


RESEARCH ARTICLE

Open Access



Comparison of RNA-seq and microarray platforms for splice event detection using a cross-platform algorithm

Juan P. Romero^{1†}, María Ortiz-Estévez^{2†}, Ander Muniategui¹, Soraya Carrancio², Fernando J. de Miguel³, Fernando Carazo¹, Luis M. Montuenga^{3,4,5,7}, Remco Loos², Rubén Pío^{3,5,6,7}, Matthew W. B. Trotter² and Angel Rubio^{1*} 

Abstract

Background: RNA-seq is a reference technology for determining alternative splicing at genome-wide level. Exon arrays remain widely used for the analysis of gene expression, but show poor validation rate with regard to splicing events. Commercial arrays that include probes within exon junctions have been developed in order to overcome this problem.

We compare the performance of RNA-seq (Illumina HiSeq) and junction arrays (Affymetrix Human Transcriptome array) for the analysis of transcript splicing events. Three different breast cancer cell lines were treated with CX-4945, a drug that severely affects splicing. To enable a direct comparison of the two platforms, we adapted EventPointer, an algorithm that detects and labels alternative splicing events using junction arrays, to work also on RNA-seq data. Common results and discrepancies between the technologies were validated and/or resolved by over 200 PCR experiments.

Results: As might be expected, RNA-seq appears superior in cases where the technologies disagree and is able to discover novel splicing events beyond the limitations of physical probe-sets. We observe a high degree of coherence between the two technologies, however, with correlation of EventPointer results over 0.90. Through decimation, the detection power of the junction arrays is equivalent to RNA-seq with up to 60 million reads.

Conclusions: Our results suggest, therefore, that exon-junction arrays are a viable alternative to RNA-seq for detection of alternative splicing events when focusing on well-described transcriptional regions.

Keywords: Alternative splicing, RNA-seq, Microarrays

Background

Alternative Splicing (AS) is known to play a major role in human biology, and the identification of transcriptional splicing patterns has potential uses for diagnosis, prognosis, and therapeutic target evaluation in the disease context [1, 2]. The development of exon microarrays enabled the transcriptomic study of differential splicing events, but PCR validation rates for identification of splice differences via microarray analysis tend to be lower than those observed for identification of

differential gene expression using similar technologies [3–5]. Junction arrays [6–10] have been proposed to overcome this problem by using oligonucleotide probe-sets that interrogate junctions between exons in the transcriptome, as well as the exons themselves.

Since the advent of next-generation sequencing (NGS), RNA-seq has become the technology of choice via which to detect and quantify alternative splicing (for a review see [11]). Various published works compare the performance of RNA-seq and expression microarrays for the analysis of gene expression [12, 13], but a thorough evaluation of both technologies in terms of their ability to detect differential AS events has yet to be presented. In the present study, we perform a comparison of

* Correspondence: arubio@tecnun.es

[†]Juan P. Romero and María Ortiz-Estévez contributed equally to this work.

¹CEIT and Tecnun, University of Navarra, Parque Tecnológico de San Sebastián, Paseo Mikeletegi 48, 20009 San Sebastián, Gipuzkoa, Spain

Full list of author information is available at the end of the article



RNA-seq technology (using the Illumina HiSeq platform) and junction arrays commercialized by Affymetrix (Human Transcriptome array, or HTA).

AS can be studied from two complementary points of view: with focus on transcripts or splicing events respectively. In the former, the subject of analysis is the transcript (or isoform), whereas in the latter, the subject(s) are the splicing events themselves.

The pipeline of the transcript-focused approach uses RNA-seq data with [14] and without known annotations in order to reconstruct the transcriptome and estimate the concentration values of the transcripts. Finally, the significance of change in absolute or relative concentrations is assessed using suitable statistical methods [15–17]. Transcript reconstruction is challenging [18] (even the better methods display transcriptome reconstruction levels below 50% when using simulated reads) and any error in reconstruction of transcript structure may be propagated to the output of statistical analysis. Moreover, the challenge of estimating isoform concentrations for genes with many transcripts yields wide confidence intervals [19].

On this basis, therefore, an event-based method appears a more suitable approach via which to compare AS detection technologies, with the additional benefit of straightforward validation using PCR. Event-based methods focus directly on the analysis of differential splicing events, rather than first attempting to estimate transcript concentration levels. These events can be classified into five canonical categories [20]: cassette exon, alternative 3', alternative 5', mutually exclusive exons and intron retention. In some cases, alternative start and termination sites are included also when defining splicing events. This approach has gained traction and several algorithms have been developed recently for detection of splicing events using RNA-seq data, including rMats, SplAdder, spliceGrapher or SGSeq [21–24]. SpliceGrapher and SGSeq detect events prior to application of separate software in order to state corresponding statistical significance, whereas rMats and SplAdder perform both detection and statistical analysis. Alongside NGS-based approaches, AS event detection methods are available for exon arrays [25], and exon-junction arrays [6, 8, 9, 26]. The latter methods display validation rates well above 50%.

The principal aim of this work is to compare RNA-seq and exon-junction microarray technologies in their ability to detect differential AS events. To do so comprehensively, and to allow as close to a direct comparison as possible, we have adapted the EventPointer [8] algorithm for application to data from both platforms, generated from the same control experiment. The control experiment comprises three distinct triple-negative breast cancer (TNBC) cell-lines, exposed in culture to a drug

known to affect the transcriptional machinery and, thereby, to induce AS events.

Further to comparative analysis of the resulting data, we conclude that both technologies show considerable concordance with high PCR validation rates, and that exon-junction microarrays have potential as an alternative to RNA-seq profiling for detection of AS events in annotated transcripts.

Results

CX-4945 is a potent and selective orally bioavailable small molecule inhibitor of casein kinase CK2 [27], which has been proposed previously as a cancer therapy [28], and which has been shown to regulate splicing in mammalian cells [29]. RNA samples taken from three distinct triple-negative breast cancer (TNBC) cell-lines, exposed to CX-4945 and also to a DMSO control, were profiled using both RNA sequencing¹ and hybridization to exon-junction microarrays (see Methods for details). We extended the EventPointer algorithm (available via Bioconductor, see Methods) for application to data from both platforms and applied it to the corresponding datasets in order to identify AS events.

Prior to the comparison of platforms for splice event detection, the data was assessed at the gene level in order to ensure signal quality and coherence. Gene expression was computed from RNA-seq data using Kallisto [14] to quantify expression as the sum of isoform concentrations for each gene. RMA [30] was used to quantify gene expression from microarray data, using annotation files from Brainarray [31]. The same version of the Ensembl Transcriptome (Ensembl v.74, GRCh 37.75) was used in both cases.

Considering each technology independently, correlation between sample replicate profiles in each cell-line and experimental condition is high for both platforms (correlation coefficient ranging from 0.988 to 0.996 in arrays and 0.996 to 0.997 in RNA-seq). When comparing profiles from the same samples between technologies, strong coherence is observed for well-expressed genes. Median correlation of gene expression between technologies on the same samples is 0.510, and gene expression patterns across all samples display correlation of 0.680 between technologies. The first one is smaller owing to the different probe affinities of the set of probes that interrogates each gene. When only the 50% most highly expressed genes are considered, the median correlation of gene expression patterns is 0.750 (Additional file 1). The gene expression correlations observed are similar to previously reported comparisons between RNA-seq and exon arrays [32]. It is important to point out that the expected correlations for gene expression are larger (either using microarrays or RNAseq) since

the number of probes/reads that interrogate a gene is larger than the ones that interrogate a splicing event.

