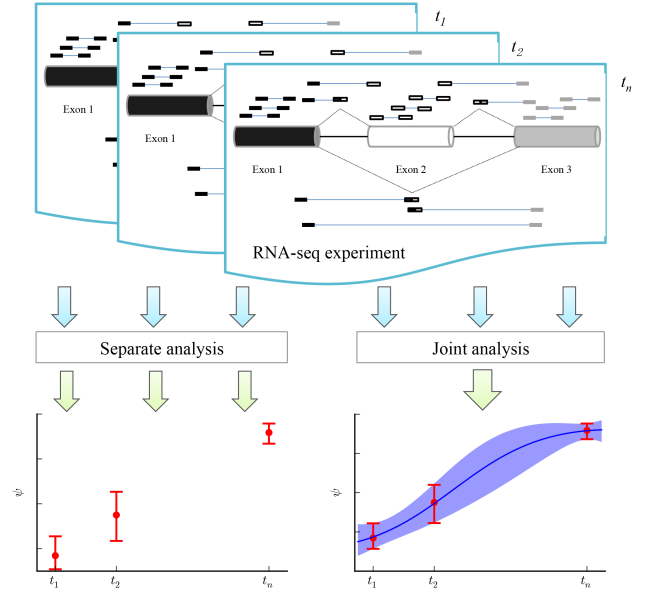
We assume in the following that the prior GP has zero mean, but this can be adjusted in a straightforward way to a more informative prior. Hyperparameters can also be sampled, however this leads to a much more complex inference problem since latent function values and hyperparameters are strongly correlated. We therefore fix  $\theta_{c,1} = 3.0$ , so that the 95% prior confidence intervals of  $\psi$  at an independent time point goes from 0.03 to 0.97, and set the second hyperparameter  $\theta_2$  empirically to account for approximately 20-40% of the duration of the experiment. A sensitivity analysis to  $\theta_2$  is provided in Supplementary Table S1 and Supplementary Figure S3. Inference of  $\theta_2$  can also be achieved by a straightforward extension of Algorithm 1 (see Supplementary Algorithm S1). However, this comes with a large additional computational cost, and in our experiments does not lead to improvements in accuracy; this is probably due to the fact that the typical RNA-seq time series is too short to carry enough information about the value of hyperparameters.

Having defined the posterior of the splicing dynamics, we introduce a Metropolis-Hasting sampler in Algorithm 1, which is a Markov chain Monte Carlo (MCMC) method, to infer the posterior of the splicing dynamics.

#### Algorithm 1 Metropolis-Hastings sampler for posterior of latent $Y$

**Require:**  $T, \mathbf{R}, \Theta, \lambda$   
**Initialize:**  $Y^{(0)}$   
**Calculate:**  $\Psi^{(0)} = \text{Softmax}(Y^{(0)})$ ;  $K = \text{GPCov}(\Theta, T)$   
**for**  $i = 0$  to  $H$  **do**  
  **Sample:**  $\mu \sim U(0, 1)$   
  **Sample:**  $Y^* \sim Q_y(Y^*|Y^{(i)}, \lambda K)$   
  **Calculate:**  $\Psi^* = \text{Softmax}(Y^*)$   
  **if**  $\mu < \min\left\{\frac{P(\Psi^*|\mathbf{R}) \times Q_y(Y^{(i)}|Y^*, \lambda K)}{P(\Psi^{(i)}|\mathbf{R}) \times Q_y(Y^*|Y^{(i)}, \lambda K)}, 1\right\}$  **then**  
     $Y^{(i+1)} \leftarrow Y^*$ ;  $\Psi^{(i+1)} \leftarrow \Psi^*$   
  **else**  
     $Y^{(i+1)} \leftarrow Y^{(i)}$ ;  $\Psi^{(i+1)} \leftarrow \Psi^{(i)}$   
  **end if**  
**end for**

Here, the proposal distribution  $Q_y$  for  $Y_c$  is a multivariate Gaussian distribution, whose mean is the last accepted  $Y_c^{(i)}$ , and the covariance matrix is defined by the fixed hyper-parameters  $\theta_c$  and the times  $T$ , but adjusted to the data itself, including the empirical variance of  $y$ , the number of isoforms, and number of time points, to ensure the 30-50% acceptance ratio. Namely,  $\hat{K}_c = \lambda K_c$ ;  $\lambda = (5\sigma_y^2)/(CT\theta_{c,1})$ , and the proposal distribution is  $\mathcal{N}(Y_c^{(i)}, \hat{K}_c)$ . Notice that, in contrast to the MISO



**Fig. 1.** A cartoon comparison between separate and joint analysis of time-series RNA-seq experiments. In the example gene, there are two isoforms with one alternative exon (the white one), and many paired-end reads are aligned to the genome for isoform quantification. The "separate analysis" estimates the isoform proportions for three time points independently, but the "joint analysis" estimate them together with a joint Gaussian process prior.

algorithm Katz *et al.* (2010), our sampler directly collapses the read identity variables, leading to considerable speedups when the number of isoforms is not too high.

