

# Leveraging transcript quantification for fast computation of alternative splicing profiles

Gael P. Alamancos<sup>1,\*</sup>, Amadís Pagès<sup>1,2,\*</sup>, Juan L. Trincado<sup>1</sup>, Nicolás Bellora<sup>3</sup>, Eduardo Eyras<sup>1,4,5</sup>

<sup>1</sup>Universitat Pompeu Fabra, E08003, Barcelona, Spain

<sup>2</sup>Centre for Genomic Regulation, E08003, Barcelona, Spain

<sup>3</sup>INIBIOMA, CONICET-UNComahue, Bariloche, Río Negro, Argentina

<sup>4</sup>Catalan Institution for Research and Advanced Studies, E08010 Barcelona, Spain

<sup>5</sup>corresponding author: [eduardo.eyras@upf.edu](mailto:eduardo.eyras@upf.edu)

\*These authors contributed equally

Running title: Fast computation of alternative splicing profiles

Keywords: RNA-Seq, splicing, splicing event

## Abstract

Alternative splicing plays an essential role in many cellular processes and bears major relevance in the understanding of multiple diseases, including cancer. High-throughput RNA sequencing allows genome-wide analyses of splicing across multiple conditions. However, the increasing number of available datasets represents a major challenge in terms of computation time and storage requirements. We describe SUPPA, a computational tool to calculate relative inclusion values of alternative splicing events, exploiting fast transcript quantification. SUPPA accuracy is comparable and sometimes superior to standard methods using simulated as well as real RNA sequencing data compared to experimentally validated events. We assess the variability in terms of the choice of annotation and provide evidence that using complete transcripts rather than more transcripts per gene provides better estimates. Moreover, SUPPA coupled with *de novo* transcript reconstruction methods does not achieve accuracies as high as using quantification of known transcripts, but remains comparable to existing methods. Finally, we show that SUPPA is more than 1000 times faster than standard methods. Coupled with fast transcript quantification, SUPPA provides inclusion values at a much higher speed than existing methods without compromising accuracy, thereby facilitating the systematic splicing analysis of large datasets with limited computational resources. The software is implemented in Python 2.7 and is available under the MIT license at <https://bitbucket.org/regulatorygenomicsupf/suppa>.

**Contact:** [eduardo.eyras@upf.edu](mailto:eduardo.eyras@upf.edu)

**Supplementary Information:** available at

[https://bitbucket.org/regulatorygenomicsupf/suppa/downloads/Supplementary\\_Data.zip](https://bitbucket.org/regulatorygenomicsupf/suppa/downloads/Supplementary_Data.zip)

## Introduction

Alternative splicing plays an important role in many cellular processes and bears major relevance in the understanding of multiple diseases, including cancer (David & Manley 2010, Ward & Cooper 2010). Numerous genome wide surveys have facilitated the description of the alternative splicing patterns under multiple cellular conditions and disease states. These studies are generally based on the measurement of local variations in the patterns of splicing, encoded as events, and have carried out using microarrays (Thorsen et al. 2008, Lapuk et al. 2010, Misquitta-Ali et al. 2011), RT-PCR platforms (Klinck et al. 2008, Simpson et al. 2008), or RNA sequencing (Pan et al. 2008, Wang et al. 2008). The description of alternative splicing in terms of events facilitates their experimental validation using PCR methods and the characterization of regulatory mechanisms using sequence analysis and biochemical approaches (Bechara et al. 2013, Raj et al. 2014); and they provide a valuable description for predictive and therapeutic strategies (Xiong et al. 2014, Hua et al. 2015). Events are generally defined as local variations of the exon-intron structure that can take two possible configurations, and are characterized by an inclusion level, also termed PSI or  $\Psi$ , which measures the fraction of mRNAs expressed from the gene that contain an specific form of the event (Venables et al. 2008, Wang et al. 2008). In terms of sequencing reads,  $\Psi$  is usually defined as the ratio of the density of inclusion reads to the sum of the densities of inclusion and exclusion reads (Wang et al. 2008, Shen et al. 2012). Initial methods to estimate  $\Psi$  values were based on reads from junction, exons or both (Pan et al. 2008, Wang et al. 2008, Sultan et al. 2008). Later methods were developed that take into account the uncertainty of quantification from single experiments (Katz et al. 2010), the comparison of two conditions (Katz et al. 2010, Griffith et al. 2010, Shen et al. 2012, Wu et al. 2011, Shi et al. 2013), as well as multiple replicates per condition (Shen et al. 2012, Brooks et al. 2011, Singh et al. 2011, Hu et al. 2013) and paired-replicates (Shen et al. 2014).

Current tools to process RNA sequencing data to study alternative splicing events can take more than a day to analyze a single sample and often require excessive storage, so they are not competitive to be applied systematically to large data sets, unless access to large computational resources is granted. In particular, methods for estimating  $\Psi$  values generally involve the mapping of reads to the genome or to a library of known exon-exon junctions, both of which require considerable time and storage. Additionally, accuracy is often achieved

at the cost of computing time. All this represents a major obstacle for the analysis of large datasets, and in particular, for the re-analysis of public data and updates with new annotations or assembly versions. More importantly, these analyses remain unfeasible at small labs with limited computational resources. On the other hand, recent developments in the quantification of known transcripts have shown that considerable accuracy can be achieved at high speed (Li et al. 2011, Roberts et al. 2013, Patro et al. 2014, Zhang et al. 2014). This raises the question of whether fast transcript abundance computation could be used to obtain accurate estimates of  $\Psi$  values for local alternative splicing events genome wide.

In this article we describe SUPPA, a computational tool to leverage fast transcript quantification for rapid estimation of  $\Psi$  values directly from the abundances of the transcripts defining each event. Using simulated data we show that  $\Psi$  values estimated by SUPPA, coupled to Sailfish or RSEM transcript quantification, are closer to the ground-truth than two standard methods, MATS and MISO. Additionally, using an experimentally validated set and matched RNA-Seq data we show that SUPPA achieves slightly superior or comparable accuracy compared with MATS and MISO. We further assess the variability in terms of the choice of annotation and provide evidence that using complete transcripts rather than more transcripts per gene in the annotation provide better estimates. Moreover, we show that SUPPA coupled with *de novo* transcript reconstruction methods does not achieve accuracies as high as using the quantification of known transcripts, but remains comparable to existing methods. Finally, speed benchmarking provides evidence that SUPPA can obtain  $\Psi$  values at a much higher speed than existing methods. We argue that coupled to a fast transcript quantification method, SUPPA provides a fast and accurate approach to systematic splicing analysis. SUPPA facilitates the accurate splicing analysis of large datasets, making possible for labs with limited computational resources to exploit data from large genomics projects and contribute to the understanding of the role of alternative splicing in cell biology and disease.

## Results

### SUPPA

SUPPA provides an effective and easy-to-use software to calculate the inclusion levels ( $\Psi$ ) of alternative splicing events exploiting transcript quantification (Figure 1A). An alternative splicing event is a local summary representation of the exon-intron structure from the transcripts that cover a given genic region, and is generally represented as a binary form, although more complex variations may happen. Accordingly, an event can be characterized in terms of the sets of transcripts that describe either form of the event, which can be denoted as  $F_1$  and  $F_2$ . For instance, for an exon-skipping event,  $F_1$  represents the transcripts that include the exon, whereas  $F_2$  represents the transcripts that skip the exon. The inclusion value ( $\Psi$ ) of an event is defined as the ratio of the abundance of transcripts that include one form of the event,  $F_1$ , over the abundance of the transcripts that contain either form of the event,  $F_1 \cup F_2$  (Venables et al. 2008, Wang et al. 2008, Katz et al. 2010, Shen et al. 2012). Given the abundances for all transcripts isoforms, assumed without loss of generality to be given in transcript per million units (TPM) (Li et al. 2010), which we denote as  $TPM_k$ , SUPPA then calculates  $\Psi$  for an event as follows:

$$\Psi = \frac{\sum_{k \in F_1} TPM_k}{\sum_{j \in F_1 \cup F_2} TPM_j} \quad (1)$$

SUPPA provides the identifiers for the transcripts that describe either form of the event, which in combination with the transcript quantification is used to obtain the  $\Psi$  values using formula (1) (Figure 1A). SUPPA is agnostic of the actual methodology for quantifying transcripts and can read the quantification from multiple experiments from a single input. SUPPA generates different alternative splicing events types from an input annotation file in GTF format: exon skipping (SE), alternative 5' and 3' splice-sites (A5/A3), mutually exclusive exons (MX), intron retention (RI), and alternative first and last exons (AF/AL) (Figure 1B). The  $\Psi$  value for an event is calculated with respect to one of the two forms of the event (Figure 1B). Further details and options of the software are given at <https://bitbucket.org/regulatorygenomicsupf/suppa/>.