Events detected by RNA-seq and junction arrays show strong qualitative and quantitative concordance, with a subset detected exclusively by one of the technologies

Figure 1 depicts the EventPointer pipeline for both profiling technologies (see original publication [8] for further detail), with CEL files (microarray) or BAM files (RNA-seq) as starting input. When building the splicing

graph, each exon is split into two nodes that correspond to its start and end genomic positions respectively (Fig. 1b). Each event is described by two alternative paths (Paths 1 and 2) and a shared reference path (Path Ref) within the splicing graph. These paths are sets of edges in the splicing graph. Paths 1 and 2 are mutually exclusive in terms of isoforms (i.e. if an isoform includes Path 1 it does not include Path 2 and vice versa) and all isoforms interrogated by the event share the reference path. Therefore, events are contained in several isoforms

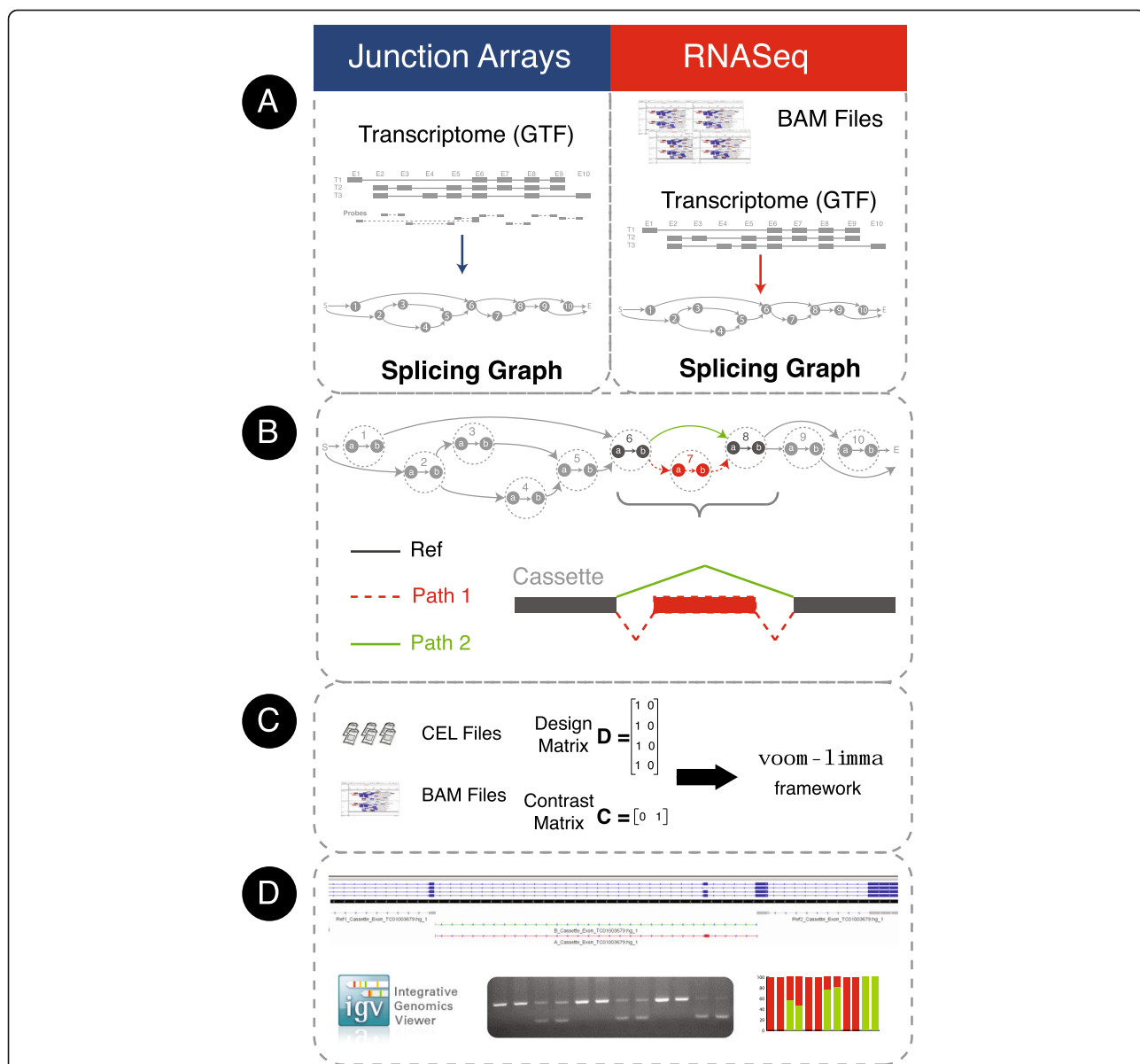


Fig. 1 EventPointer overview for junction arrays and RNA-Seq data. **a** The CEL or BAM files are the input data for each technology. The splicing graph for each gene is built using the array annotation files or directly using the sequenced reads. **b** Each node in the splicing graph is splitted into two nodes that correspond to the start and end positions in the genome respectively. EventPointer identifies events within each gene and annotates the type of event. In the figure, among the events in the gene, an exon cassette is highlighted. **c** Statistical significance of the events is computed. **d** Finally, the top-ranked events are validated using PCR and the results visualized in IGV

(at least two). A simple example would be the cassette exon shown in Fig. 1b: the reference path is composed by the edge that links nodes 6a and 6b (i.e. the coverage of exon 6 or the signal in the probe-set of the array that interrogates this exon) and the edge that links nodes 8a and 8b (coverage of exon 8). All these measurements are summarized into one average value. Path 1 includes the edges in the path 6b-7a-7b-8a (coverage of exon 7 and its flanking junctions) and Path 2 is the edge that links 6b and 8a (coverage of the skipping junction). EventPointer distinguishes between events as corresponding to: cassettes; alternative 5'; alternative 3'; mutually exclusive exons; alternative first exons; and alternative end exons. Complex events that do not match any of these categories are denoted as such.

Three different cell-lines were profiled, each exposed to CX-4945 and DMSO respectively across five replicates. AS differences were tested using a linear model which controlled for cell-line differences. Using the read coverage (or probe-set signal) for each path, a statistical analysis based on voom-limma [17, 33] is applied to determine the significance of each event via comparison between alternative path signals (see Methods for details). In addition to the statistical analysis, we compute the Percent Splice Index (PSI or Ψ) [34], an estimate of relative isoform concentrations that map to paths 1 and 2 for each event. In a cassette event, if the exon is retained, Ψ is equal to one. If it is skipped, Ψ is equal zero. If both isoforms that retain and skip the exon are present, Ψ is the ratio between the expression of the isoforms that retain the exon and the overall expression of the isoforms that skip or retain the exon. Ψ has become the standard method to quantify splicing events.

In order to identify well-expressed events (more likely to be biologically significant and less prone to validation error), the comparison of AS detection was performed on a subset of the data with expression above a set threshold (see Methods for details). In brief, a junction coverage threshold was applied to the RNA-seq data (default 2 FPKM) and a threshold on expression percentile

applied to the microarray data (default probe-set expression greater than 25% of probes in any sample profile).

Table 1 displays the number of detected events after setting a threshold on the expression for both technologies, the number of differentially spliced events (p value < 0.001) with their corresponding False Discovery Rate (FDR) detected via application of EventPointer to RNA-seq and microarray data respectively. The statistical analysis compares differential AS in profiles from cell-lines treated with CX-4945 and with DMSO control. As may be expected, setting more stringent expression thresholds yields fewer events detected with better FDR on both platforms. The FDR is the estimated proportion of false discoveries (i.e. events are assumed to have differential splicing and that do not). For example, an FDR of $5e-4$ means that 0.05% of the selected events are expected to be false positives.

Table 1 shows that fixing p -value to $1e-3$ yields False Discovery Rates (FDRs) less than 1% for both technologies. The expected proportion of AS events appears high ($1-\pi_0$ approx. 46%) [35], i.e. more than 46% of the events have its splicing patterns altered, which reflects the anticipated strong effect of compound exposure on the splicing machinery. It is also apparent that, for a similar number of detected AS events, the FDR corresponding to RNA-seq analysis is smaller.

Events detected by both technologies (referred to as "matched events" hereon) were defined by a stringent criterion in which nucleotide sequences of paths identified via one technology must be a subset of sequences identified via the other, yielding 6222 matched events. When reporting correspondence and divergence between AS events, below, the following naming convention is used: R^+ represents number of events deemed significantly altered in RNA-seq analysis (p value $< 1e-3$); R^- represents number of events deemed not significantly altered in RNA-seq analysis (p value > 0.2). M^+ and M^- are the counterpart terms used to describe microarray results. Events not detected by each technology are labelled $R\emptyset$ and $M\emptyset$ respectively.

Table 1 Number and statistical significance of detected AS events using both RNA-seq and array technologies

Expression Threshold	Detected Events	Significant events	FDR for significant
RNASeq			
Junction coverage > 6 FPKM	9277	4526	2.7e-4
Junction coverage > 2 FPKM	34,961	13,780	4.7e-4
Junction coverage > 2/3 FPKM	92,986	29,443	7.0e-4
Exon-junction arrays			
Signal > 50%	10,114	2385	9.2e-4
Signal > 25%	31,506	6197	1.37e-3
No threshold	92,405	11,761	3.45e-3

for different expression thresholds, default filters are junction coverage greater than 2 FPKM for RNA-seq and probe-set signal greater than top 25% quantile for microarray

A subset of matched events is significant in both technologies (R^+M^+) and shows coherent change in the corresponding Ψ . There are also significant events detected by only one of the technologies ($R^+M\emptyset$ and $R\emptyset M^+$). The summary of findings is presented in Table 2.

Table 2 shows that the FDR of the events detected only by RNA-seq is similar to that for events detected by both platforms (4.56e-4 vs. 4.96e-4). In other words, the reliability of events discovered only by RNA-seq is similar to that of events identified by both technologies. In the case of the arrays, the FDR of matched events is three times smaller than for those discovered solely by the arrays (1.99e-3 vs 6.23e-4 i.e. $R\emptyset M^+$ events are less reliable than R^+M^+ events for the same p -value threshold. In addition, Table 2 shows that the number of significant events that are RNA-seq specific ($R^+M\emptyset$) is larger than the number of significant events detected only by arrays ($R\emptyset M^+$) (10,617 vs 3297 events).