For each gene, the initial MCMC chain contains 1000 iterations. Then the Geweke diagnostic  $Z$  score (Geweke, 1991) is applied to check the convergence of  $Y$ , using the first 10% and the last 50% iteration of the sampled chain. If  $|Z| > 2$ , then 100 more iterations will be added until the criterion is passed.

### 2.4 Reads probability and bias correction

DICSeq supports both single-end and paired-end reads. Here we describe the situation of paired-end reads; for single-end reads, just change the fragment length into read length. Given a read (pair)  $R_n$  mapping to an isoform  $c$ , the reads probability  $P(R_n|I_n = c)$  could be defined by taking information of the fragment length  $l_f$ , the alignment quality  $\text{mapq}$ , and the reads position  $p$ , as follow,

$$P(R_n|I_n) = P(l_f|I_n)P(p|I_n, l_f)P(R_n|\text{mapq}) \quad (6)$$

Here, we apply a Gaussian distribution to model the distribution of fragment length. The parameters (mean and variance) could be either set by user or learnt from the data itself. In some species, most reads can be very well mapped to a single position of the genome, and we could simply use uniquely mapped reads. However, in some other species, such as yeast, which contains many paralogs, there are higher chances to align a read to multiple positions. In the latter case for keeping multiply aligned reads, the  $\text{mapq}$  score will be taken into account, as  $P(R_n|\text{mapq}) = 1 - 10^{-\text{MAPQ}/10}$ , and we take the score of the better aligned mate for reads pair.

The reads position could be assumed to come from a uniform distribution, or could explicitly model sequence and position biases. In both cases, we could describe the probability as follows,

$$P(p|I_n = c, l_f) = \frac{b_c(p)}{\sum_{j=1}^{l_k - l_f + 1} b_c(j)} \quad (7)$$

where  $b_c(p)$  is relative weight of a position  $p$ . For uniform distribution,  $b_c(p) \equiv 1$ , so that  $P(p|I_n, l_f) = 1/(l_k - l_f + 1)$ . For the bias distribution, we employed the bias correction model that was proposed by Roberts *et al.* (Roberts *et al.*, 2011) to correct the position and sequence bias.

Briefly, Roberts *et al.*'s model of position bias tries to estimate which fractional position is preferred for sequencing. Thus, 20 bins from the beginning to the end of the isoform were used to count aligned reads, and isoforms are also divided into 5 groups based on their length. The sequence bias correction model tries to estimate the occurrence of a read with a surrounding sequence of each end from -8 to +12 nucleotides. A variable length Markov models were used to reduce the combinations of the 21 nucleotides, resulting in 774 parameters, as in (Roberts *et al.*, 2011). In DICEseq, we estimate these parameters empirically from the genes with only one isoform.

Empirically, we observed that correcting for biases did not significantly alter the results of our analyses, see Supplementary Table S2.

## 2.5 Gene annotation, input datasets and processing

Simulated reads in fastq format were generated from Spanki v0.5.0 (Sturgill *et al.*, 2013). It is based on the human gene annotation and genome sequences which were downloaded from GENCODE with release 22. In addition to exclusively keeping protein coding genes, we further removed those genes that only have one isoform (for these the problem is trivial) or overlap with others. Note that overlapping genes can also be accommodated by considering an extended isoform identification problem, whereby the identity of the read also includes the gene it belongs; this however requires a modification of the annotation file and was not considered for the purposes of illustrating our algorithm. Consequently, 90,759 isoforms from 11,426 genes were included for simulation. We randomly generated isoform ratios for each gene at 8 time points, with an assumption of either Gaussian process or first-order dynamics. Then the randomly generated isoform ratios were multiplied with the fixed library reads-per-kilobase (RPK, ranging from 50 to 1,600), to further define the number of isoform specific reads for the Spanki simulator.

4tU-seq data sets are available from the Gene Expression Omnibus (GEO; accession number GSE70378). The yeast gene annotation and genome sequences were downloaded from Ensembl with version R64-1-1, and all 309 intron-containing genes were included for analysis.

