

RNA sequencing data: hitchhiker's guide to expression analysis

Koen Van den Berge*¹, Katharina Hembach*², Charlotte Sonesson*^{2,3}, Simone Tiberi*², Lieven Clement^{†1}, Michael I. Love^{‡4}, Rob Patro^{‡5}, Mark D. Robinson^{‡2#}

5 ¹ Bioinformatics Institute and Department of Applied Mathematics, Computer Science and Statistics, Ghent University, 9000 Ghent, Belgium

² Institute of Molecular Life Sciences and SIB Swiss Institute of Bioinformatics, University of Zurich, 8057 Zurich, Switzerland

10 ³ current affiliation: Friedrich Miescher Institute for Biomedical Research and SIB Swiss Institute of Bioinformatics, 4058 Basel, Switzerland

⁴ Department of Biostatistics and Department of Genetics, UNC-Chapel Hill, Chapel Hill, NC, USA

⁵ Department of Computer Science, Stony Brook University, NY, USA

* these authors contributed equally

‡ these authors contributed equally

15 # correspondence to: mark.robinson@imls.uzh.ch

Abstract

Gene expression is the fundamental level at which the result of various genetic and regulatory programs are observable. The measurement of transcriptome-wide gene expression has convincingly switched from microarrays to sequencing in a matter of years. RNA sequencing (RNA-seq) provides a quantitative and open system for profiling transcriptional outcomes on a large scale and therefore facilitates a large diversity of applications, including basic science studies, but also agricultural or clinical situations. In the past 10 years or so, much has been learned about the characteristics of the RNA-seq datasets as well as the performance of the myriad of methods developed. In this review, we give an overall view of the developments in RNA-seq data analysis, including experimental design, with an explicit focus on quantification of gene expression and statistical approaches for differential expression. We also highlight emerging data types, such as single-cell RNA-seq and gene expression profiling using long-read technologies.

<p>Abbreviations:</p> <p>DE - differential expression / differentially expressed</p> <p>DGE - differential gene expression</p> <p>DTE - differential transcript expression</p> <p>DTU - differential transcript usage</p> <p>GLM - generalized linear model</p> <p>DS - differential splicing</p> <p>LRT - likelihood ratio test</p> <p>TPM - transcripts per million</p> <p>EM - expectation maximization</p>	<p>SJ - splice junction</p> <p>bp - base pairs</p> <p>NB - negative binomial</p> <p>ML - maximum likelihood</p> <p>MM - method of moments</p> <p>FDR - false discovery rate</p> <p>VB - variational Bayes</p> <p>LRTS - long-read transcript sequencing</p> <p>APL - approximate profile likelihood</p>
--	---

Introduction: overview of the RNA sequencing assay

30 *“After that it gets a bit complicated, and there’s all sort of stuff going on in dimensions thirteen to twenty-two that you really wouldn’t want to know about. All you really need to know for the moment is that the universe is a lot more complicated than you might think, even if you start from a position of thinking it’s pretty damn complicated in the first place.” - from Mostly Harmless by Douglas Adams*

35 Molecular biologists are using gene expression studies to get a snapshot of the set of RNA molecules present in a biological system, which ultimately dictates what cells are doing or are capable of. The original RNA sequencing (RNA-seq) protocols, describing the sequencing of complementary DNA (cDNA) fragments on a large scale from a population of cells, were published over 10 years ago [1–5]. Since then, the system has been optimized for different types and qualities of starting material, as well as different research questions, and many distinct and mature protocols are available.