Figure 2 depicts a Sankey diagram of the relationship between matched events. An event is declared to be significant (in either technology) if the p value is smaller than 1e-3. It is declared non-significant if the p value is larger than 0.2 and inconclusive otherwise. It is apparent that many events that are significant for RNA-seq are not detected by arrays, but also that events significantly detected via arrays are not detected by RNAseq. Most matched events are consistent across technologies: significant events for one technology are also significant for the other.

We also considered the FDR for different types of splicing events in both technologies. As shown in Fig. 3, alternative 3' (5'), start and end sites have larger FDR than cassette exons, i.e. they are harder to measure. There were too few matched mutually exclusive events to estimate accurately FDR for this type of events.

PCR validation rates are over 80% in both technologies

PCR validation was performed on a subset of predicted AS events drawn from each of the subsets discussed in previous sections, i.e. events detected by one or both technologies. PCRs were performed on:

1. Five top-ranked events detected by both technologies (*topRNA* and *topArrays*) regardless of the matching with the other technology
2. Five top-ranked events detected by one technology ($R^+M\emptyset$ and $R\emptyset M^+$)
3. Five top-ranked events significant in one technology (R^+M^- and R^-M^+)
4. Five top ranked events detected by both technologies (R^+M^+)

These potential 35 validations are in fact 29 since there is overlap in the top-ranked events of different categories. The characteristics of validated events (genome location, event type, etc.) and links to the corresponding PCR images are included in Additional file 2. PCR for events in non-coherent classes (R^+M^- , R^-M^+) required up to 40 PCR cycles and were harder to validate in general. The corresponding GTF files to browse these events in IGV [36] are included in Additional file 3. All the results are summarized in Table 3.

Figure 4 shows the Ψ estimates and the PCR bands for two of the top-ranked events in R^+M^+ (gene names *DONSON* and *MELK*), with clear concordance of the splice index, Ψ , across the three technologies despite use of end-point (i.e. non-quantitative) PCR. Similar figures for events in the other AS categories are included in the additional material (Additional file 1: Figures S2 to S8).

Statistics and Ψ for matched events are similar

Figure 5 shows the increment of the Ψ value estimated by EventPointer for events detected by both technologies. Correlation for AS events is over 0.90, and z-values of the statistical test are also similar (Additional file 1: Figure S1). PCR figures also show high coherence between the estimated Ψ using both technologies, especially for RNA-seq, and the PCR results (Fig. 4 and Additional file 1: Figure S2 to S8).

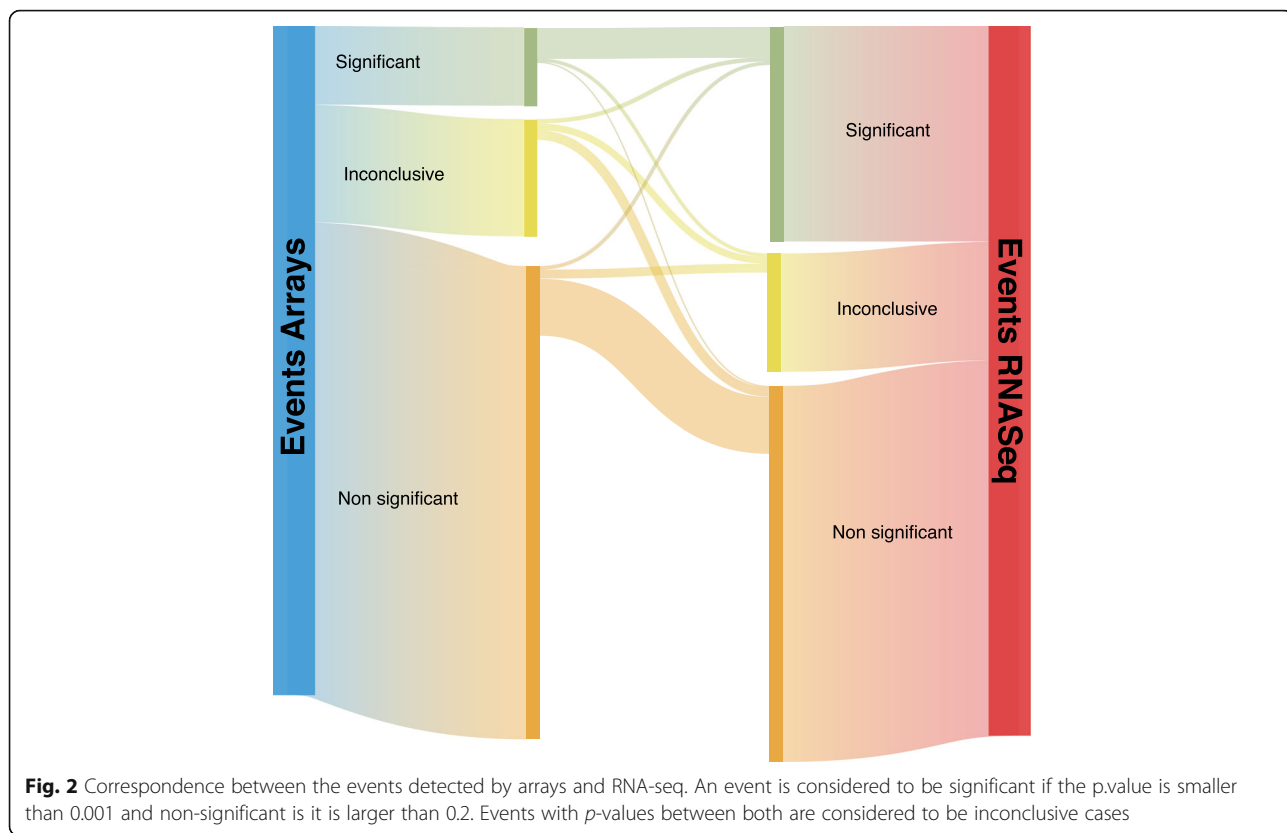
Both technologies detect a similar distribution of AS types

Figure 6a shows the number and type of AS events detected by the EventPointer algorithm on data from both

Table 2 Number of AS events detected per technology, alongside statistical significance of events against distinct thresholds

Matched Events					
Matched Events	Significant in both (R^+M^+)		FDR (RNASeq)	FDR (arrays)	
6222	1324		4.96e-4	6.23e-4	
	$R^+M\emptyset$			$R\emptyset M^+$	
Expression Threshold	Detected Events	FDR	Expression Threshold	Detected Events	FDR for significant
Junction coverage > 6 FPKM	2973	2.44e-4	Signal > 50%	1016	1.46e-3
Junction coverage > 2 FPKM	10,617	4.56e-4	Signal > 25%	3297	1.99e-3
Junction coverage > 2/3 FPKM	25,063	6.90e-4	No threshold	7581	4.85e-3

Where thresholds not shown, default filters were employed (junction coverage > 2 FPKM for RNA-seq; upper quartile probe signal for microarrays)



profiling technologies. The number of detected cassette exons using arrays is smaller than that using sequencing (*p* value <1e-16, test for equality of proportions). In fact, after matching the events detected by both technologies, a large proportion of the cassette exons in RNA-seq appear as complex in microarrays (see Fig. 6b). The reason

for this disparity is the complexity of the reference transcriptome used in the HTA array. For this analysis, we used the transcriptome provided by Affymetrix, which includes a range of annotation sources, e.g. RefSeq, Vega, Ensembl, MGC (v10), UCSC known genes and other sources for non-coding isoforms. The underlying