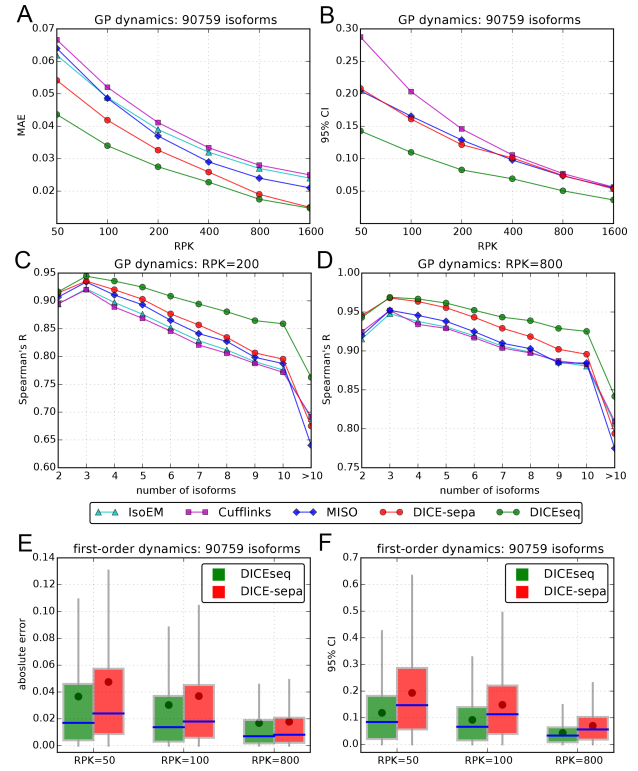
Circadian RNA-seq and microarray data sets on mouse liver were downloaded from GEO: GSE54652. The gene annotation and genome sequences were downloaded from GENCODE with release M6. Based on the annotation, we included 55,440 isoforms from 10,553 multiple-isoform, non-overlap, protein-coding genes. Processed microarray data (Zhang *et al.*, 2014), which are based on Affymetrix MoGene 1.0 ST, were employed for validation of the isoform estimate from RNA-seq. The microarray probe ids were mapped to GENCODE ids by Ensembl BioMart, leaving 30534 isoforms from 9755 genes for study.

All above RNA-seq data sets were downloaded in fastq format, and first aligned to corresponding Genome sequences above via HISAT 0.1.6-beta (Kim *et al.*, 2015), in paired-end mode with default setting.

## 3 Results

### 3.1 Methods comparison using simulated reads

In order to assess the performance of DICEseq, we compared it with three commonly used methods in their latest version: IsoEM v1.1.4 (Nicolae *et al.*, 2011), MISO v0.5.3 (Katz *et al.*, 2010), and Cufflinks v2.2.1 (Trapnell *et al.*, 2010). We also report results for a variant of DICEseq which ignores temporal correlations (DICE-sepa). Notice that DICE-sepa is essentially the same as MISO as a model, only differing in the estimation



**Fig. 2.** Comparison of accuracy between methods using simulated reads. (a) Mean absolute error between estimated isoform proportion and the truth. (b) 95% confidence interval of the estimates. (c) Influence of the number of isoforms on the estimates when RPK=200. (d) Influence of the number of isoforms on the estimates when RPK=800. The simulation is based on GP dynamics assumption for (a-d). (e) Boxplot of absolute error between estimated isoform proportion and the truth. (f) Boxplot of 95% confidence interval of the estimates. The round dot is the mean. The simulation is based on first-order dynamics assumption for (e-f).

procedure and prior (collapsed M-H sampler and softmax of a Gaussian). Simulated reads for 11,426 human protein coding genes, accounting for a total of 90,759 distinct isoforms, were generated by Spanki v0.5.0 (Sturgill *et al.*, 2013) with coverage from RPK of 50 to 1600 for 8 time points. We initially induced a temporal correlation between isoform proportions at different time points by enforcing the assumption of Gaussian process dynamics. All methods used paired-end reads, with the exception of MISO, which provided better performance in these experiments using single-end reads (see Supplementary Figure S2 and Table S3). We focus here on comparing the accuracy of the various methods; for a comparison of computational performance see Supplementary Figure S1.

We first studied the accuracy of each method at different coverage levels. We report average accuracy by computing the mean absolute error (MAE) between inferred isoform ratios and the truth from all the 90,759 isoforms of the 11,426 genes and 8 time points. Figure 2A shows that all methods return accurate estimates, and that the errors generally decrease with the increase of coverages. As expected, DICEseq is able to exploit effectively the temporal information, providing a significantly lower mean absolute error than the other methods, an advantage which is particularly marked at lower coverage. In a real RNA-seq time series experiment, many genes are likely to have relatively low coverage in at least one time point (see section 3.3 and 3.4 for our real data experiments), therefore the improved performance of DICEseq is likely to be important in quantifying isoforms for a substantial fraction of genes.



A second, often very important, metric is the confidence intervals associated with the predictions. These can be useful when deciding e.g. which genes to include in downstream analyses as in (Barrass *et al.*, 2015). We examined the average size of the confidence intervals for the three Bayesian methods DICEseq, MISO and Cufflinks as we vary the simulated coverage levels. As expected, confidence intervals shrink as we increase coverage for all three methods, however DICEseq clearly is able to provide more confident predictions at all coverage levels (Figure 2B). DICEseq is particularly strong at lower coverage; this is important, as often the confidence of an estimate is used to select genes which are further analysed (Barrass *et al.*, 2015).