## Accuracy analysis with simulated data

Transcript abundances and corresponding paired-end reads were simulated using FluxSimulator (Griebel et al. 2012) with the RefSeq annotation as reference (Methods) (Supplementary Table 2). The reference set for accuracy analysis was built using events in genes with only two alternative transcripts in the RefSeq annotation that did not overlap any other events. In these cases, the  $\Psi$  of the event is identical to the relative abundance of one of the two transcripts. The ground-truth  $\Psi$  values were then defined to be the relative abundances of the transcripts isoforms in these genes, where the transcript abundances were taken to be the simulated abundances. Simulated RNA-Seq reads were mapped to the genome and used to calculate  $\Psi_{\text{MISO}}$  and  $\Psi_{\text{MATS}}$  values with MISO (Katz et al. 2010) and MATS (Shen et al. 2011), respectively (Methods). The same simulated reads were also used to quantify transcript abundances with Sailfish (Patro et al. 2014) and RSEM (Li et al. 2011), and  $\Psi_{\text{Sailfish}}$  and  $\Psi_{\text{RSEM}}$  values were then calculated with SUPPA (Methods). Only genes with total transcripts per million (TPM) abundance, calculated as the sum of the TPM of its transcripts, greater than 1 were considered. This resulted in a set of 144 events (Supplementary Data 1). Comparing the four sets of estimated  $\Psi$  values with the ground-truth, the  $\Psi_{\text{Sailfish}}$  and  $\Psi_{\text{RSEM}}$  values calculated with SUPPA show the highest correlations (Table 1) (Figure 2A). Moreover, calculating how different the estimated  $\Psi$  values are from the ground-truth, SUPPA  $\Psi$  values ( $\Psi_{\text{Sailfish}}$  and  $\Psi_{\text{RSEM}}$ ) show the closest behaviour, followed by MISO and MATS, which behave similarly (Figure 2B).

## Accuracy analysis with experimentally validated events

To further validate the calculation of  $\Psi$  values with SUPPA, we used a set of 163 alternative splicing events validated by RT-PCR in MDA-MB-231 cells under two conditions: with overexpression of the splicing factor ESRP1 (ESRP1) and with an empty vector (EV) (Shen et al. 2012). We used the RNA-Seq data obtained from the same samples (Shen et al. 2012) to predict the  $\Psi$  values as before. From both RNA-Seq datasets we quantified the RefSeq transcripts using Sailfish and RSEM, and calculated the SUPPA  $\Psi_{\text{Sailfish}}$  and  $\Psi_{\text{RSEM}}$  values. RNA-Seq reads were mapped to the genome to run MISO and MATS to obtain the corresponding  $\Psi$  values (Methods). From the 163 validated events, we finally compared those

60 that were present in the RefSeq annotation and for which we had  $\Psi$  values for all methods (Supplementary Data 2). Sailfish+SUPPA and RSEM+SUPPA show an overall slightly better correlation than the other methods for the ESRP1 sample, whereas RSEM+SUPPA and MATS show the best correlations for the EV sample (Table 1) (Figure 3A). Although RSEM+SUPPA shows the highest correlations in all cases, Sailfish+SUPPA correlations are comparable to the rest. Calculating the absolute difference between the estimated and the experimental  $\Psi$  values for each event, we observe that SUPPA, either combined with Sailfish or RSEM, is more accurate than MISO and MATS (Figure 3B). Performing the same analysis using the Ensembl annotation, comparing a total of 91 events common to all approaches, we observe a general decrease of accuracy in all methods (Table 1) (Supplementary Figure 1).

### **Variability associated to replicates and annotation choice**

To study how the choice of annotation may impact the accuracy of  $\Psi$  estimation, we obtained RNA from two biological replicates for the cytosolic fractions of MCF7 and MCF10 cells and performed sequencing using standard protocols. Correlation between replicates of the SUPPA  $\Psi$  values, using quantification with Sailfish on the RefSeq annotation, is high in all comparisons (Person R ~0.86-0.89) (Supplementary Figure 2). Furthermore, restricting this analysis to genes with TPM>1, calculated as the sum of the TPMs from all transcripts in each gene, the correlation between replicates increases (Pearson R~0.95-0.97) (Supplementary Figure 2). We then compared the results obtained using SUPPA with the quantifications on the RefSeq and Ensembl transcripts. SUPPA  $\Psi$  values were calculated using both replicas of the cytosolic MCF7 RNA-Seq data (similar results were observed for MCF10, data not shown). The comparisons were performed using the 9301 (MCF7, replica 1) and 9287 (MCF7, replica 2) events that were found in both annotations, and not overlapping with other events from the same replica. We observe variability in the estimation of  $\Psi$  between annotations that does not depend on the difference in the number of transcripts used for  $\Psi$  calculation (Figure 4A). Similarly, this variability is also independent of the difference in the number of transcripts annotated in the gene in which the event is contained (Supplementary Figure 3). Moreover, the disparity in  $\Psi$  estimates is also independent of the mean expression of the gene in which the event is contained (Figure 4B). On the other hand, the dispersion of  $\Psi$  estimates comparing replicas and using the same annotation decreases with the mean

expression of the gene (Figure 4C), which at low expression it is comparable to the dispersion for  $\Psi$  estimates as a result of differences in annotation (Figures 4A and 4B).

### **Annotation-free estimation**

The previous analyses suggest that incomplete annotations may lead to inaccurate transcript quantification, which will have in turn a negative impact on the  $\Psi$  estimates by SUPPA. Methods for *de novo* transcript reconstruction facilitate the discovery of new transcripts missing from the annotation and the completion of existing ones from RNA-Seq reads (Trapnell et al. 2010, Li et al. 2011b, Li et al. 2011c, Li et al. 2012, Mezlini et al. 2012, Behr et al. 2013, Tomescu et al. 2013, Rossell et al. 2014, Maretty et al. 2014). As these methods produce an annotation of transcripts and their corresponding abundances, their output can be used with SUPPA to calculate alternative splicing events and their  $\Psi$  values. They thus provide an opportunity to assess whether a *de novo* prediction of transcripts structures and subsequent quantification from RNA-Seq data may lead to more accurate  $\Psi$  values than using a fixed annotation. To test this, we run Cufflinks with the *de novo* options with RNA-Seq data from the ESRP1 and EV samples (Methods). Using the resulting annotation, we calculated all possible alternative splicing events and their contributing transcripts with SUPPA. Similarly, we calculate the  $\Psi$  with MATS and MISO using the same reads mapped to the genome, this time guided by the Cufflinks annotation. We then compared the  $\Psi$  values obtained for the events in common with the experimentally validated set (Shen et al. 2012): 82 for ESRP1 and 47 for EV (Supplementary Data 3). We observe that for all approaches the correlation of  $\Psi$  values decreases (Table 2). The  $\Psi_{\text{Cufflinks}}$  values obtained with SUPPA (Figure 5) (Table 2) are comparable to the values obtained using the Ensembl annotation (Table 1). Moreover, we recalculated the transcript quantification using Sailfish on the Cufflinks annotations, but found no improvement (Table 2).

### **Speed benchmarking**

The time needed by each methodology to obtain the  $\Psi$  values from a FastQ file depends on multiple different steps. To make a comparative assessment of computation times we broke



down the benchmarking into three different tasks, equivalent to the three necessary steps for the SUPPA analysis. The first step involves the calculation of alternative splicing events from an annotation file, which only needs to be carried once for a give annotation. To calculate 66577 alternative splicing events from the Ensembl 75 annotation (37494 genes, 135521 transcripts), SUPPA *generateEvents* took 20 minutes, whereas to calculate 16714 alternative splicing events from the RefSeq annotation (25937 genes, 48566 transcripts), it took 3 minutes.

The second step consists in the assignment of reads to transcripts and/or genomic positions. For the purpose of speed benchmarking of read-assignment to transcripts, although transcript abundance estimation includes extra computation steps, we considered the transcript quantification by Sailfish to be approximately equivalent to the read mapping to a reference genome. To perform this speed comparison we used the synthetic data (45 millions of paired-end reads) and both (ESRP1 and EV) RNA-Seq samples from the MDA-MB-231 cells pooled together (256 millions of single-end reads), and used STAR (Dobin et al. 2012) and TopHat as a comparison. Sailfish and STAR are the fastest to assign reads to their likely molecular sources, compared to TopHat and RSEM (Figure 6A).