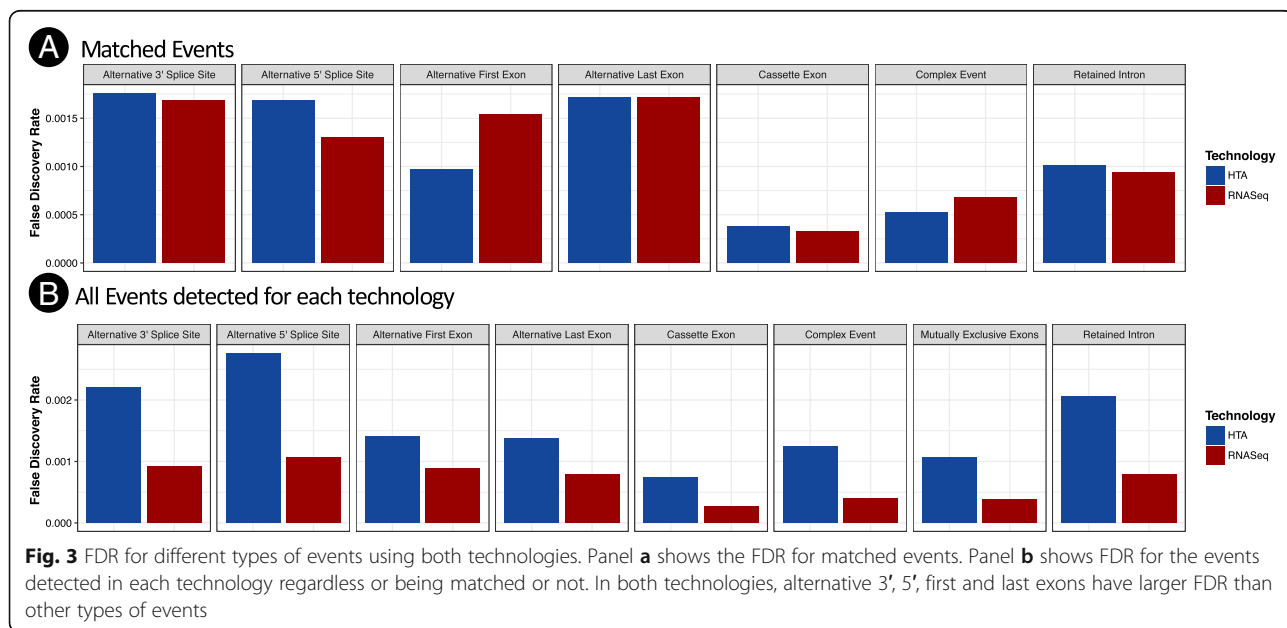


Table 3 PCR validation for RNA-seq and microarray technologies across events detected by one or both technologies

AS Event Category	RNA-seq	Arrays
Top-ranked events (topRNA, topArrays)	5/5	5/5
Significant in RNA-seq and not detected by arrays (R + M∅)	5/5	–
Significant in arrays and not detected by RNA-seq (R∅M+)	–	5/5
Detected by both. Significant in RNA-seq, not significant in arrays (R + M-)	5/5	
Detected by both. Significant in arrays, not significant in RNA-seq (R-M+)	3/5	
Detected by both. Significant and coherent events (R + M+)	5/5	

Values reported are validations / events selected

transcriptome for HTA includes such a variety of isoforms that many detected AS events are labelled as complex.

In addition, the proportion of retained introns is smaller for RNA-seq (p value $<1e-16$, test for equality of proportions), perhaps owing to the coverage required to include a region as expressed by SGSeq (defaults to 0.5 FPKM) which may exclude weakly expressed introns.

Power of arrays to detect events is approximately equivalent to shallow RNA-seq

The comparisons above suggest that that RNA-seq – at the depth of sequencing deployed here - detects a larger number of AS events at lower FDR.

We subsampled the initial RNA-seq data to 30% and 10% of the input, yielding approximately 30 million and 10 million reads respectively. Using these subsampled data, we estimated their FDR (see Table 4). Interpolating

the FDR for them, the FDR using junction arrays is equivalent to the FDR of an RNA-seq experiment with sequencing depth of approximately 20 million reads.

We hypothesize that the performance of the arrays could be greatly improved by removing bad-performing probesets. The RMA summarization algorithm withstands the presence of a few outlier probes. In fact, we have identified some probes that cross-hybridize in several loci of the transcriptome. However, if most of the probes that interrogate either of the paths or the reference do not perform well, the whole estimate of the splicing event will be compromised. Some of these cases can be detected since the signal of the events do not show internal coherence with the model (i.e. they show a large relative error if the weighted sum of the signals in Paths 1 and 2 and the reference Path are compared). These bad probesets are somehow expected: the design of junction probes has strong limitations since there is

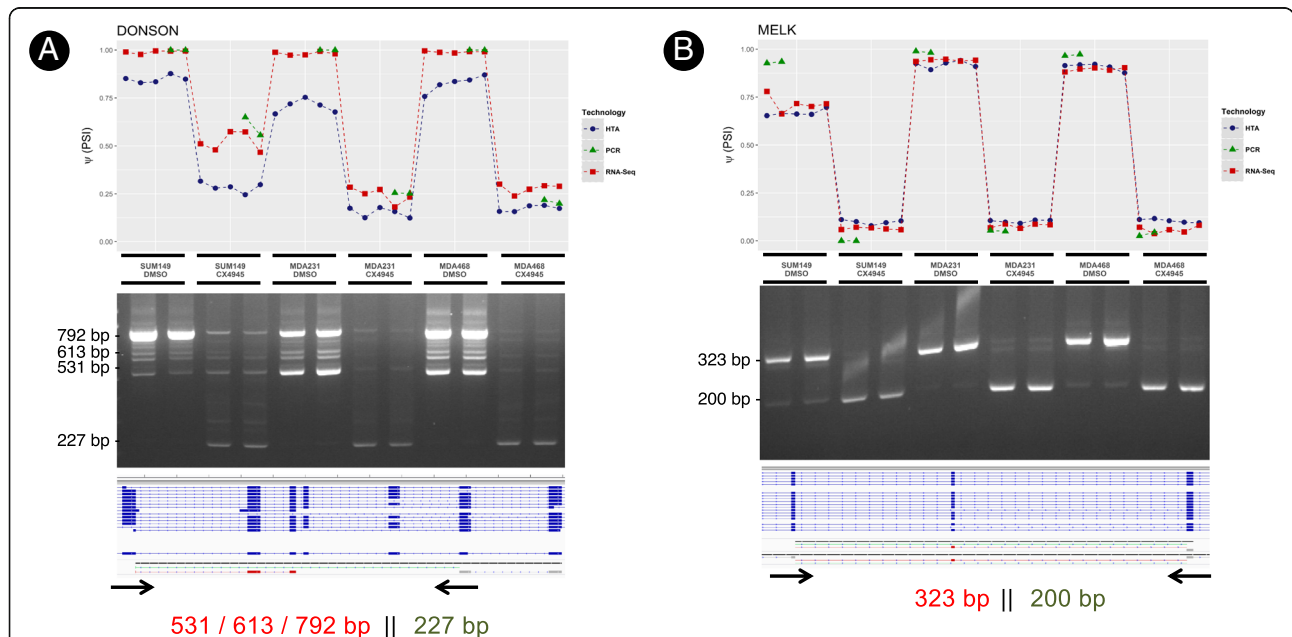
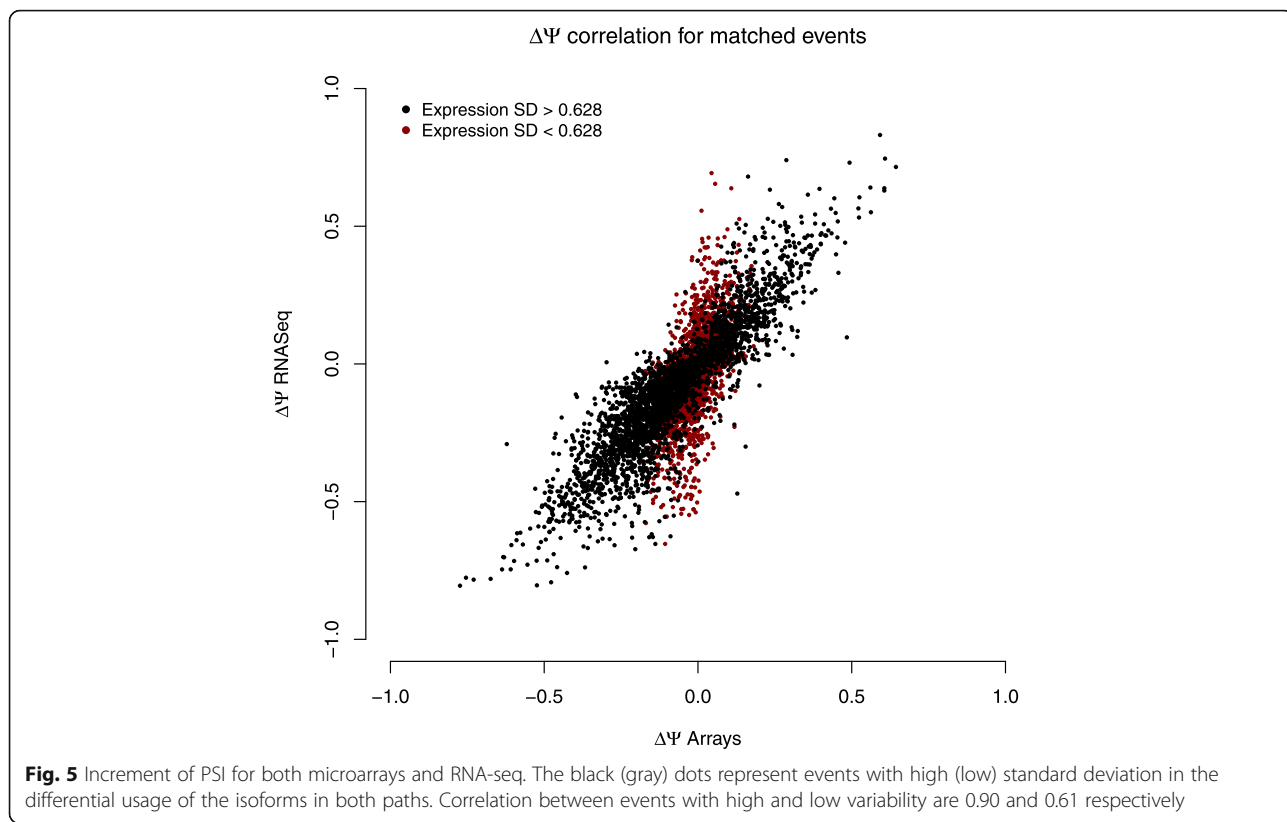