Thirdly, we investigate the influence of isoform number on the quality of the estimate at a specified coverage level. By selecting the genes with a specific number of isoform, Figure 2C (RPK=200) and 2D (RPK=800) both show that the rank correlation (Spearman's correlation) coefficient between the estimated isoform proportions and the truth generally decreases as the number of isoform increases. This is expected, because the presence of more isoforms reduces the number of uniquely assignable reads. Once again, we see that including temporal information can yield significantly improved estimates, with DICEseq yielding an improvement in rank correlation of more than five percentage points for genes with many isoforms ( $>8$ ).

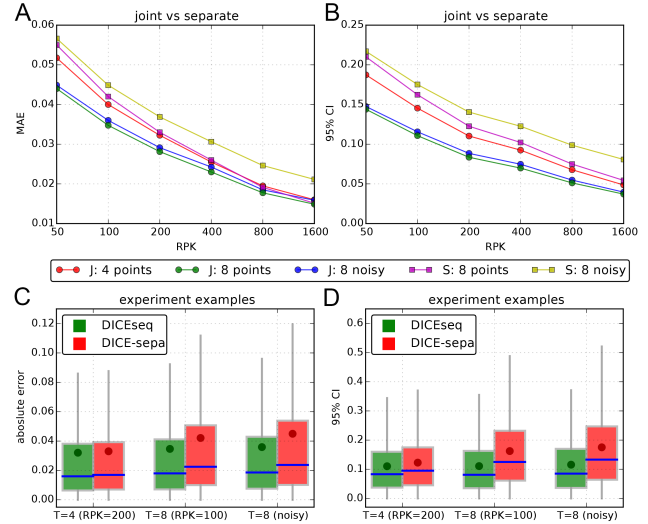
Finally, we investigate the robustness of DICEseq to model mismatch. To do so, we generated time series data where the isoform proportions vary according to a first-order dynamical system (rather than a Gaussian process), a commonly used modelling hypothesis (Eser *et al.*, 2016). Figure 2E-F clearly shows that incorporating temporal information yields a considerable improvement, even under model mismatch. This improvement is particularly marked at low coverages. Notice that the mean accuracy (represented by a dot in the box plots) is very similar to the one obtained under the GP assumption (Figure 2A). Additional simulations varying hyperparameters were also performed (see Supplementary file), and Supplementary Table S1 again shows robustness to mis-specification of the hyperparameters.

In summary, the results of these simulation studies show that DICEseq can provide accurate reconstruction of isoform proportions, and can successfully leverage temporal information to provide more accurate and confident predictions at low coverage and for higher numbers of isoforms.

### 3.2 Design of time-series RNA-seq experiments

Incorporating temporal information in the analysis of time series experiments is desirable in principle, because it provides experimentalists with a further direction for experimental design. Intuitively, resources can be invested in either improving the accuracy of each time point (by sequencing deeper), or by collecting more time points. This is an important trade-off, and it can only be achieved if the data is analysed jointly. To address these questions, we compared DICEseq versus DICE-sepa as we vary coverage levels and number of time points, by simulating reads as in the previous section (under GP assumption). In Figure 3A, we clearly see again that with the coverage increasing, all MAE decrease. In the joint model, the MAE largely decreases when more time points (i.e., 8) are used, especially for the case with low coverage.

These results highlight the importance of the analysis method for experimental design: while with the non-temporal model DICE-sepa increasing coverage is the only way to improve accuracy, methods that incorporate temporal information can benefit both from an increase in coverage and an increase in sampling frequency. Broadly speaking, we see that a doubling of the sampling frequency is roughly equivalent to a doubling of the sequencing depth, with the obvious advantage that a finer temporal information is provided. Figure 3C and 3D show an example of this trade-off: 4 time points and higher coverage of  $RPK = 200$  give



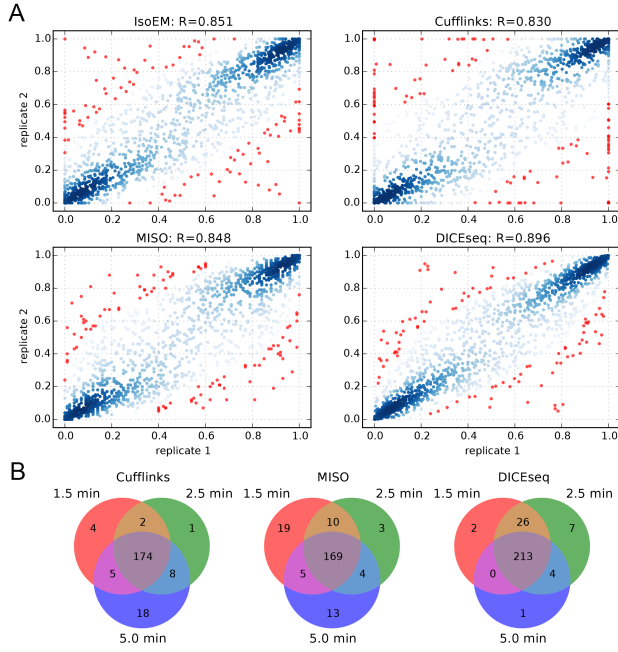
**Fig. 3.** Comparison between experiment design on time points and coverages. (a) Mean absolute error between estimated isoform proportion and the truth for different experiments. "S" and "J" means DICEseq separate and joint mode, respectively, and the "noisy" means the RPK=25 at the 5th point for all (a-d). All simulation here is based on GP dynamics assumption. (b) 95% confidence interval of the estimates. (c) Boxplot of absolute error between estimated isoform proportion and the truth for three example experiments. The round dot is the mean. "T=4" and "T=8" means the 4 and 8 number of time points. The "noisy" example was also conducted at RPK=100. (d) Boxplot of 95% confidence interval of the estimates.