The third and final step is the  $\Psi$  calculation from either transcript quantification (SUPPA) or from the mapped reads (MISO and MATS). SUPPA *psiPerEvent* operation took less than a minute to produce an output size of 1Mb for 16714 events and was >1000 times faster than MISO and MATS on the same datasets (Figure 6B). In summary, the total time from the raw reads in FastQ format to the  $\Psi$  values for Sailfish + SUPPA against the RefSeq annotation-derived events took 214 (~ 3,5 mins) and 4022 (~ 1h) seconds for the synthetic and the MDA-MB-231 samples, respectively. We conclude that when used in conjunction with Sailfish, SUPPA is much faster than MISO and MATS, even if an ultra-fast aligner such as STAR (Dobin et al. 2012) is used for read mapping to the genome.

## Discussion

We have described SUPPA, a tool to calculate alternative splicing events from a given annotation and to estimate their  $\Psi$  values from the quantification of the transcripts that define the events. Using synthetic and experimental data, we have shown that SUPPA accuracy is

generally comparable to and sometimes higher than other frequently used methods. Importantly, SUPPA can obtain  $\Psi$  values at a much higher speed without compromising accuracy. Moreover, SUPPA needs very little configuration, requires a small number of command lines for preprocessing and running and has no dependencies on Python libraries.

Although RNA-Seq data presents a number of systematic biases that need correction for accurate transcript quantification (Hansen et al. 2010, Li et al. 2010, Roberts et al. 2011), we did not observe differences in the accuracy of SUPPA when comparing corrected or uncorrected transcript quantification with Sailfish (data not shown). In fact, previous reports have already indicated that bias correction in RNA-Seq data does not influence the estimation of  $\Psi$  values (Shen et al. 2012, Zhao et al. 2013). On the other hand, we did observe that there is variability in the estimation of  $\Psi$  values associated to the choice of annotation. In the benchmarking using experimental data, using Ensembl annotation provides slightly worse accuracy than using RefSeq annotation, and this behaviour is consistent amongst all the tested methods. Interestingly, the observed variability between annotations does not depend on the difference in the number of transcripts per gene, on the number of transcripts used to describe the events, or on the expression of the gene in which the event is contained. On the other hand, the observed variability is comparable to the expected variability for lowly expressed genes between biological replicates. Such variability is in fact also frequently observed in transcript quantification methods (Patro et al. 2014, Maretty et al. 2014). It should be noted that RefSeq annotation includes less transcripts per gene than Ensembl, but these transcripts are mostly full-length mRNAs. In particular, RefSeq transcripts generally include complete untranslated regions, which generally hold a large contribution of the reads coming from a transcript, whereas a large proportion of Ensembl transcripts may be incomplete. These facts, together with our results, suggest that the completeness of the transcript structures, rather than the number of transcripts in genes, is determinant for an accurate estimate of transcript abundances, and consequently, for the correct estimate of event  $\Psi$  values with SUPPA.

Although SUPPA at the moment SUPPA generates the most common types of events, its model can be potentially expanded to more complex events, possibly involving more than two possible conformations. However, these complex events may not always be easy to test experimentally. On the other hand, the complexity may not always have to do with the number of possible conformations, but rather with a binary change that cannot be easily

described in terms of just one or two exon boundary changes, as described recently for the gene QKI in lung adenocarcinomas (Sebestyen et al. 2015). We argue that a large proportion of the relevant splicing variation can be encapsulated with the binary events described by SUPPA and that more complex variations may be better described using transcript isoform changes (Sebestyen et al. 2015). Although SUPPA is limited to the splicing events available in the gene annotation, events can be expanded with novel transcript variants obtained by other means, like *de novo* transcript reconstruction and quantification methods. In this case we observed accuracies similar to the tests performed with the Ensembl annotations but lower than when using the RefSeq annotations. Moreover, performing quantification on the reconstructed transcripts using a different method does not improve the accuracy, indicating there is still a limitation on how well we can recover the right exon-intron structures *de novo* from RNA-Seq.

As transcript reconstruction and quantification methods improve in accuracy and methods for RNA sequencing increase their efficiency and reliability, our knowledge of the census of RNA molecules in cells will keep on progressing. Although single molecule sequencing methods may eventually lead to the abandonment of transcript reconstruction methods, they are still costly and error prone, and quantification still relies on short read sequencing. Transcript quantification methods will therefore continue to be an essential component in the description of the abundance of RNA molecules in cells. As fast reliable methods still depend on the annotation, future efforts may perhaps focus on improving transcript annotations under multiple conditions. In parallel to these advances, the local description of alternative splicing in terms of events will remain a valuable description RNA variability in genes in the context of studies of RNA regulation (Bechara et al. 2013, Raj et al. 2014) and of predictive and therapeutic strategies (Xiong et al. 2014, Hua et al. 2015).

In summary, when coupled to a fast transcript quantification method, SUPPA outperforms other methods in speed without compromising the accuracy. This is of special relevance when analyzing large amount of samples. Accordingly, SUPPA facilitates the systematic analyses of alternative splicing in the context of large-scale projects using limited computational resources. We conclude that SUPPA provides a method to leverage fast transcript quantification for efficient and accurate alternative splicing analysis for a large number of samples.

## Methods

### Alternative splicing events

The Ensembl annotation (Release 75) (Flicek et al. 2014) and the RefSeq annotation (NM\_ and NR\_ transcripts) (Pruitt et al. 2014) (assembly hg19) were downloaded in GTF format from the Ensembl FTP server and the UCSC genome table browser, respectively. All annotations on chromosomes other than autosomes or sex chromosomes were removed. In total, 37,494 genes and 135,521 transcripts were obtained for the Ensembl annotation, while 25,937 genes and 48,566 transcripts were obtained for the RefSeq annotation. We applied SUPPA to each annotation to obtain 16714 and 66577 events from RefSeq and Ensembl, respectively, including exon skipping (SE), alternative 5' and 3' splice-sites (A5/A3), mutually exclusive exons (MX), and intron retention (RI) events (Supplementary Table 1). Alternative first (AF) and last exons (AL) were not included in the analysis but can be also computed with SUPPA. Each event has a unique identifier that includes the gene symbol, the type of event, and the coordinates and strand that characterize the event:

*<gene\_id>;<event\_type>:<seqname>:<coordinates\_of\_the\_event>:strand*

where *gene\_id*, *seqname* and *strand* are obtained directly from the input annotation in GTF, *seqname* is the field 1 from the GTF file, generally the chromosome. The field *coordinates\_of\_the\_event* is defined by start and end coordinates that define the *event\_type* (SE, MX, A5, A3, RI, AF, AL).

### RNA sequencing data

A total of 45 million 2x50bp paired-end simulated reads were generated using FluxSimulator (Griebel et al. 2012) (parameter file described in Supplementary Table 2). RNA sequencing data from (Shen et al. 2012) was also used, corresponding to ESRP1-overexpression (ESRP1) and empty-vector (EV) experiments in MBA-MD-231 cells, available from the short read archive (SRA) under id SRX122589. Moreover, RNA sequencing was also performed in duplicate on cytosolic fractions of MCF7 and MCF10 cells using standard protocols (Supplementary Material), available at SRA under id SRP045592.

## Read mapping and PSI quantification

Read mapping to the genome was performed with the MATS pipeline (Shen et al. 2012), which uses TopHat (Trapnell et al. 2009) and an input annotation to map the reads. Reads mapping to *de novo* splice junctions were allowed, and those reads mapping to more than one genomic position were filtered out. For benchmarking, the same annotation used for transcript quantification was also used for read mapping to the genome in each of the comparisons (RefSeq, Ensembl or *de novo* Cufflinks). The mapping pipeline was run on simulated and real RNA-Seq reads. Mapped reads for each of the datasets were used with MATS, to obtain  $\Psi_{\text{MATS}}$  values for the different alternative splicing events (Supplementary Table 3). Similarly, mapped reads in SAM format were converted to BAM format and then sorted with samtools (Li et al. 2009) and analysed with MISO (Katz et al. 2010) to calculate the  $\Psi_{\text{MISO}}$  values for each of the datasets (Supplementary Table 4).