40

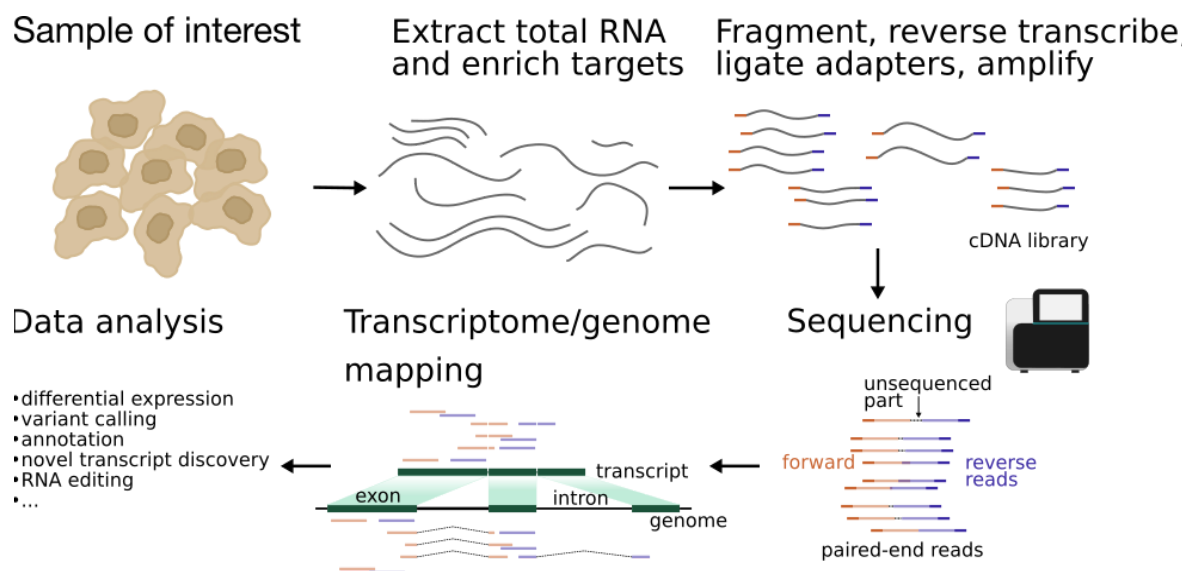


Figure 1: Overview of the experimental steps in a RNA-seq protocol. The cDNA library is generated from isolated RNA targets, sequenced and the reads are mapped against a reference genome or transcriptome. Downstream data analysis depends on the goal of the experiment and can include, among other things, assessing differential expression, variant calling or genome annotation.

A basic overview of the main steps in a standard RNA-seq experiment is given in Figure 1. The first step is the extraction and purification of RNA from a sample of interest followed by an enrichment of target RNAs. Most commonly used is poly(A) capture, to select for polyadenylated RNAs, or ribosomal depletion, to deplete ribosomal and transfer RNAs that are highly abundant in a cell (approximately 95% of total RNA) [6] and are usually not of primary interest [7]. The selected RNAs are then chemically or enzymatically fragmented to molecules of appropriate

45

size (e.g., Illumina TruSeq: 300-500 bp). Current dominant systems (e.g., Illumina) only sequence DNA; single-stranded target RNAs are thus reverse transcribed to cDNA (first-strand), the RNA is degraded and the first-strand cDNA is complemented to a double strand. Adapter sequences are either ligated to the 3' and 5' end of the double-stranded cDNA or used as primers in the reverse transcription reaction. The final cDNA library consists of cDNA inserts flanked by an adapter sequence on each end. In the last step, the cDNA library is amplified by polymerase chain reaction (PCR) using parts of the adapter sequences as primers.

For Illumina sequencing, the library is loaded onto a flow cell where the cDNAs bind to short oligonucleotides complementary to the adapter sequence. Bridge amplification creates dense "clonal clusters" of each cDNA loaded [8]. The sequence of each cluster is determined by a process called sequencing by synthesis [9]: single-stranded templates are read as the complementary strand is generated. A single fluorescently labeled deoxynucleoside triphosphate (dNTP) is added in each step. The label acts as a terminator and prevents the incorporation of more than one dNTP at the same time. After the fluorescent label has been imaged, it is enzymatically cleaved and the next dNTP can bind to the chain. Base calls are inferred directly from the measured fluorescent signal intensity.