indistinguishable results for the joint model to 8 time points and lower coverage of  $RPK = 100$  (first two pairs in Figure 3C/D).

Another potential advantage of incorporating temporal information is to improve robustness of the estimation against noise/low coverage at some time points. This aspect is particularly important as of course coverage level for a particular gene is largely determined by the gene's expression level, therefore genes with a large dynamic range of expressions during the time series will necessarily have some time points with low coverage. To simulate this situation, we generated time series with very low coverage ( $RPK = 25$ , termed "noisy") in the 5th time point. From the "noisy" case in Figure 3, we could see that the joint model dramatically reduces the variation compared to the separated model. Thus, incorporating time information in the joint model leads to a more robust estimation, facilitating isoform estimation for genes with dynamic expression levels and providing a possibility to combine low coverage with high coverage time points for time series libraries.

### 3.3 RNA splicing dynamics with 4tU-seq data

Recently, biotin labelling combined with RNA-seq has become an important tool to study the kinetics of RNA transcription and splicing with high temporal resolution (Windhager *et al.*, 2012; Veloso *et al.*, 2014; Fuchs *et al.*, 2014). These experiments naturally produce RNA-seq data sets with high temporal resolution; furthermore, at very early time points, labelled RNA may be of low abundance, resulting in high uncertainty estimates. Here we use a recent data set with high temporal resolution to probe the suitability of DICEseq as an analysis tool for biotin labelled RNA-seq; the data was produced by our collaborators in the Beggs and Granneman labs at the Wellcome Trust Centre for Cell Biology in Edinburgh (Barrass *et al.*, 2015). The data set consists of approximately 50M mapped reads; roughly 50% of genes have a coverage of  $RPK < 120$  in at least one time point.

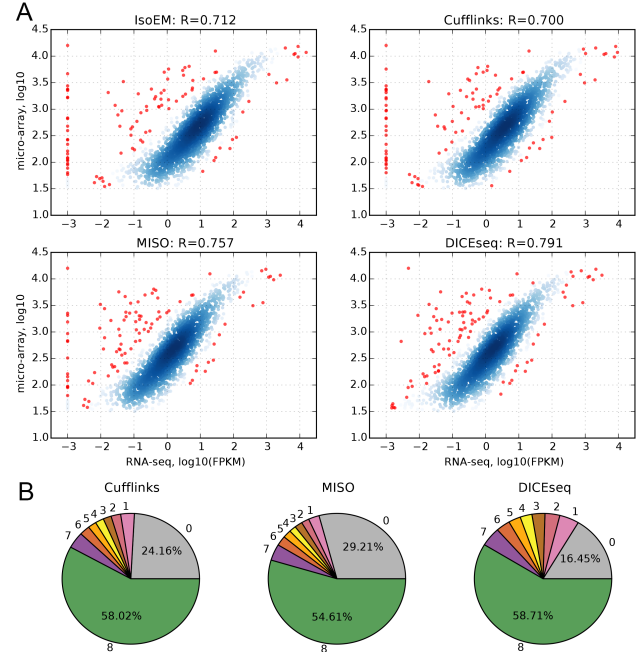


**Fig. 4.** Analysis of time series 4tU-seq data. (a) The Pearson's correlation between two replicates. (b) The number of genes whose 95% confidence interval < 0.3.

To assess accuracy of our method, we compare the correlation between two replicates for 309 intron-containing genes at 1.5, 2.5 and 5.0 minute. Figure 4A shows that IsoEM, Cufflinks and MISO all result in a good correlation between replicates, with Pearson's correlation coefficient varying between 0.83 and 0.85; DICEseq further improves with a Pearson's correlation coefficient of 0.896, outperforming by between 4 and 6 percentage points existing methods (all  $p$ -values <  $1e-5$  under the Fisher  $r$ -to- $z$  transform test (Diedenhofen and Musch, 2015)). The improvement is particularly marked if we consider the lowest expressed genes (see Table 1): on the lower third of the expression range, DICEseq still obtains a Pearson correlation of 0.860, while the other methods achieve much lower correlations, ranging from 0.657 (Cufflinks) to 0.775 (IsoEM). This is remarkable since, as there are only three time points, the improvement obtained by taking temporal information into account could be expected to be limited. Notice in particular that, while IsoEM and particularly Cufflinks sometimes give deterministic estimates in one replicate but not on the other (red points on the boundaries of the square in Figure 4A), this problem does not occur with DICEseq, presumably due to the stronger regularisation enforced by the temporal correlations.