Sailfish (Patro et al. 2014) and RSEM (Li et al. 2011) were used to quantify all transcripts in the Ensembl and RefSeq annotations using the simulated and the real RNA-Seq datasets. The FASTA sequences of the transcripts corresponding to the same annotation as the GTF described earlier, were downloaded and used to generate the Sailfish index, selecting a *k*-mer size of 31 to minimize the number of reads assigned to multiple transcripts. Sailfish was then run using the FASTQ files for each read set and uncorrected and corrected (for sequence composition bias and transcript length) TPMs were calculated (Patro et al. 2014). RSEM was run as described previously (Li et al. 2011). The *psiPerEvent* operation of SUPPA was used to calculate the  $\Psi_{\text{Sailfish}}$  and  $\Psi_{\text{RSEM}}$  values from the transcript quantifications obtained by Sailfish and RSEM, respectively, for the alternative splicing events generated before, using the simulated and real datasets. The number of events for which SUPPA estimated a  $\Psi_{\text{Sailfish}}$  or  $\Psi_{\text{RSEM}}$  values are given in Supplementary Tables 5 and 6. For the purpose of benchmarking, the PSI values obtained from SUPPA ( $\Psi_{\text{Sailfish}}$  and  $\Psi_{\text{RSEM}}$ ), from MATS ( $\Psi_{\text{MATS}}$ ) and from MISO ( $\Psi_{\text{MISO}}$ ) for those events identified by all methods in each of the experiment, were compared with the simulated or the experimental values. Details of the commands used to run the different analyses are provided in Supplementary Tables 7-10. Supplementary data files with the alternative splicing events used in each one of the comparisons tested can be found at

[https://bitbucket.org/regulatorygenomicsupf/suppa/downloads/Supplementary\\_Data.zip](https://bitbucket.org/regulatorygenomicsupf/suppa/downloads/Supplementary_Data.zip)

## Cufflinks analysis

The BAM files from 2 the MBA-MD-231 datasets were used to run Cufflinks (Trapnell et al. 2010) in order to generate and quantify transcriptome annotations *de novo*. The same read mapping as before was used. A total of 47211 transcripts were predicted and quantified for the ESRP1 dataset, whereas 37699 transcripts were predicted and quantified for the EV dataset. SUPPA *generateEvents* operation was then run on the GTF annotation generated by Cufflinks to calculate all the exon skipping events. This produced a total of 2566 and 2139 exon skipping events for the ESRP1 and EV datasets, respectively. Finally, SUPPA *psiPerEvent* operation was used to calculate the  $\Psi_{\text{Cufflinks}}$  values from the transcript quantification obtained by Cufflinks. For MISO and MATS, reads were mapped with the MATS pipeline using the Cufflinks annotation as input, and  $\Psi_{\text{MATS}}$  and  $\Psi_{\text{MISO}}$  were estimated as before. Additionally, Cufflinks reconstructed transcripts were used with SUPPA to quantify them from the same RNA-Seq data and to calculate  $\Psi$  values with SUPPA as before. The events common to all methods and coinciding with the experimentally validated ones were used for the benchmarking.

## Time benchmarking

All tools were run on the same node of an Oracle Grid Engine cluster, with 98Gb of RAM memory and 24 AMD Opteron (1.4 GHz) processors. All tools were run in multi-threaded mode when possible, but time reported is the actual cumulative time the process used across all CPUs.

## Acknowledgements

The authors wish to thank S. Mount, S. Janga, Y. Barash, M. Robinson for useful discussions and especially Y. Xing for useful discussions and for sharing sequencing and RT-PCR data. This work was supported by the Spanish Government [BIO2011-23920, CSD2009-00080], by the Sandra Ibarra Foundation for Cancer [FSI-2013] and partly by the Spanish National Institute of Bioinformatics (INB).

## References

Bechara EG, Sebestyén E, Bernardis I, Eyraş E, Valcárcel J. RBM5, 6, and 10 differentially regulate NUMB alternative splicing to control cancer cell proliferation. *Mol Cell*. 2013 Dec 12;52(5):720-33.

Behr J, Kahles A, Zhong Y, Sreedharan VT, Drewe P, Rátsch G (2013). MITIE: Simultaneous RNA-Seq- based transcript identification and quantification in multiple samples. *Bioinformatics* 2013, 29:2529–2538.

Brooks AN, Yang L, Duff MO, Hansen KD, Park JW, Dudoit S, Brenner SE, Graveley BR. (2011). Conservation of an RNA regulatory map between *Drosophila* and mammals. *Genome Res*. 21(2):193-202.

David CJ, Manley JL.(2010) Alternative pre-mRNA splicing regulation in cancer: pathways and programs unhinged. *Genes Dev*. 24(21):2343-64.

Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013 Jan 1;29(1):15-21.

Flicek P, et al. Ensembl 2014. *Nucleic Acids Res*. 2014 Jan;42(Database issue):D749-55.

Griebel T, Zacher B, Ribeca P, Raineri E, Lacroix V, Guigó R, Sammeth M. Modelling and simulating generic RNA-Seq experiments with the flux simulator. *Nucleic Acids Res*. 2012 Nov 1;40(20):10073-83.

Griffith M, Griffith OL, Mwenifumbo J, Goya R, Morrissy AS, Morin RD, Corbett R, Tang MJ, Hou YC, Pugh TJ, Robertson G, Chittaranjan S, Ally A, Asano JK, Chan SY, Li HI, McDonald H, Teague K, Zhao Y, Zeng T, Delaney A, Hirst M, Morin GB, Jones SJ, Tai IT, Marra MA. (2010). Alternative expression analysis by RNA sequencing. *Nat Methods* 7(10):843-7.

Hansen KD, Brenner SE, Dudoit S (2010) Biases in Illumina transcriptome sequencing

caused by random hexamer priming. *Nucleic Acids Res* 38(12):e131.

Hu Y, Huang Y, Du Y, Orellana CF, Singh D, Johnson AR, Monroy A, Kuan PF, Hammond SM, Makowski L, Randell SH, Chiang DY, Hayes DN, Jones C, Liu Y, Prins JF, Liu J. (2013). DiffSplice: the genome-wide detection of differential splicing events with RNA-seq. *Nucleic Acids Res.* 41(2):e39.

Hua Y, Liu YH, Sahashi K, Rigo F, Bennett CF, Krainer AR. Motor neuron cell-nonautonomous rescue of spinal muscular atrophy phenotypes in mild and severe transgenic mouse models. *Genes Dev.* 2015 Jan 12.

Katz Y, Wang ET, Airoidi EM, Burge CB. Analysis and design of RNA sequencing experiments for identifying isoform regulation. *Nat Methods.* 2010 Dec;7(12):1009-15.

Klinck, R., Bramard, A., Inkel, L., Dufresne-Martin, G., Gervais-Bird, J., Madden, R., Paquet, E.R., Koh, C., Venables, J.P., Prinos, P., Jilaveanu-Pelms, M., Wellinger, R., Rancourt, C., Chabot, B., Elela, S.A. (2008) Multiple alternative splicing markers for ovarian cancer. *Cancer Res.*, **68**, 657-63.

Lapuk, A., Marr, H., Jakkula, L., Pedro, H., Bhattacharya, S., Purdom, E., Hu, Z., Simpson, K., Pachter, L., Durinck, S., Wang, N., Parvin, B., Fontenay, G., Speed, T., Garbe, J., Stampfer, M., Bayandorian, H., Dorton, S., Clark, T.A., Schweitzer, A., Wyrobek, A., Feiler, H., Spellman, P., Conboy, J., Gray, J.W. (2010) Exon-level microarray analyses identify alternative splicing programs in breast cancer. *Mol. Cancer Res.*, **8**, 961-74.

Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R; 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and SAMtools. *Bioinformatics.* 2009 Aug 15;25(16):2078-9.

Li B, Ruotti V, Stewart RM, Thomson JA, Dewey CN. RNA-Seq gene expression estimation with read mapping uncertainty. *Bioinformatics.* 2010 Feb 15;26(4):493-500.

Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics.* 2011 Aug 4;12:323.



Li W, Feng J, Jiang T (2011b). IsoLasso: a LASSO regression approach to RNA-Seq based transcriptome assembly. *J Comput Biol* 2011, 18:1693–1707.

Li JJ, Jiang C-R, Brown JB, Huang H, Bickel PJ (2011c): Sparse linear modeling of next-generation mRNA sequencing (RNA-Seq) data for isoform discovery and abundance estimation. *Proc Nat Acad Sci USA* 2011, 108:19867–19872.

Li W, Jiang T: Transcriptome assembly and isoform expression level estimation from biased RNA-seq reads. *Bioinformatics* 2012, 28:2914–2921.

Maretty L, Sibbesen J, Krogh A. (2014). Bayesian transcriptome assembly. *Genome Biol.* 15(10):501.