Fig. 4 Estimated PSI (for RNA-seq, microarrays and PCR image analysis), PCR bands, the reference HTA transcriptome and the alternative paths of the *DONSON* (panel a) and *MELK* (panel b) genes in R^*M^+ . Each of the points represents the same replicate in either of the three technologies. The last numbers shown are expected bands for the selected primers. If the number is shown to the left side of the double bars, the band corresponds to Path 1 of the event (long path). If shown to the right side, corresponds to Path 2 (short path)



no room to select a probe with certain standards of quality (GC content, no cross-hybridization against the genome of the transcriptome, etc.). Owing to these probesets, a number of events are not being measured accurately. We have included in the additional material Additional file 1: Figure S9 to illustrate bad and good performing probesets: in panel A, it is shown an event with internal coherence and in panel B, an event with

bad internal coherence. We have also included Additional file 1: Table S10 that shows the FDR for genes with large coherence (small relative error) and small coherence (large relative error). The FDR for genes with large internal coherence is 2 times bigger than for events with weak internal coherence.

The events that are not matched with RNA-seq are enriched in these pathological cases as shown in Table 2.

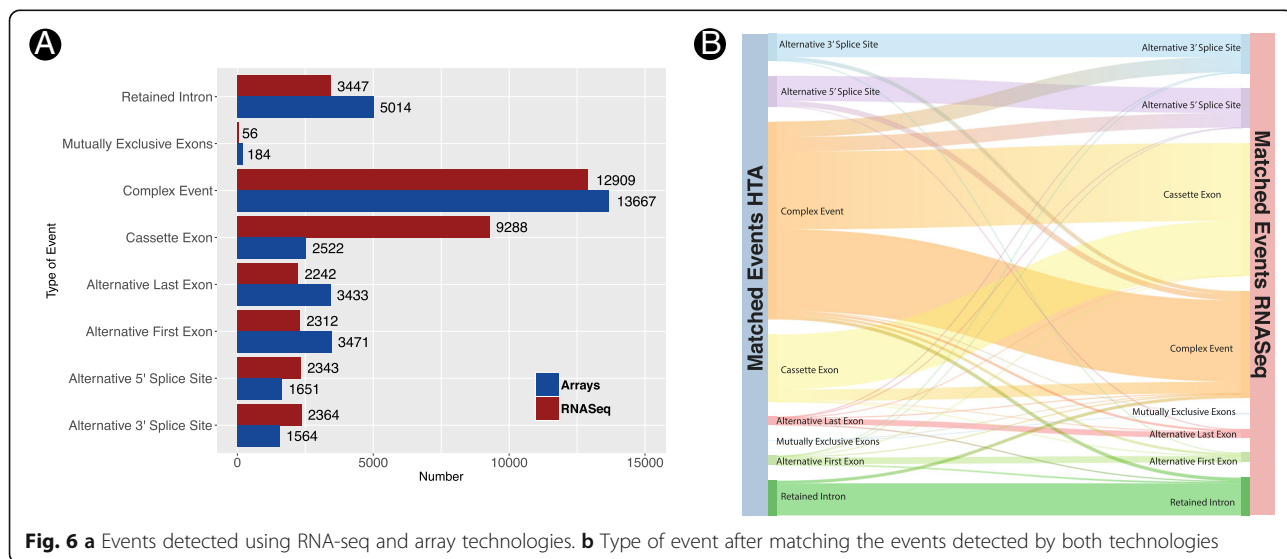


Table 4 Obtained FDR values after subsampling the number of reads in the RNA-seq experiment

Decimation percentage	FDR
Decimated 10%	1.86e-3
Decimated 30%	6.97e-4

On the contrary, the expected false discovery rate for matched events finds HTA arrays to be equivalent to RNA-seq with a depth of approximately 60 million reads. A proper filtering of the probes identifying events prone to errors could ideally take the arrays closer to the RNA-seq performance with this depth.

Discussion

The main aim of this work was to quantitatively compare the performance of RNA-seq and junction array technologies to detect splicing events. To do this in a balanced manner, we adapted our algorithm EventPointer, originally developed for HTA arrays, to work also on RNA-seq data.

This study highlights:

- The creation of a real-world cell exposure dataset specifically relevant for the study of alternative splicing.
- Adaptation of an existing AS event detection algorithm to a cross-platform method to enable comparative application, and addition of percent splice index method.
- Strong correlation of splicing event detection in regions covered by both technologies, validated by PCR on a subset of top ranked events identified by both and each platform respectively.
- Benefits of RNA-seq in terms of coverage and flexibility, as expected, and higher validation rates in case of disagreement between technologies.
- Good performance of HTA arrays, estimated by approximation to be equivalent to relatively shallow RNA-seq in transcript regions covered.

Top-ranked events detected by each platform technology and estimates of relative event occurrence ($\Delta\Psi$) were validated by PCR. The relative occurrence estimates were also strongly correlated, close to 0.90 for events detected by both technologies. In addition to enabling comparison of the two profiling platforms, these results suggest also that the estimates themselves are a relevant addition to the original EventPointer algorithm.

We relied on SGSeq to build the splicing graph. Using other algorithms (such as Spladder) could impact the detected events especially for weakly expressed genes. According to a recent review of some of the authors [37], AltAnalyze is the only alternative that provides the

analysis of splicing events both for arrays and RNA-seq. AltAnalyze characterizes each event by several values that must be integrated. For example, in a cassette exon, the signals of the probes in the exon, the flanking junctions and the skipping junction (and the equivalent coverage values for RNAseq) should be integrated to get a single figure of merit. We found it difficult to perform this integration since we are not the developers of this method. Nevertheless, using this method could be also informative to compare both platforms.

As might be expected in the absence of physical probe-sets, over 10,000 statistically significant events were identified by RNA-seq alone, the top ranked of which were validated via PCR. Approximately 3300 events were detected using microarrays but not detected using RNA-seq. In this case, some (3/5) of the top ranked events were validated and correspond to well-expressed genes. Those which did not may reflect the specific technical biases of each technology (cross hybridization of the probes, multi mapping reads, GC dependence, etc.)

A recent study [38] also compared RNAseq and arrays. This study was focused on patient derived samples instead of cell-lines as we did. This study pinpoints differences in the output between both technologies and states it in the title. In our case, the main divergence that we found between both technologies appears in the events that are not matched. Interestingly, the coherence between matched events -shown by Fig. 5 and the validated PCRs- is very strong: if an event is detected by both technologies, the increment of Ψ and its statistical significance (see Additional file 1: Figure S1) are very similar.

RNA-seq has inherent advantages over microarrays, including the ability to detect unlimited novel events. Furthermore, sensitivity can be improved by increasing sequencing depth. Another advantage of RNA-seq is its better approximation of gene/transcript concentrations (e.g. allowing to state a threshold based on the expression of an event). On the other hand, arrays were able to detect some weakly expressed events missed by RNA-seq and, in general across the comparisons, performed similarly to RNA-seq when treating well-expressed and well-defined transcriptional regions. As expected, a similar algorithmic approach applied to both platforms consumed less time and memory resources when treating microarray data than when treating RNA-seq data (See Table 5).

Conclusions

In conclusion, comparison of RNA-seq and junction microarrays using a cross-platform algorithm suggests that both technologies provide accurate identification of splice events. Moreover, predictions by both

Table 5 Resources required for both technologies. Analysis was performed on 16 cores (Intel Xeon E5–2670 @ 2.60 GHz) with 64 GB of RAM Linux server running 64-bit CentOS distribution

	Computing time		Memory requirements		Storage requirements	
	RNA-seq	HTA	RNA-seq	HTA	RNA-seq	HTA
Mapping to the transcriptome (STAR)	11.5 h	–	32Gb	–	1023 GB	–
Splicing graph generation (SGSeq)	2 days (5 cores)	14 h	8 Gb per core	5 Gb	70.3 Mb	2 Gb
Event detection	7 min 16 s (10 cores)		1.2 Gb per core		643.6 Mb	
Statistical analysis	1 min 43 s	3 min 06 s	2 Gb	< 1Gb	6.2 Mb	11 Mb

technologies tend to correlate strongly and yield similar results when compared by Ψ estimates and PCR. RNA-seq holds a clear advantage in terms of flexibility, and stronger PCR validation of events detected in one platform but not the other. As compared, HTA microarrays are shown nevertheless to provide a reasonable alternative to relatively shallow RNA-seq in the transcriptional regions that they reference.