To further explore the usefulness of DICEseq, we consider the confidence intervals reported by the various methods. Isoform quantification methods are often used as an initial step in kinetic analyses of individual transcripts; in order to reduce false positives, genes with unreliable isoform estimates (as determined by thresholding on the confidence intervals) are discarded. When quantifying isoforms in isolation, some genes are then discarded just because one of the time points have lower expression level. Therefore, we computed the number of transcripts that pass a frequently used threshold (95%CI < 0.3) for further analysis (Barras *et al.*, 2015). Figure 4B illustrates the results, showing that at all time points around 20% more genes are retained using a joint analysis, compared to methods that analyse data points in isolation.

To summarise, our results on a real yeast kinetic data set confirm that DICEseq yields significantly more reproducible and confident results than



**Fig. 5.** Analysis of circadian time series data. (a) The Pearson's correlation between the measurement of RNA-seq and microarray. (b) The proportion of genes whose 95% confidence interval < 0.3 in a certain number of time points (index on the external side of the circle).

existing state-of-the-art methods, highlighting the value of incorporating temporal information in the analysis of time series real data.

### 3.4 Circadian dynamics of alternative splicing

As a second real-data example, we turned to a recent data set investigating circadian control of gene expression in mouse. Due to the day-night oscillations, many biological processes, including gene expression, show circadian rhythms. Recently, Zhang *et al.* (Zhang *et al.*, 2014) systematically studied circadian gene expression on 12 mouse tissues using high-temporal resolution microarrays and RNA-seq, and found that 43% protein coding genes oscillate in at least one of the 12 tissues; here we focus on data from liver. The RNA-seq here has a comparably low time resolution, as eight time points were collected over a period of 48 hours; we expect therefore that the advantages of incorporating time information may be less pronounced in this scenario. In total, there are between 67M and 105M uniquely mapped reads in each experiment; on average of 8 time points, 50% of genes have all isoforms with RPK < 70; 75% of genes have all isoform with RPK < 400.

To assess the performance of the various methods, we used the microarray data set to validate the isoform estimates from RNA-seq. Unfortunately, only about one hundred microarray probes map to a unique annotated isoform (out of 30,534 annotated isoforms which map to at least one microarray probe); in other words, most microarray probes map to multiple isoforms within a gene. Thus, we used the estimated isoform proportions, together with the total numbers of reads mapped to each gene, to quantify the gene expression level (as FPKM), and then compared the resulting estimate from RNA-seq with the microarray measurement with Pearson's correlation coefficient. In Figure 5A, we see that the estimates obtained from RNA-seq using all methods have a high correlation with the direct measurements from the microarrays. Still, DICEseq shows a significantly improved correlation from 0.757 to 0.791 ( $p$ -value <  $1e-5$ , Fisher  $r$ -to- $z$  transform test); in particular, very low expressed isoforms

(outliers in the left end of the plot) show a much better quantification with DICEseq than with the other methods, probably due to the sharing of the temporal information, which is also evidenced by the medium third genes in Table 1.

We further measured 95% confidence intervals (CI) of all the 55,440 isoforms at the 8 time points, and quantified the fraction of isoform quantifications that pass the threshold  $95\%CI < 0.3$ . In Figure 5B we see that all Bayesian methods (Cufflinks, MISO and DICEseq) give confident estimates for between 50 and 60 % of isoforms at all time points. Once again, DICEseq estimates are more confident, thanks to the value of temporal information sharing at low coverages, even though the advantage is more modest in this data set.

To summarise, our results in this low-frequency RNA-seq time series data set show that even in this case DICEseq produces quantitatively better estimates of isoform ratios, even though the value of sharing temporal information is more limited here due to the weaker correlations between time points.

Table 1. Robust performance of DICEseq in lower or medium coverage. "All" means all annotated genes; "1/3 low" and "1/3 mid" respectively mean lowest and medium 1/3 genes in coverage. The scores are Pearson's correlation coefficients between two replicates (4tU-seq) or two techniques (circadian).