Mezlini AM, Smith EJ, Fiume M, Buske O, Savich G, Shah S, Aparicion S, Chiang D, Goldenberg A, Brudno M: (2012). iReckon: Simultaneous isoform discovery and abundance estimation from RNA-seq data. *Genome Res* 2012, 23:519–529.

Misquitta-Ali CM, et al. (2011). Global profiling and molecular characterization of alternative splicing events misregulated in lung cancer. *Mol Cell Biol.* **31**(1):138-50.

Pan Q, Shai O, Lee LJ, Frey BJ, Blencowe BJ. Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing. *Nat Genet.* 2008 Dec;40(12):1413-5.

Patro R, Mount SM, Kingsford C. Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nat Biotechnol.* 2014 May;32(5):462-4.

Pruitt KD, Brown GR, Hiatt SM, Thibaud-Nissen F, Astashyn A, Ermolaeva O, Farrell CM, Hart J, Landrum MJ, McGarvey KM, Murphy MR, O'Leary NA, Pujar S, Rajput B, Rangwala SH, Riddick LD, Shkeda A, Sun H, Tamez P, Tully RE, Wallin C, Webb D, Weber J, Wu W, DiCuccio M, Kitts P, Maglott DR, Murphy TD, Ostell JM. RefSeq: an update on mammalian reference sequences. *Nucleic Acids Res.* 2014 Jan;42(Database issue):D756-63.

Raj B, Irimia M, Braunschweig U, Sterne-Weiler T, O'Hanlon D, Lin ZY, Chen GI, Easton LE, Ule J, Gingras AC, Eyraş E, Blencowe BJ. A global regulatory mechanism for activating an exon network required for neurogenesis. *Mol Cell*. 2014 Oct 2;56(1):90-103.

Roberts A, Trapnell C, Donaghey J, Rinn JL, Pachter L. Improving RNA-Seq expression estimates by correcting for fragment bias. *Genome Biol*. 2011;12(3):R22.

Roberts A, Pachter L. Streaming fragment assignment for real-time analysis of sequencing experiments. *Nat Methods*. 2013 Jan;10(1):71-3.

Rossell, D., Attolini, C. S. O., Kroiss, M., & Stöcker, A. (2014). Quantifying alternative splicing from paired-end RNA-sequencing data. *The annals of applied statistics*, 8(1), 309-330.

Sebestyén E, Zawisza M, Eyraş E. Detection of recurrent alternative splicing switches in tumor samples reveals novel signatures of cancer. *Nucleic Acids Res*. 2015 Jan 10. doi: 10.1093/nar/gku1392

Shen S, Park JW, Huang J, Dittmar KA, Lu ZX, Zhou Q, Carstens RP, Xing Y. MATS: a Bayesian framework for flexible detection of differential alternative splicing from RNA-Seq data. *Nucleic Acids Res*. 2012 Apr;40(8):e61.

Shen S, Park JW, Lu ZX, Lin L, Henry MD, Wu YN, Zhou Q, Xing Y. rMATS: Robust and flexible detection of differential alternative splicing from replicate RNA-Seq data. *Proc Natl Acad Sci U S A*. 2014 Dec 5. pii: 201419161.

Shi Y, Jiang H. rSeqDiff: detecting differential isoform expression from RNA-Seq data using hierarchical likelihood ratio test. *PLoS One*. 2013 Nov 18;8(11):e79448.

Simpson CG, Fuller J, Maronova M, Kalyna M, Davidson D, McNicol J, Barta A, Brown JW. (2008). Monitoring changes in alternative precursor messenger RNA splicing in multiple gene transcripts. *Plant J*. 53(6):1035-48.

Singh D, Orellana CF, Hu Y, Jones CD, Liu Y, Chiang DY, Liu J, Prins JF. (2011). FDM: a

graph-based statistical method to detect differential transcription using RNA-seq data. *Bioinformatics* 27(19):2633-40.

Sultan M, Schulz MH, Richard H, Magen A, Klingenhoff A, Scherf M, Seifert M, Borodina T, Soldatov A, Parkhomchuk D, Schmidt D, O'Keefe S, Haas S, Vingron M, Lehrach H, Yaspo ML. (2008). A global view of gene activity and alternative splicing by deep sequencing of the human transcriptome. *Science* 321(5891):956-60.

Thorsen, K., Sørensen, K.D., Brems-Eskildsen, A.S., Modin, C., Gaustadnes, M., Hein, A.M., Kruhøffer, M., Laurberg, S., Borre, M., Wang, K., Brunak, S., Krainer, A.R., Tørring, N., Dyrskjøt, L., Andersen, C.L., ørntoft, T.F. (2008) Alternative splicing in colon, bladder, and prostate cancer identified by exon array analysis. *Mol. Cell. Proteomics*, 7, 1214-24.

Tomescu AI, Kuosmanen A, Rizzi R, Mäkinen V, Veli M (2013). A novel min-cost flow method for estimating transcript expression with RNA-Seq. *BMC Bioinformatics* 14:S15.

Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics*. 2009 May 1;25(9):1105-11.

Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L. (2010). Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*. 28(5):511-5.

Venables JP, Klinck R, Bramard A, Inkel L, Dufresne-Martin G, Koh C, Gervais-Bird J, Lapointe E, Froehlich U, Durand M, Gendron D, Brosseau JP, Thibault P, Lucier JF, Tremblay K, Prinos P, Wellinger RJ, Chabot B, Rancourt C, Elela SA. Identification of alternative splicing markers for breast cancer. *Cancer Res*. 2008 Nov 15;68(22):9525-31.

Wang ET, Sandberg R, Luo S, Khrebtkova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP, Burge CB. Alternative isoform regulation in human tissue transcriptomes. *Nature*. 2008 Nov 27;456(7221):470-6.

Ward, A. and Cooper, T. (2010) The pathobiology of splicing. *J. Pathol.*, 220, 152–163.

Wu J, Akerman M, Sun S, McCombie WR, Krainer AR, Zhang MQ. SpliceTrap: a method to quantify alternative splicing under single cellular conditions. *Bioinformatics*. 2011 Nov 1;27(21):3010-6.

Xiong HY, Alipanahi B, Lee LJ, Bretschneider H, Merico D, Yuen RK, Hua Y, Gueroussov S, Najafabadi HS, Hughes TR, Morris Q, Barash Y, Krainer AR, Jovic N, Scherer SW, Blencowe BJ, Frey BJ. RNA splicing. The human splicing code reveals new insights into the genetic determinants of disease. *Science*. 2015 Jan 9;347(6218):1254806.

Zhang Z & Wang W. (2014) RNA-Skim: a rapid method for RNA-Seq quantification at transcript level. *Bioinformatics* **30**(12):i283-i292.

Zhao K, Lu ZX, Park JW, Zhou Q, Xing Y (2013) GLiMMPS: Robust statistical model for regulatory variation of alternative splicing using RNA-Seq data. *Genome Biol* 14(7):R74.

## Tables

<i>Annotation</i>	<i>Dataset</i>	<b>Sailfish+SUPPA</b>		<b>RSEM+SUPPA</b>		<b>MATS</b>		<b>MISO</b>	
		<i>Pearson</i>	<i>Spearman</i>	<i>Pearson</i>	<i>Spearman</i>	<i>Pearson</i>	<i>Spearman</i>	<i>Pearson</i>	<i>Spearman</i>
<b><i>RefSeq</i></b>	<b><i>Synthetic</i></b>	0,971	0,959	0,987	0,978	0,833	0,819	0,862	0,815
	<b><i>ESRP1</i></b>	0,778	0,769	0,795	0,779	0,763	0,753	0,701	0,715
	<b><i>EV</i></b>	0,767	0,766	0,808	0,823	0,805	0,815	0,765	0,782
<b><i>Ensembl</i></b>	<b><i>ESRP1</i></b>	0,633	0,627	0,682	0,673	0,723	0,715	0,708	0,691
	<b><i>EV</i></b>	0,608	0,613	0,664	0,668	0,794	0,790	0,747	0,774

**Table 1.** First row: Correlation values (Spearman and Pearson) between the estimated and ground-truth  $\Psi$  values using simulated data. The comparison involves 144 events (Supplementary Data 1). Second and third rows: correlation values (Spearman and Pearson R) between the estimated  $\Psi$  values from ESRP1-overexpressed (ESRP) and empty-vector (EV) RNA-Seq datasets and the RT-PCR validation for the same samples (Shen et al. 2012). This comparison involves the 60 events that were in the RefSeq annotation and had a  $\Psi$  value from every method (Supplementary Data 2).