Methods

Sample preparation

Triple negative breast cancer cell lines MDA-MB-231 and MDA-MB-468 were obtained from ATCC (Manassas, VA) with identification numbers HTB-26 and HT-132 respectively and SUM149 was purchased from Asterand plc (Detroit, MI). All cell lines were grown according to the suppliers' recommendation. CK2 inhibitor CX-4945 (Selleckchem, Houston, TX) was dissolved in DMSO and stored frozen at -80°C until used.

To induce splicing events, cells were grown to $\sim 70\%$ confluence and treated with $1\ \mu\text{M}$ CX-4945 or DMSO during 12 h in a total of 5 replicates per condition. Total RNAs were isolated using the RNeasy Mini Kit (Qiagen, Germantown, MD) according to the manufacturer's protocol. Integrity of RNA was quantified using the Agilent 2100 Bioanalyzer (Agilent Biosystems, Foster City, CA). Samples were labeled and hybridized in Human Transcriptome arrays (HTA) by the Genomics Core Facility of the Center for Applied Medical Research (CIMA) following manufacturer's instructions.

RNAseq was performed in the Center for Cooperative Research in Biosciences (CICBiogune) using the Illumina HiSeq2000 sequencing technology, HiSeq Flow Cell v3 and TruSeq SBS Kit v3. $2\ \mu\text{g}$ of RNA of each sample was sent for this purpose. The run type was strand specific, multiplexed with paired-end reads of 100 nucleotides each. The amount of RNA for hybridization and validation purposes was $5\ \mu\text{g}$.

STAR 2.4.0 h1 was used to align the reads against the human genome. The reference genome was Ensembl v.74, GRCh 37.75. The output were sorted BAM files. All the other parameters were set to the default values.

The average sequencing depth was 49 million reads (9.8 billion nucleotides sequenced per sample).

The microarray data preprocessing was performed using the aroma.affymetrix framework using the standard RMA algorithm applied to probesets of the paths [30]. In addition, we used both platforms to quantify expression at the gene level. Results are shown in the Additional file 1. Gene expression was computed from RNA-seq data using Kallisto [14] to quantify expression using a pseudo-alignment method. Kallisto returns an estimate of the expression of all the isoforms for each gene. The overall expression of the gene was simply computed by summing up the expression of each the its isoforms. We estimated the expression of the arrays using the RMA algorithm in the aroma.affymetrix framework using a Brainarray reference file of the Ensembl 74 transcriptome.

Event pointer for RNAseq

EventPointer is an R package to identify, classify and analyze alternative splicing events using microarrays and RNA-Seq data. The software is available for download at Bioconductor. A thorough description of EventPointer for microarrays can be found in [8]. This method has been extended to RNA-seq.

The concepts for detection, classification and statistical analysis are shared in EventPointer for the analysis of both technologies. The main difference of EventPointer for RNA-seq compared with that of microarrays are the ones associated with the type of input data (CEL or BAM files). The R code for the analysis is available at <https://github.com/jpromeror/SplicingComparison>.

EventPointer requires a splicing graph -a directed graph used to represent the structure of the different isoforms of a given gene [39] - as input to detect splicing events. EventPointer for RNA-seq uses SGSeq [24] to build the corresponding splicing graphs from BAM files. The complexity of the splicing graph can be controlled in SGSeq by setting different thresholds in the expression values of splicing junctions of the splicing graph (by default set to 2 FPKM). For RNA-seq, the splicing graphs are constructed for every single experiment. On

the contrary, in the case of microarrays the same splicing graph (and the corresponding CDF) is used for all the experiments run on the same type of microarray (HTA or, more recently Clariom-D).

The input data for the statistical analysis is different in both technologies: signal values of the probes in microarrays and counts in RNA-seq. In order to deal with reads, Voom [17] is applied to preprocess the RNA-seq count data. The statistics to deal with the processed RNA-seq data is identical to the one used for microarray data and hence, the same statistical tests -based on limma [33]- are applied to both technologies.

As output, EventPointer provides a table with the following information associated to each detected alternatively spliced event: gene identifier, genomic position, type of event, statistical parameters and $\Delta\Psi$ values. Additionally, EventPointer generates a “Gene Transfer Format” (GTF) file that can be used with the Integrative Genomics Viewer (IGV) [36] to view the structures of each detected alternative splicing event. This visualization facilitates the interpretation of the detected events and the design of primers for the validation of the events using standard PCR.

Estimation of PSI

We have included a novel algorithm to estimate Ψ that can be applied to both RNA-seq and microarrays. Assuming that the signal of a probe-set in microarrays and the number of reads within a region of the transcriptome in RNA-Seq depend on the product of an affinity value of the probe-set (or the equivalent length in RNA-seq) and the concentration of the interrogated isoforms in the paths, the following equation holds

$$S_i = a_i \cdot t_i \tag{1}$$

where S_i is the measured expression value of path i , a_i is the affinity of the probes or equivalent length of the path i and t_i is the concentration of the isoforms mapped to path i . The affinity values (or equivalent lengths) and concentration values are assumed to be unknown and must be estimated from the data.

Particularizing the above equation to each of the paths and taking into account that the concentration of the isoforms in the reference path must be the sum of those of paths 1 and 2, the following equations are obtained:

$$S_1 = a_1 \cdot t_1 \tag{2}$$

$$S_2 = a_2 \cdot t_2 \tag{3}$$

$$S_R = a_R \cdot t_R = a_R(t_1 + t_2) \tag{4}$$

In turn, the signal value of the reference path can be expressed as the sum of the signal values of paths 1 and 2 as follows,

$$S_R = a_R a_1^{-1} S_1 + a_R a_2^{-1} S_2 = u S_1 + v S_2 \tag{5}$$

where u and v represent the fraction of the affinities of the mapped probe-set (or equivalent lengths) in the reference path and paths 1 or 2 respectively. The values of u and v can be estimated from signal data.

Dividing eq. (2) with eq. (4) we get,

$$\frac{S_1}{S_R} = \frac{a_1 t_1}{a_R(t_1 + t_2)} \tag{6}$$

Combining eqs. (5) and (6), the desired equation of the Percent Spliced Index (Ψ) used in EventPointer is obtained:

$$\Psi = \frac{t_1}{t_1 + t_2} = \frac{u S_1}{S_R} = \frac{u S_1}{u S_1 + v S_2} \tag{7}$$

Note that Ψ can be directly obtained from signal values once u and v are known. This equation does not require the estimation of the affinities (difficult to predict accurately) to compute Ψ . On the contrary, it simply requires to estimate u and v from signal values using eq. (5). In the case of RNA-seq, the equivalent lengths are known a priori and hence u and v . However, using this approach has an advantage: the estimates of these lengths can accommodate the potential lack of uniformity of the reads.

Note that u and v must be positive, similar between them and close to one. The first affirmation is trivial since affinity values (or equivalent lengths) are always positive. In microarrays, probe-sets are composed by several probes and their overall affinity are expected to be similar to each other, since these affinities are a median of the average of the affinities of the probes that build up them. Therefore $a_1 \approx a_2 \approx a_R$, and $u \approx v \approx 1$. A similar reasoning can be applied to RNA-seq, if using coverage instead of read counts, since the coverage of the reference path is expected to be close to the sum of the coverages of paths 1 and 2.

These two fractions can be estimated from eq. (5) by using non-negative least squares as follows:

$$\begin{aligned} & \min \|Ax - b\|_2 \\ & s.t. x \geq 0, x \in \mathbb{R}^2, A \in \mathbb{R}^{m \times n} \end{aligned} \tag{8}$$

where,

$$\begin{aligned}
 A &= \begin{bmatrix} \text{Signal P1} & \text{Signal P2} \\ \lambda & -\lambda \\ \lambda & 0 \\ 0 & \lambda \end{bmatrix}; x = \begin{bmatrix} u \\ v \end{bmatrix}; b \\
 &= \begin{bmatrix} \text{Signal R} \\ 0 \\ \lambda \\ \lambda \end{bmatrix} \tag{9}
 \end{aligned}$$

The penalty factor λ is added to force the equation to fulfill the previous considerations: u and v must be similar and close to 1. In our results, we found that the estimates were not sensitive to the specific value of λ if there is differential alternative splicing. If the relative usage of both paths is similar and therefore, Ψ is constant, the results are more sensitive to the value of λ . This fact is shown in Fig. 4: the correlation is much better for events that show variability in the relative expression of both paths.