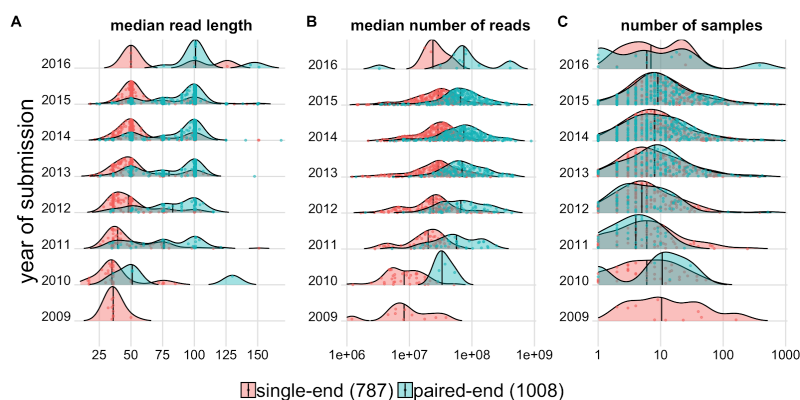


Figure 2: Ridge plot showing the progression of read length, depth and sample size in Sequence Read Archive (SRA) projects from the recount package [10]. The projects are separated by the submission year of the biosample. (A) Median read length of all samples per project and year. The color of the ridge indicates the library type of the project and the black vertical line marks the median. There are 787 single-end and 1008 paired-end projects. Each point represents one project. (B) Median number of reads across all samples per project and year. (C) Number of samples in each project.

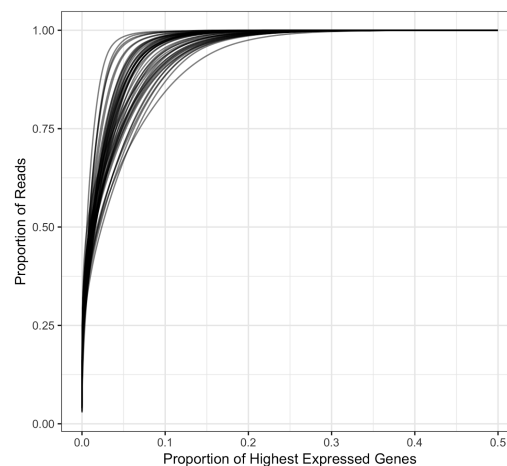
cDNA libraries can be sequenced in one of two modes: single-end or paired-end. In single-end mode, only one end of the cDNA insert is sequenced whereas in paired-end mode, both ends are sequenced, yielding two reads in opposite orientation, one from each end.

There are protocols for unstranded and stranded RNA-seq [11,12], where the latter preserves information about the coding strand of each fragment, which is useful in compact genomes or with expressed RNAs that originate from opposite strands of the same genomic locus. One possibility to construct a stranded library is to use dUTPs in the generation of the second strand

cDNA and to degrade the dUTP labeled cDNA before PCR amplification [13]. Other protocols use alternative adapters to distinguish between 5' and 3' ends of the RNA [14].

75 RNA-seq has greatly evolved over time, with early experiments having reads around 35bp long and modern (Illumina-based) experiments typically employing 50bp (single-end) or 100bp (paired-end) reads (Figure 2A). Most RNA-seq experiments comprise between 10 and 100 million reads, with a trend towards deeper sequencing over time (Figure 2B). The number of samples per project remained constant over the years with a median of around 8 samples (Figure 2C). Rapid enhancements in sequencing technology have enabled not only longer read lengths (e.g., Illumina MiSeq at 250-300bp) and much higher throughput for the same cost, but
80 also much lower amounts of required starting material. Meanwhile, third generation technologies, such as Pacific Biosciences (PacBio) and Oxford Nanopore Technologies (ONT), allow the sequencing of single molecules and have now been used for sequencing full-length transcripts on a transcriptome-wide scale [15]. Further developments are summarized below in the section 'Long Read Transcriptome Sequencing'. In addition, *single cell* RNA-seq is a rapidly
85 emerging technique that can be used to sequence the sparse transcriptome of individual cells. Some of the early developments in this area are captured in the section 'Single Cell Transcriptome Sequencing'.

Figure 3: Cumulative proportion of reads amongst the top expressed genes. The X-axis orders genes according to the total number of reads they attract and the Y-axis displays the cumulative fraction of total reads. Each line represents a single RNA-seq sample. Counts were downloaded from recount2 [10] and 50 samples were randomly selected from Accession SRP060416.