	Cufflinks+SUPPA		Cufflinks+Sailfish+SUPPA		MATS		MISO	
	<i>Pearson</i>	<i>Spearman</i>	<i>Pearson</i>	<i>Spearman</i>	<i>Pearson</i>	<i>Spearman</i>	<i>Pearson</i>	<i>Spearman</i>
<b><i>ESRP1</i></b>	0,613	0,627	0,549	0,582	0,622	0,610	0,605	0,611
<b><i>EV</i></b>	0,659	0,650	0,597	0,602	0,709	0,701	0,630	0,604

**Table 2.** Correlation values (Spearman and Pearson R) between the estimated  $\Psi$  values from ESRP1-overexpressed (ESRP) and empty-vector (EV) RNA-Seq datasets and the RT-PCR validation for the same samples (Shen et al. 2012). This comparison involves 83 events in the ESRP1 sample and 47 in the EV sample, which can be found in Supplementary Data 3.

## Figure legends

**Figure 1. SUPPA pipeline.** (A) SUPPA calculates possible alternative splicing events with the operation *generateEvents* from an annotation, which can be obtained from a database or built from RNA-Seq data using a transcript reconstruction method. For each event, the transcripts contributing to either form of the event are stored and the calculation of the  $\Psi$  value per sample for each event is performed using the transcript abundances per sample (TPMs) (Methods). From one or more transcript quantification files, SUPPA calculates for each event the  $\Psi$  value per sample with the operation *psiPerEvent*. SUPPA can use transcript quantification values obtained from any method. (B) Events generated from the annotation are given an unique identifier that includes a code for the event type (SE, MX, A5, A3, RI, AF, AL) and a set of start (s) and end (e) coordinates that define the event (shown in the figure) (Methods). In the figure, the form of the alternative splicing event that includes the region in black is the one for which the relative inclusion level ( $\Psi$ ) is given: for SE, the PSI indicates the inclusion of the middle exon; for A5/A3, the form that minimizes the intron length; for MX, the form that contains the alternative exon with the smallest start coordinate (the left-most exon) regardless of strand; for RI, the form that retains the intron; and for AF/AL, the

form that maximizes the intron length. The gray area denotes the alternative form of the event.

**Figure 2. Benchmarking with simulated data.** (A) Correlation of the ground-truth  $\Psi$  values (Methods) with those estimated with Sailfish+SUPPA using simulated data. The blue line and gray boundaries are the fitted curves with the LOESS regression method. (B) Cumulative distribution of the absolute difference between the ground-truth  $\Psi$  values and the ones estimated with Sailfish+SUPPA (SAILFISH), RSEM+SUPPA (RSEM), MISO and MATS. The lines describe the proportion of all events tested (Cumulative percent, y-axis) that are predicted at a given maximum absolute difference from the ground-truth value ( $\Delta\Psi$ , x-axis).

**Figure 3. Benchmarking using experimentally validated events.** (A) Correlation of the experimental  $\Psi$  values with those estimated with Sailfish+SUPPA in MDA-MB-231 cells with (ESRP1, left panel) and without (EV, right panel) ESRP1 overexpression. Experimental  $\Psi$  values were obtained using RT-PCR (Shen et al. 2012) and estimated  $\Psi$  were obtained from RNA-Seq data from the same samples (Shen et al. 2012). The blue line and gray boundaries are the fitted curves with the LOESS regression method. (B) Cumulative distribution of the absolute difference between the same experimental  $\Psi$  values and the ones estimated with Sailfish+SUPPA (SAILFISH), RSEM+SUPPA (RSEM), MISO and MATS from RNA-Seq data from the same samples (Shen et al. 2012). The lines describe the proportion of all events (Cumulative percent, y-axis) that are calculated at a given maximum absolute difference from the RT-PCR value ( $\Delta\Psi$ , x-axis).

**Figure 4. Annotation dependencies.** Boxplots of the difference of  $\Psi$  values estimated by SUPPA for Ensembl and RefSeq annotations from Sailfish quantification (y axis) as a function of (A) the difference in the number of transcripts defining each event in Ensembl and RefSeq or as a function of (B) the mean expression of the gene in which the event is contained. The x-axis in (B) is grouped into 10 quantiles according in the  $\log_{10}(\text{TPM})$  scale. The variability is represented for both replicates (7C1 and 7C2) of the cytosolic RNA-Seq

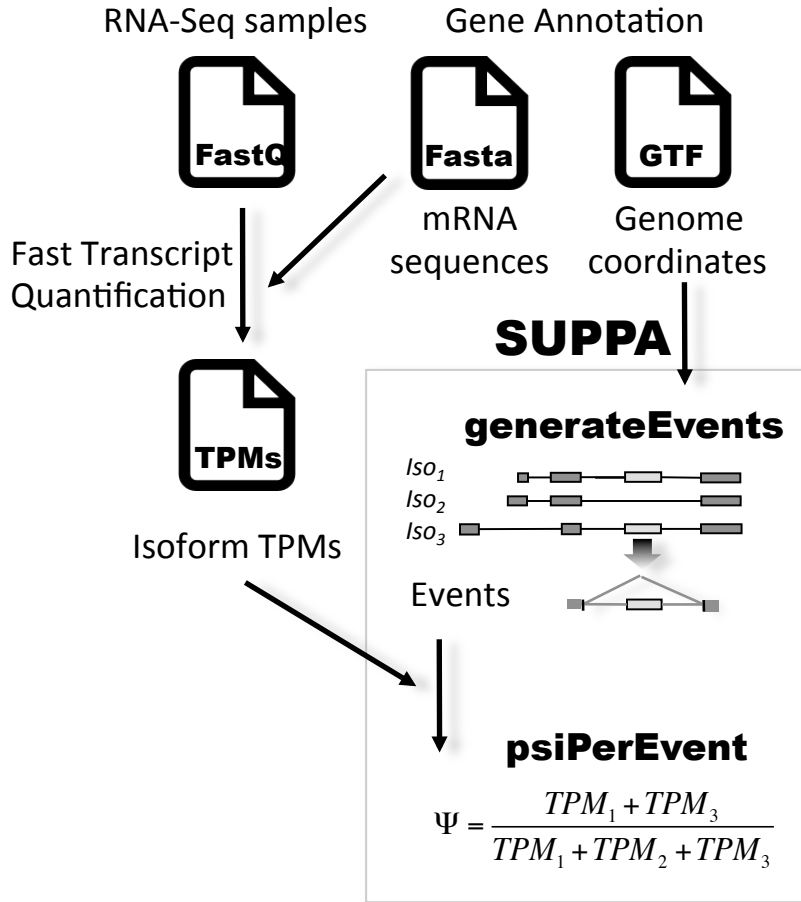
data from MCF7 cells. **(C)** Boxplots of the distribution of  $\Psi$  differences between replicates for the estimates from the Ensembl (left panel) and RefSeq (right panel) annotations as a function of the mean expression genes (x-axis), grouped into 10 quantiles in the  $\log_{10}(\text{TPM})$  scale, using genes with  $\text{TPM} > 0$ . Mean expression is calculated as the average of the  $\log_{10}(\text{TPM})$  for the each gene in the two replicates for (C) or for the each gene in the two annotations in (B).

**Figure 5. Annotation-free PSI estimation.** Correlation of the experimental  $\Psi$  values with those estimated with Cufflinks *de novo* + SUPPA in MDA-MB-231 cells with (ESRP1, left panel) and without (EV, right panel) ESRP1 overexpression. Experimental  $\Psi$  values were obtained by RT-PCR (Shen et al. 2012) and estimated PSIs were obtained from RNA-Seq data in the same samples (Shen et al. 2012). The blue line and gray boundaries are the fitted curves using the LOESS regression method.

**Figure 6. Speed benchmarking.** **(A)** Time performance for read assignment/mapping to transcript/genome positions by RSEM, Sailfish, STAR and TopHat on the synthetic as well as the ESRP1 and EV RNA-Seq datasets separately (Methods). RSEM and Sailfish include the transcript quantification operation. **(B)** Time performance for the  $\Psi$  value calculation from the already mapped reads (MATS, MISO) or quantified transcripts (SUPPA). ESRP1 and EV samples were pooled for this benchmarking (MDA-MB-231). MATS time includes the calculation of the  $\Delta\Psi$  between samples, which we could not separate from the  $\Psi$  calculation (Shen et al. 2012). All tools were run in multi-threaded mode when possible. Time reported for all cases is the actual cumulative time the process used across all threads (Methods).