	IsoEM	Cufflinks	MISO	DICEseq
4tU-seq, all	0.851	0.830	0.848	0.896
4tU-seq, 1/3 low	0.775	0.657	0.757	0.860
circadian, all	0.712	0.700	0.757	0.791
circadian, 1/3 mid	0.336	0.296	0.408	0.513

## 4 Discussion

The advent of RNA-seq technologies has revolutionised the study of mRNA splicing, and provided a powerful stimulus for the development of computational biology methods (Katz *et al.*, 2010; Trapnell *et al.*, 2010; Nicolae *et al.*, 2011; Glaus *et al.*, 2012). Recent years have seen a more wide-spread use of RNA-seq technology for the analysis of dynamical biological processes, resulting in a marked increase of biological studies adopting RNA-seq within a time series experimental design. In this article, we presented DICEseq, the first method to jointly estimate the dynamics of the splicing isoform proportions from time series RNA-seq data. A comparison of DICEseq to a selection of popular state-of-the-art methods shows that DICEseq has excellent accuracy and good computational performance; in particular, DICEseq can effectively pool information across time points to improve isoform quantification at low coverages, giving more accurate and confident predictions. Our analysis also points to the importance of coverage versus temporal sampling trade-offs in designing dynamic RNA-seq experiments; while our analysis focussed on time series experiments, we expect similar considerations to hold for other structured designs, such as dose response experiments. In this light, the use of methods which can capture structural information, such as DICEseq, may lead to a rethink of biological experimental designs for a broad class of experiments. Our application to two diverse biological data sets shows that DICEseq can be an effective tool on real biological investigations, leading to improved performance and more reproducible results.

Methodologically, DICEseq builds on a fertile line of research using GPs to model transcriptional dynamics. GPs have been used to study the dynamical behaviour of gene expression in various contexts, from transcriptional regulation (Lawrence *et al.*, 2006) to identifying the time intervals of differential expression with time series microarray data (Stegle *et al.*, 2010). More recently, Äijö *et al* used a latent GP with negative

binomial observation noise to study the profiles of gene expression during Th17 cell differentiation with time course RNA-seq (Äijö *et al.*, 2014). To our knowledge, this is the first time GPs have been proposed within the context of isoform estimation.

While we believe DICEseq offers a valuable new tool for the analysis of dynamic RNA-seq data, it also opens several novel lines of investigation. Firstly, the Gaussian process prior, which is based on a general regression, could be extended to more general dynamic splicing modeling, e.g., a first-order linear dynamic system for RNA splicing kinetics, and an oscillatory system for circadian or cell-cycle studies. All of these could be incorporated in a straightforward way as parametric mean functions in a GP framework, however it would also be of interest to explicitly model the noise correlations they induce. DICEseq could be useful in elucidating RNA processing from biotin labelled RNA-seq, as attempted e.g. in de Pretis *et al* (de Pretis *et al.*, 2015). More generally, DICEseq could provide a flexible Bayesian framework for explaining RNA-seq data from other observations, and aid studies attempting to link splicing with other genetic and epigenetic factors.

## Acknowledgements

We would like to thank Prof. Jean D Beggs, Dr. Sander Granneman, Dr. Jane E A Reid, Dr. J David Barrass, Dr Edward Wallace and Dr. Yichuan Zhang for fruitful discussions. G.S. acknowledges support from the European Research Council under grant MLCS306999. Y.H. is supported by the University of Edinburgh through a Principal Career Development scholarship.

## References

- Äijö, T., Butty, V., Chen, Z., Salo, V., Tripathi, S., Burge, C. B., Lahesmaa, R., and Lähdesmäki, H. (2014). Methods for time series analysis of RNA-seq data with application to human Th17 cell differentiation. *Bioinformatics*, **30**(12), i113–i120.
- Bar-Joseph, Z., Gitter, A., and Simon, I. (2012). Studying and modelling dynamic biological processes using time-series gene expression data. *Nature Reviews Genetics*, **13**(8), 552–564.
- Barrass, J. D., Reid, J. E., Huang, Y., Hector, R. D., Sanguinetti, G., Beggs, J. D., and Granneman, S. (2015). Transcriptome-wide RNA processing kinetics revealed using extremely short 4tU labeling. *Genome Biology*, **16**(1), 1–17.
- Blencowe, B. J. (2006). Alternative splicing: new insights from global analyses. *Cell*, **126**(1), 37–47.
- de Pretis, S., Kress, T., Morelli, M. J., Melloni, G. E., Riva, L., Amati, B., and Pelizzola, M. (2015). INSPEC: A Computational Tool to Infer mRNA Synthesis, Processing and Degradation Dynamics from RNA and 4sU-seq Time Course Experiments. *Bioinformatics*, page btv288.
- Diedenhofen, B. and Musch, J. (2015). Cocor: A comprehensive solution for the statistical comparison of correlations. *PloS one*, **10**(4), e0121945.
- Eser, P., Wachutka, L., Maier, K. C., Demel, C., Boroni, M., Iyer, S., Cramer, P., and Gagneur, J. (2016). Determinants of RNA metabolism in the Schizosaccharomyces pombe genome. *Molecular Systems Biology*, **12**(2).
- Fuchs, G., Voichek, Y., Benjamin, S., Gilad, S., Amit, I., and Oren, M. (2014). 4sUDRB-seq: measuring genomewide transcriptional elongation rates and initiation frequencies within cells. *Genome Biology*, **15**(5), R69.
- Geweke, J. (1991). *Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments*, volume 196. Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN, USA.
- Glaus, P., Honkela, A., and Rattray, M. (2012). Identifying differentially expressed transcripts from RNA-seq data with biological variation. *Bioinformatics*, **28**(13), 1721–1728.
- Graveley, B. R. (2001). Alternative splicing: increasing diversity in the proteomic world. *TRENDS in Genetics*, **17**(2), 100–107.
- Honkela, A., Peltonen, J., Topa, H., Charapitsa, I., Matarese, F., Grote, K., Stunnenberg, H. G., Reid, G., Lawrence, N. D., and Rattray, M. (2015). Genome-wide modeling of transcription kinetics reveals patterns of RNA production delays. *Proceedings of the National Academy of Sciences*, **112**(42), 13115–13120.
- Kanitz, A., Gypas, F., Gruber, A. J., Gruber, A. R., Martin, G., and Zavolan, M. (2015). Comparative assessment of methods for the computational inference of transcript isoform abundance from RNA-seq data. *Genome Biology*, **16**(1), 1–26.