Design aspects of RNA-seq

90 The basics of experimental design apply equally for RNA-seq as they do for other scientific experiments (e.g., see [16]). For example, whether the desired experiment is a simple two-group design or a full factorial design, considerations towards *randomizing* experimental units to treatments and to avoid confounding factors (e.g., via *blocking* over batches) should be respected. In an experiment that must be run in multiple batches (e.g., limited number of samples per run), it is critical to represent every experimental condition in each batch, so that,

95 when comparing conditions, differences within a batch can be averaged over in the statistical modeling.

Specific aspects to be considered while designing an RNA-seq experiment include the number of replicates and the depth of sequencing. Ultimately in modern genomic experiments, where resources (e.g., material from subjects) are scarce and the RNA-seq experiment in itself is a hypothesis-generating tool, the first driver of sample size is budget. For better or worse, many RNA-seq studies use as few as three replicates per condition (Figure 2C), near to the minimum required to do any statistical analysis.

Sample size calculators can compute the required number of samples to achieve a user defined power for detecting differential gene expression [17–20]. However, the user must define many parameters, such as the expected alignment rate, the desired power, the significance level and the log-fold-change of differentially expressed genes. A recent study came to the conclusion that the recommended sample sizes vary from tool to tool, even when estimates from pilot data are available [21]. Another issue with sample size calculators is that it might not be obvious how to precisely define the outcome: do we want to find as many DE genes as possible? Do we want a certain power for the lowly expressed genes or the highly expressed ones? In many cases, RNA-seq experiments are exploratory and thus a means to further experimentation.

Nonetheless, there is a tradeoff between the number of samples and the sequencing depth in terms of discovery performance. Increasing the number of reads might seem always beneficial, but a large proportion of the reads originate from a small pool of highly expressed genes and, there is effectively no signal saturation. Figure 3 highlights that more than 80% of reads are attributed to the 10% most expressed genes, acknowledging that transcript length also plays a role [22]. An increased number of reads only marginally increases the coverage of lowly expressed genes and therefore the statistical power to detect differential expression (DE) does not improve considerably, especially if the experiment already comprises ~10 million reads per sample [23]. In most cases, the budget is better spent on replicates. For example, Schurch et al. show that a higher number of replicates is required to identify DE of genes with low fold change and ideally at least 6 replicates per condition should be used [24].

There are options for additional capture of genes with low expression, but this implies additional labor and cost. In targeted RNA-seq (RNA CaptureSeq), specific regions are first captured by probes that are complementary to the region of interest and these selected regions are prepared and sequenced [25,26]. After capture, the quantitative nature of the assay is maintained [25]; such capture is especially useful in degraded samples (e.g., patient material stored in paraffin blocks) where the poly(A) tails may not be present.

RNA-seq applications

130 Clearly, the popularity of RNA-seq is driven by its large number of applications. One obvious application area is genome annotation. Even the well studied transcriptomes of humans or

135 model organisms such as mice, zebrafish or fruit flies are not complete. Thus, transcriptomics is used to annotate novel transcriptional events, such as exon skipping, alternative 3' acceptor or 5' donor sites or intron retention and to understand their usage in normal, developmental or pathological conditions. Transcriptomic studies identified previously unknown phenomena, such as microexons [27], cryptic exons [28], "skiptic" exons [29], circular RNAs [30], enhancer RNAs [31], fusion genes [32] and so-called epi-transcriptomics involving RNA base modifications [33].

140 One of the main application areas is that of gene regulation. RNA-seq enables the comparison of gene / transcript / exon expression between different tissues, cell types, genotypes, stimulation conditions, time points, disease states, growth conditions and so on. Ultimately, the goal of such comparisons is to identify the set of genes that *change* in expression, hopefully leading to some understanding of the molecular pathways that are used or altered or the regulatory components that are utilized.

145 Gene expression has been employed for the molecular sub-classification of cancer since the early days of microarrays [34]. RNA-seq offers this same capacity, but at higher resolution and can include, for example, categorization by splicing [35]. Not surprisingly, there is considerable interest in using RNA-seq in clinical applications, to augment or corroborate the information that genome sequencing gives [36,37].

150 Among others, further applications include spatial transcriptomics, where cellular positional information is maintained in the preparation of cDNA fragments [38], host-pathogen interactions via "dual RNA-seq", where the transcriptomes of both host and pathogen are simultaneously assayed [39], the analysis of genetic variation among expressed genes [40], RNA editing events [41], characterization of long non-coding RNAs [42] and meta-transcriptomics [43].