Figure 1

(A)



(B)

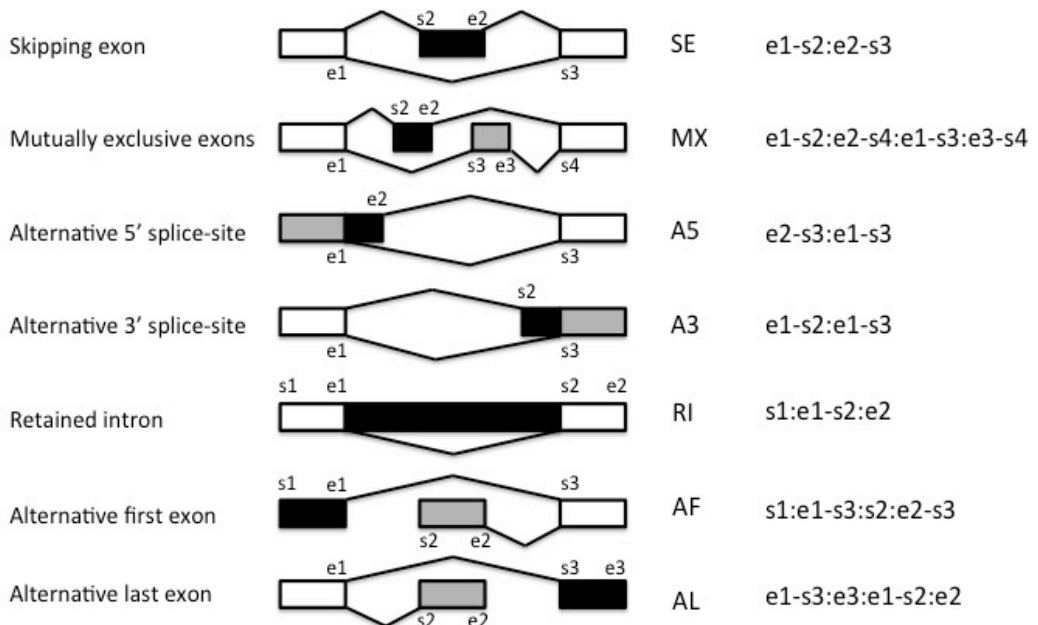
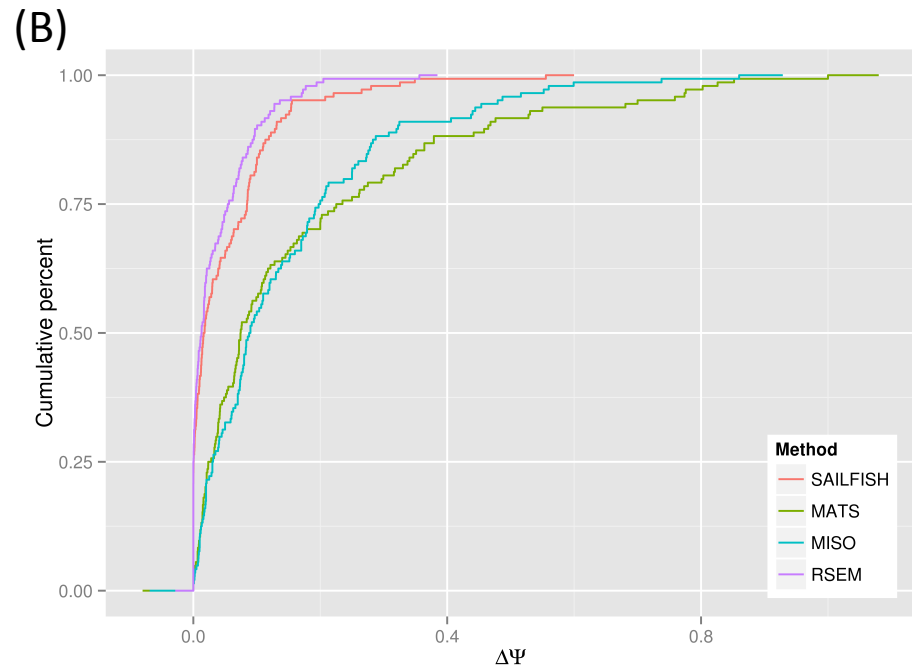
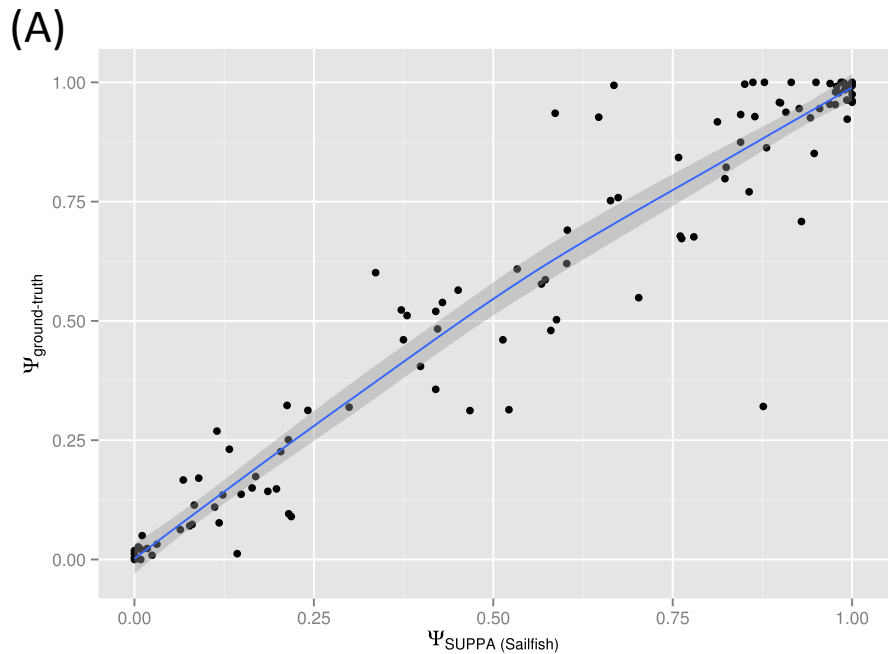