- Katz, Y., Wang, E. T., Airoidi, E. M., and Burge, C. B. (2010). Analysis and design of RNA sequencing experiments for identifying isoform regulation. *Nature Methods*, **7**(12), 1009–1015.
- Kim, D., Langmead, B., and Salzberg, S. L. (2015). HISAT: a fast spliced aligner with low memory requirements. *Nature Methods*, **12**(4), 357–360.
- Lawrence, N. D., Sanguinetti, G., and Rattray, M. (2006). Modelling transcriptional regulation using Gaussian processes. In *Advances in Neural Information Processing Systems*, pages 785–792.
- Nicolae, M., Mangul, S., Mandoiu, I. I., and Zelikovsky, A. (2011). Estimation of alternative splicing isoform frequencies from RNA-Seq data. *Algorithms for Molecular Biology*, **6**(1), 9.
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian processes for machine learning*. MIT Press, Cambridge, MA, USA.
- Roberts, A., Trapnell, C., Donaghey, J., Rinn, J. L., Pachter, L., *et al.* (2011). Improving RNA-Seq expression estimates by correcting for fragment bias. *Genome Biology*, **12**(3), R22.
- Scotti, M. M. and Swanson, M. S. (2016). RNA mis-splicing in disease. *Nature Reviews Genetics*, **17**(1), 19–32.
- Stegle, O., Denby, K. J., Cooke, E. J., Wild, D. L., Ghahramani, Z., and Borgwardt, K. M. (2010). A robust Bayesian two-sample test for detecting intervals of differential gene expression in microarray time series. *Journal of Computational Biology*, **17**(3), 355–367.
- Sturgill, D., Malone, J. H., Sun, X., Smith, H. E., Rabinow, L., Samson, M.-L., and Oliver, B. (2013). Design of RNA splicing analysis null models for post hoc filtering of Drosophila head RNA-Seq data with the splicing analysis kit (Spanki). *BMC Bioinformatics*, **14**(1), 320.
- Trapnell, C., Williams, B. A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M. J., Salzberg, S. L., Wold, B. J., and Pachter, L. (2010). Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnology*, **28**(5), 511–515.
- Tuomela, S., Salo, V., Tripathi, S. K., Chen, Z., Laurila, K., Gupta, B., Äijö, T., Oikari, L., Stockinger, B., Lähdesmäki, H., *et al.* (2012). Identification of early gene expression changes during human Th17 cell differentiation. *Blood*, **119**(23), e151–e160.
- Veloso, A., Kirkconnell, K. S., Magnuson, B., Biewen, B., Paulsen, M. T., Wilson, T. E., and Ljungman, M. (2014). Rate of elongation by RNA polymerase II is associated with specific gene features and epigenetic modifications. *Genome Research*, **24**(6), 896–905.
- Wang, E. T., Sandberg, R., Luo, S., Khrebtkova, I., Zhang, L., Mayr, C., Kingsmore, S. F., Schroth, G. P., and Burge, C. B. (2008). Alternative isoform regulation in human tissue transcriptomes. *Nature*, **456**(7221), 470–476.
- Wang, Z., Gerstein, M., and Snyder, M. (2009). RNA-Seq: a revolutionary tool for transcriptomics. *Nature Reviews Genetics*, **10**(1), 57–63.
- Windhager, L., Bonfert, T., Burger, K., Ruzsics, Z., Krebs, S., Kaufmann, S., Malterer, G., L'Hernault, A., Schilhabel, M., Schreiber, S., *et al.* (2012). Ultrashort and progressive 4sU-tagging reveals key characteristics of RNA processing at nucleotide resolution. *Genome Research*, **22**(10), 2031–2042.
- Zhang, R., Lahens, N. F., Ballance, H. I., Hughes, M. E., and Hogenesch, J. B. (2014). A circadian gene expression atlas in mammals: Implications for biology and medicine. *Proceedings of the National Academy of Sciences*, **111**(45), 16219–16224.