The residuals of this model can be used to test if the additive model of eqs. 2, 3 and 4 holds. We computed the relative error of the residuals as follows:

$$\epsilon = \frac{\|(u \cdot \text{Signal P1} + v \cdot \text{Signal P2}) - \text{Signal R}\|_2}{\|\text{Signal R}\|_2}$$

If the relative error is large, the additive model does not fit the data and, therefore, the estimates are expected to be less reliable. In order to test this, we divided the events according to the relative error. The events with top 50% relative error have FDR two times larger than the bottom 50% as shown in Additional file 1: Table S10.

Statistical analysis

The comparison and analysis of the profiling data was done using a linear model. The design matrix was built considering both the cell line and treatment with CX4945 as factors. The interaction between cell line type and treatment was not considered.

The selected contrasts test for the difference between control samples (DMSO) and drug exposed ones (CX4945) controlling for the cell-type. The complete experimental design in the form of design and contrast matrices is included in Additional file 1: Table S9.

EventPointer includes several statistical methods to state the significance of an event. In this experiment, the events are considered to be statistically significant if there is a change in the expression of the isoforms associated to each of the alternative paths, this change occurs in opposite direction, i.e. opposite signs for the fold changes and the summarized p.value is significant (p value < 0.001).

In order to compare the arrays with different sequencing depths, we subsampled the RNAseq data to 30 and 10 million reads and rerun the whole pipeline with these

data. The FDR for 30 million reads was better than using arrays. On the contrary, using 10 million reads the FDR was worse than using arrays. Interpolating both data, the FDR for arrays is similar to a depth of 20 million reads.

Filters used to include the events

For arrays, the signal of the probe-sets interrogating each of the alternative paths involved in a splicing event, must be expressed more than a certain threshold in at least one sample. This threshold is the 25% quantile of the expression of the signal in the reference paths for all the events included in the array. For RNAseq, the edges of the splicing graph (junction reads) are included only if their expression is at least 2 FPKM in at least one sample (SGSeq defaults).

Matching of the events using different technologies

Let's assume that A_R and A_M are, possibly non-contiguous, regions of the genome that correspond to path A using either technology (A_R for RNA-seq and A_M for HTA). B_R and B_M have a similar description for path B and R_R and R_M for the reference path in each technology. Two events are considered to match if any of the following two expressions is true:

$$((A_R \subset A_M) | (A_M \subset A_R)) ((B_R \subset B_M) | (B_M \subset B_R)) \times ((R_R \cap R_M) \neq \emptyset) \tag{10}$$

$$((A_R \subset B_M) | (B_M \subset A_R)) ((B_R \subset A_M) | (A_M \subset B_R)) \times ((R_R \cap R_M) \neq \emptyset) \tag{11}$$

In these expressions, $(x \subset y)$ is true if the genomic region x is a subset of the genomic region y (the nucleotide sequence of x is a substring of the nucleotide sequence in y). Besides, the operators “|” and “&” and the logical OR and AND operations. If $(x \subset y) | (y \subset x)$, then one of the regions is contained in the other are considered to be “compatible”. On the other hand, $(x \cap y) \neq \emptyset$ means that regions x and y overlap in the genome. Therefore, the first expression is true if both paths A_R and A_H are compatible, B_R and B_M are compatible and R_R and R_M overlap. The second expression is true if path A_R and path B_M are compatible and also path A_M and B_R are compatible and, again, and R_R and R_M overlap.

Within an event, the longer path in the transcriptome is assigned the name “A” and the other the B. The second eq. (11) takes into account that, in some few cases, the name of the paths can be switched in both technologies.

PCR validation

For each splicing event, an end-point PCR was run using primers designed in the exons that flank the event of interest. RNA was retro-transcribed and the PCR was performed and analyzed as previously described [40]. Primers used are shown in Additional file 2.

Endnotes

¹Average sequencing depth for RNA-seq was approx. 98 million (paired-end, stranded protocol), yielding on average approx. 49 million fragments per sample.

Additional files

Additional file 1: Vignette on the comparison based on expression analysis. **Figure S1** to **S9**. Experiment design and contrast matrices.

Table S10. (PDF 7766 kb)

Additional file 2: Excel file with characteristics of the validated events and PCR primers. (XLSX 27 kb)

Additional file 3: Compressed file including the GTF files generated for the matched events. (ZIP 15 kb)

Abbreviations

AS: Alternative Splicing; CDF: Chip Definition File; FDR: False Discovery Rate; GTF: Gene Transfer Format; HTA: Human Transcriptome Array; NGS: Next-Generation Sequencing; PCR: Polymerase Chain Reaction; TNBC: Triple-Negative Breast Cancer

Acknowledgements

The authors are grateful to Francisco J. Planes, Iñigo Apaolaza, Juan Ferrer and Xabier Cendoya for their comments on the preparation of this manuscript.

Funding

The work performed and described was funded by Celgene Research SL, part of Celgene Corporation. Author FC was partially supported by a Basque Government predoctoral Grant [PRE_2016_1_0194]. LMM and RP were partially funded by Spanish Ministry of Economy and Innovation and Fondo de Investigación Sanitaria-Fondo Europeo de Desarrollo Regional [P114/00806, P116/01821], CIBERONC and AECC Scientific Foundation [GCB14-2170]. AR was partially supported by the Provincial Council of Gipuzkoa through the MINEDRUG project. Celgene Research SL contributed in the design of the study and collection, analysis and interpretation of data and in writing the manuscript. The rest of the funding body had no role in the aforementioned steps.

Availability of data and materials

EventPointer for both RNA-seq and microarrays is available at Bioconductor [<https://bioconductor.org/packages/release/bioc/html/EventPointer.html>]. All the RNA-seq and microarray data are available in a SuperSeries at Gene Expression Omnibus, accession number GSE104974.

Authors' contributions

Conception and design: JPR, MOE, MT, AR. Development of methodology: JPR, AR, AM. Acquisition of data (provided cell-lines treatment, provided sequencing and hybridization, performed PCR validations, provided facilities, etc.): SC, FJM, RP, MT. Analysis and interpretation of data (e.g., statistical analysis, biostatistics, computational analysis): JPR, AR, RL, MOE. Writing, review, and/or revision of the manuscript: JPR, MOE, AM, LM, RL, RP, AR, MT. Administrative, technical, or material support (i.e., reporting or organizing data, constructing vignettes): FC, JPR, AM, AR. Study supervision: RL, AR, MT. All authors read and approved the final manuscript.

Ethics approval and consent to participate

Not applicable as cell lines were obtained from commercially available suppliers. Cell lines MDA-MB-231 and MDA-MB-468 were obtained from ATCC (Manassas, VA) with identification numbers HTB-26 and HT-132 respectively and cell line SUM149 was purchased from Asterand plc (Detroit, MI).

Consent for publication

Not applicable.

Competing interests

A. Rubio, J.P. Romero and F. Carazo are being funded by Affymetrix in an independent project. The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹CEIT and Tecnun, University of Navarra, Parque Tecnológico de San Sebastián, Paseo Mikeletegi 48, 20009 San Sebastián, Gipuzkoa, Spain. ²Celgene Institute for Translational Research Europe, Celgene Corporation, Parque Científico y Tecnológico Cartuja 93, Centro de Empresas Pabellón de Italia, Isaac Newton, 4, E-41092 Seville, Spain. ³Program in Solid Tumors and Biomarkers, CIMBA, University of Navarra, Avda. Pío XII, 55, E-31008 Pamplona, Navarra, Spain. ⁴Department of Histology and Pathology, University of Navarra, Campus Universitario, 31009 Pamplona, Navarra, Spain. ⁵IdiSNA, Navarra Institute for Health Research, Recinto de Complejo Hospitalario de Navarra, Irunlarrea 3, 31008 Pamplona, Navarra, Spain. ⁶Department of Biochemistry and Genetics, University of Navarra, Campus Universitario, 31009 Pamplona, Navarra, Spain. ⁷CIBERONC, Centro de Investigación Biomédica en Red, Instituto de Salud Carlos III, Calle Monforte de Lemos 3-5, Pabellón 11. Planta 0, 28029 Madrid, Spain.