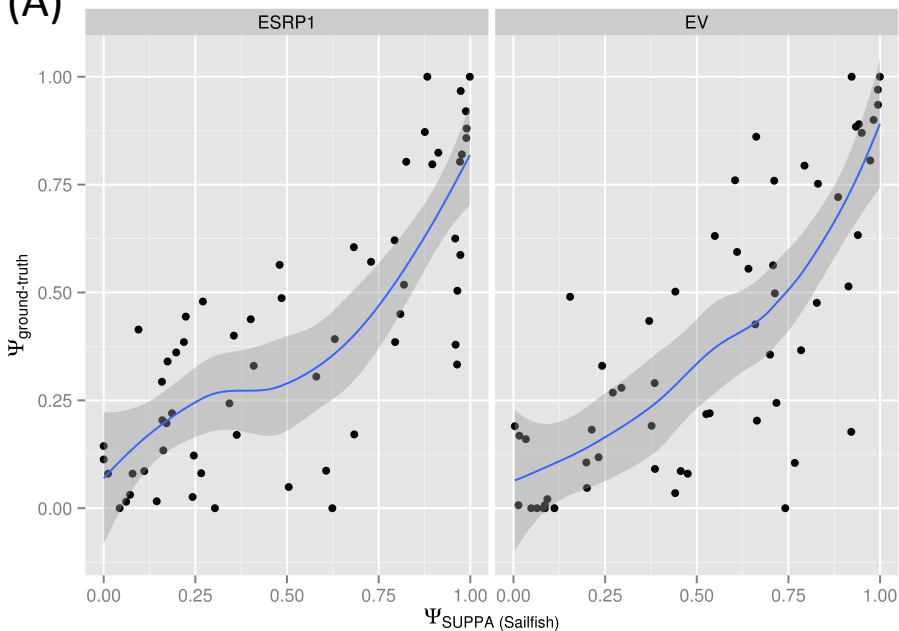


Figure 2

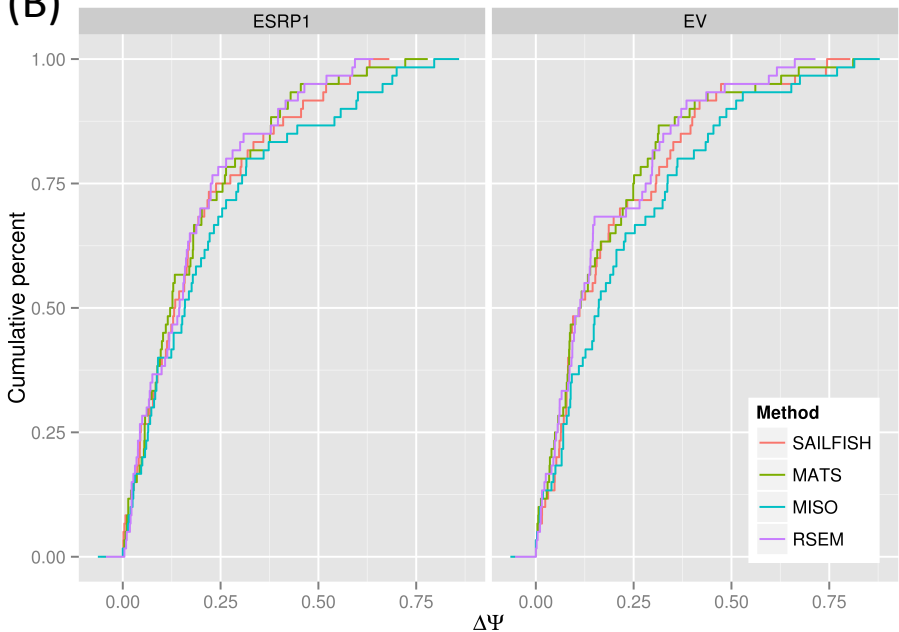


# Figure 3

## (A)

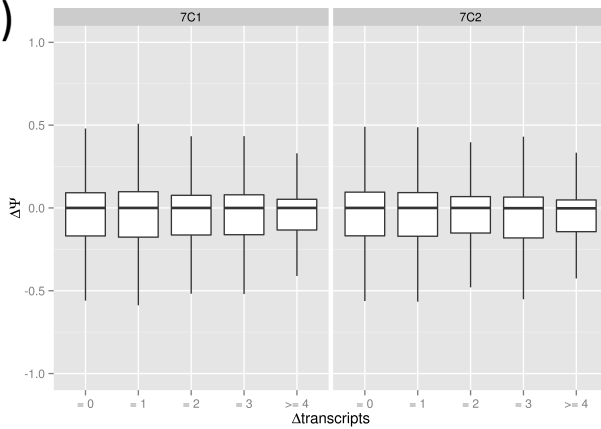


## (B)

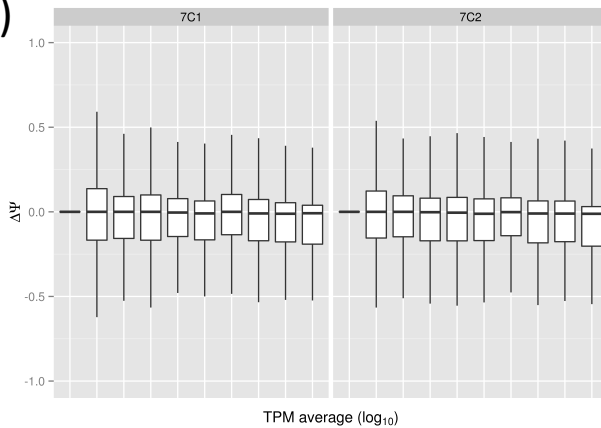


# Figure 4

## (A)



## (B)



## (C)

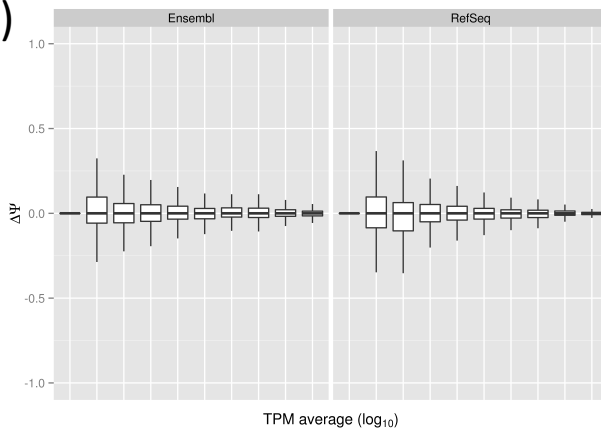


Figure 5

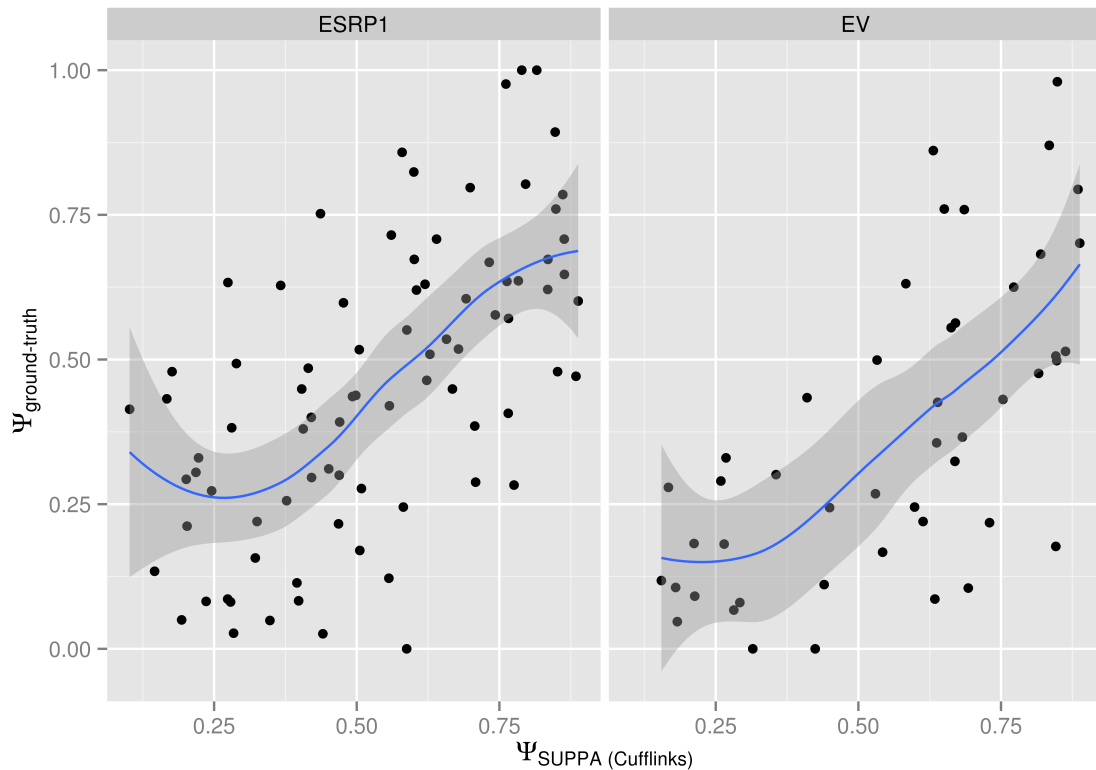
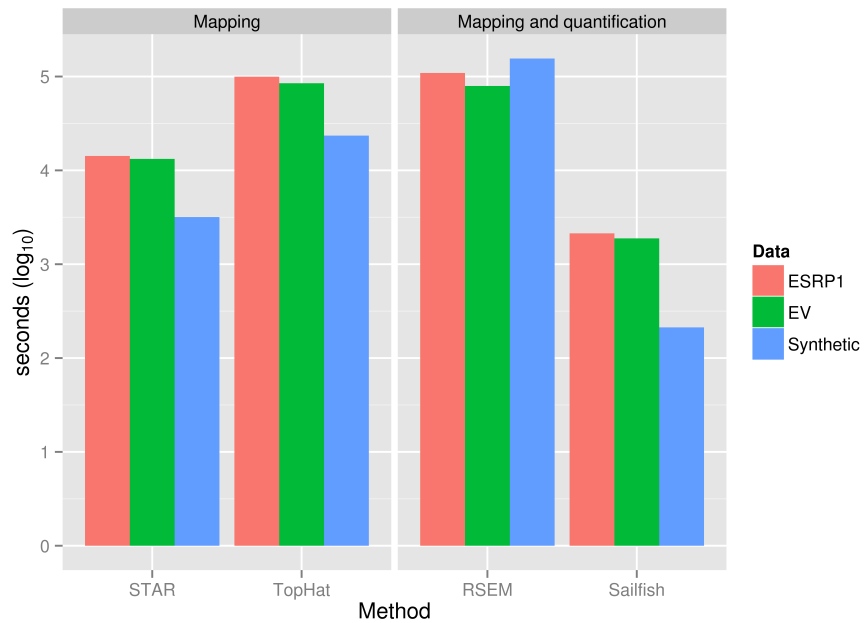


Figure 6

(A)



(B)