Received: 22 March 2018 Accepted: 17 September 2018

Published online: 25 September 2018

References

- Pajares MJ, Ezponda T, Catena R, Calvo A, Pio R, Montuenga LM. Alternative splicing: an emerging topic in molecular and clinical oncology. *The lancet oncology*. 2007;8(4):349–57.
- Sveen A, Kilpinen S, Ruusulehto A, Lothe R, Skotheim R. Aberrant RNA splicing in cancer; expression changes and driver mutations of splicing factor genes. *Oncogene*. 2015;35:2413–27.
- Gardina PJ, Clark TA, Shimada B, Staples MK, Yang Q, Veitch J, Schweitzer A, Awad T, Sugnet C, Dee S, Davies C, Williams A, Turpaz Y. Alternative splicing and differential gene expression in colon cancer detected by a whole genome exon array. *BMC Genomics*. 2006;7:325.
- Kwan T, Benovoy D, Dias C, Gurd S, Serre D, Zuzan H, Clark TA, Schweitzer A, Staples MK, Wang H, Blume JE, Hudson TJ, Sladek R, Majewski J. Heritability of alternative splicing in the human genome. *Genome Res*. 2007; 17:1210–8.
- Yeo GW, Xu X, Liang TY, Muotri AR, Carson CT, Coufal NG, Gage FH. Alternative splicing events identified in human embryonic stem cells and neural progenitors. *PLoS Comput Biol*. 2007;3:e196.
- Shen S, Warzecha CC, Carstens RP, Xing Y. MADS+: discovery of differential splicing events from Affymetrix exon junction array data. *Bioinformatics*. 2010;26:268–9.
- De Miguel FJ, Sharma RD, Pajares MJ, Montuenga LM, Rubio A, Pio R. Identification of alternative splicing events regulated by the oncogenic factor SRSF1 in lung cancer. *Cancer Res*. 2014;74:1105–15.
- Romero JP, Muniategui A, De Miguel FJ, Aramburu A, De Miguel F. EventPointer: an effective identification of alternative splicing events using junction arrays background keyword. *BMC Genomics*. 2016;17(1):467.
- Seok J, Xu W, Davis RW, Xiao W. RASA: robust alternative splicing analysis for human Transcriptome arrays. *Sci Rep*. 2015;5:11917.
- Sood S, Szkop KJ, Nakhuda A, Gallagher IJ, Murie C, Brogan RJ, Kaprio J, Kainulainen H, Atherton PJ, Kujala UM, Gustafsson T, Larsson O, Timmons JA. iGEMS: an integrated model for identification of alternative exon usage events. *Nucleic Acids Res*. 2016;44(11):e109.

11. Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet.* 2009;10:57–63.
12. Wang C, Gong B, Bushel PR, Thierry-Mieg J, Thierry-Mieg D, Xu J, Fang H, Hong H, Shen J, Su Z, Meehan J, Li X, Yang L, Li H, Labaj PP, Kreil DP, Megherbi D, Gaj S, Caiment F, van Delft J, Kleinjans J, Scherer A, Devanarayan V, Wang J, Yang Y, Qian H-R, Lancashire LJ, Bessarabova M, Nikolsky Y, Furlanello C, et al. The concordance between RNA-seq and microarray data depends on chemical treatment and transcript abundance. *Nat Biotechnol.* 2014;32(9):926.
13. Fumagalli D, Blanchet-Cohen A, Brown D, Desmedt C, Gacquer D, Michiels S, Rothé F, Majaj S, Salgado R, Larsimont D, Ignatiadis M, Maetens M, Piccart M, Detours V, Sotiriou C, Haibe-Kains B. Transfer of clinically relevant gene expression signatures in breast cancer: from Affymetrix microarray to Illumina RNA-sequencing technology. *BMC Genomics.* 2014;15:1008.
14. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol.* 2016;34:525–7.
15. Garber M, Grabherr M, Guttman M, Trapnell C. Computational methods for transcriptome annotation and quantification using RNA-seq. *Nat Methods.* 2011;8:469–77.
16. Anders S, Reyes A, Huber W. Detecting differential usage of exons from RNA-seq data. *Genome Res.* 2012;22:2008–17.
17. Law CW, Chen Y, Shi W, Smyth GK. Voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.* 2014;15:R29.
18. Steijger T, Abril JF, Engström PG, Kokocinski F, Hubbard TJ, Guigó R, Harrow J, Bertone P, Consortium R, et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Methods.* 2013;10:1177–84.
19. Soneson C, Matthes KL, Nowicka M, Law CW, Robinson MD. Isoform prefiltering improves performance of count-based methods for analysis of differential transcript usage. *Genome Biol.* 2016;17:12.
20. Black DL. Mechanisms of alternative pre-messenger RNA splicing. *Annu Rev Biochem.* 2003;72:291–336.
21. Shen S, Park JW, Lu Z, 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.* 2014;111(51):E5593–601.
22. Kahles A, Ong CS, Zhong Y, Ratsch G. SplAdder: identification, quantification and testing of alternative splicing events from RNA-Seq data. *Bioinformatics.* 2016;32(12):1840–7.
23. Rogers MF, Thomas J, Reddy AS, Ben-Hur A. SpliceGrapher: detecting patterns of alternative splicing from RNA-Seq data in the context of gene models and EST data. *Genome Biol.* 2012;13:R4.
24. Goldstein LD, Cao Y, Pau G, Lawrence M, Wu TD, Seshagiri S, Gentleman R. Prediction and quantification of splice events from RNA-Seq data. *PLoS One.* 2016;11:e0156132.
25. Zimmermann K, Jentsch M, Rasche A, Hummel M, Leser U. Algorithms for differential splicing detection using exon arrays: a comparative assessment. *BMC Genomics.* 2015;16(1):136.
26. Emig D, Salomonis N, Baumbach J, Lengauer T, Conklin BR, Albrecht M. AltAnalyze and DomainGraph: analyzing and visualizing exon expression data. *Nucleic Acids Res.* 2010;38(SUPPL. 2):W755–62.
27. Siddiqui-Jain A, Drygin D, Streiner N, Chua P, Pierre F, O'Brien SE, Bliesath J, Omori M, Huser N, Ho C, et al. CX-4945, an orally bioavailable selective inhibitor of protein kinase CK2, inhibits prosurvival and angiogenic signaling and exhibits antitumor efficacy. *Cancer Res.* 2010;70:10288–98.
28. Chon HJ, Bae KJ, Lee Y, Kim J. The casein kinase 2 inhibitor, CX-4945, as an anti-cancer drug in treatment of human hematological malignancies. *Front Pharmacol.* 2015;6:70.
29. Kim H, Choi K, Kang H, Lee S-Y, Chi S-W, Lee M-S, Song J, Im D, Choi Y, Cho S. Identification of a novel function of CX-4945 as a splicing regulator. *PLoS One.* 2014;9:e94978.
30. Irizarry RA, Hobbs B, Collin F, Beazer-Barclay YD, Antonellis KJ, Scherf U, Speed TP. Exploration, normalization, and summaries of high density oligonucleotide array probe level data. *Biostatistics.* 2003;4:249–64.
31. Dai M, Wang P, Boyd AD, Kostov G, Athey B, Jones EG, Bunney WE, Myers RM, Speed TP, Akil H, Watson SJ, Meng F. Evolving gene/transcript definitions significantly alter the interpretation of GeneChip data. *Nucleic Acids Res.* 2005;33:e175.
32. Raghavachari N, Barb J, Yang Y, Liu P, Woodhouse K, Levy D, O'Donnell CJ, Munson PJ, Kato GJ. A systematic comparison and evaluation of high density exon arrays and RNA-seq technology used to unravel the peripheral blood transcriptome of sickle cell disease. *BMC Med Genet.* 2012;5(1):28.
33. Smyth GK. *Limma: linear models for microarray data.* In *Bioinformatics and computational biology solutions using R and Bioconductor.* New York, NY: Springer. 2005:397–420.
34. Wang ET, Sandberg R, Luo S, Khrebukova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP, Burge CB. Alternative isoform regulation in human tissue transcriptomes. *Nature.* 2008;456(November):470–6.
35. Benjamini Y, Hochberg Y, Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc B.* 1995;57:289–300.
36. Thorvaldsdóttir H, Robinson JT, Mesirov JP. Integrative genomics viewer (IGV): high-performance genomics data visualization and exploration. *Brief Bioinform.* 2013;14:178–92.
37. Carazo F, Romero JP, Rubio A. Upstream analysis of alternative splicing: a review of computational approaches to predict context-dependent splicing factors. *Brief Bioinform.* 2018.
38. Nazarov PV, Muller A, Kaoma T, Nicot N, Maximo C, Birembaut P, Tran NL, Dittmar G, Vallar L. RNA sequencing and transcriptome arrays analyses show opposing results for alternative splicing in patient derived samples. *BMC Genomics.* 2017;18:443.
39. Heber S, Alekseyev M, Sze S-H, Tang H, Pevzner PA. Splicing graphs and EST assembly problem. *Bioinformatics.* 2002;18(Suppl 1):S181–8.
40. de Miguel FJ, Pajares MJ, Martinez-Teroba E, Ajona D, Morales X, Sharma RD, Pardo FJ, Rouzaut A, Rubio A, Montuenga LM, Pio R. A large-scale analysis of alternative splicing reveals a key role of QKI in lung cancer. *Mol Oncol.* 2016;10:1437–49.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