155 Despite the many use cases for bulk RNA-seq, there are applications where single cell resolution is desired, especially when studying heterogeneous tissues that consist of more than one cell type. While bulk RNA-seq can be computationally deconvoluted to estimate the composition of cells present [44], it is not possible to discover new cell types or perform cell-type-specific analyses with bulk RNA-seq and thus single cell RNA-seq opens the door to new applications.

160

Outline

165 The focus of this review is on data analysis aspects, with a view towards equipping the reader with an overall view of the computational steps involved (focus on DE), but also shining a light on various statistical and computational challenges and the range of approaches that have been proposed to address them. Not surprisingly, the review is largely geared towards Illumina-based RNA-seq data on model organisms, as that is the dominant application area. There are already excellent reviews for major application or computational areas, such as *de novo* (or reference-based) transcriptome assembly [45], allele-specific expression analyses [46], expression quantitative trait loci mapping [47], splicing [48], analysis of gene regulatory

170 networks [49] or pathway analyses [50,51]. In a large majority of applications, the overarching goal is to identify DE, whether that be at the gene, transcript or exon level. The set of DE entities provides a snapshot into the molecular underpinnings of a stimulus, a disease condition, a genetic mutation or any other perturbation being interrogated. In most cases, DE is only an intermediate (though critical) step to understanding the biological system under study.

175 The review is organized as follows. First, we discuss ‘Alignment and Quantification’, where RNA-seq reads are placed in the context of the genome and/or annotation catalogs and the relative expression level of each target is assessed. Following quantification, we split the rather broad topic of DE into ‘Basics of differential expression’, which lays the foundation for the current frameworks, and ‘Variants of differential expression’, to highlight the diverse conceptual tools available to run the discovery process. Following this, we discuss two rapidly evolving research areas, namely ‘Single Cell Transcriptome Sequencing’ and ‘Long Read Transcriptome Sequencing’, which have both experienced considerable activity in the last few years.

Alignment and Quantification

185 After an experiment has been designed and executed, the analyst is presented with files containing potentially millions to billions of short cDNA fragments. Following sufficient quality control of the sequencing reactions, *alignment* to a reference genome or (*de novo* assembled) transcriptome is one of the critical steps in translating the raw data into something quantitative.

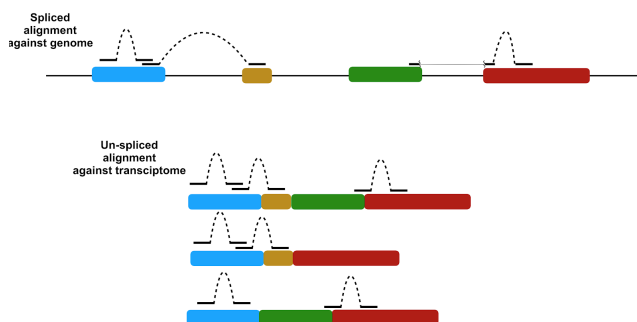


Figure 4: Illustration of spliced alignment of RNA-seq fragments to a genome (top) and direct alignment to a transcriptome (bottom). Reads are designated by thick solid lines, while dashed arcs represent the pairing relationship between paired-end reads. This illustration depicts alignment to a single 4 exon gene consisting of 3 distinct transcripts. In the spliced alignment (top), the left read of the rightmost pair is a junction-spanning alignment to the red-green exon boundary. In the direct alignment to the transcriptome (bottom), one observes how the same alignment (e.g., the alignment to the blue exon) is repeated for each transcript.

190 Because the sequenced fragments are derived from cDNA corresponding to fully (or partially) spliced transcripts, reads will often span the boundaries of splice junctions (SJs), resulting in so-called “junction-spanning” reads (Figure 4). This results in contiguous read sequences whose constituent sub-sequences may be separated by tens of thousands of nucleotides on the genome. This poses a considerable computational challenge, as the position of splice junctions in spanning reads needs to be accurately identified for a read to be properly aligned. There are two main approaches for handling spliced reads, each having its own challenges and benefits: a spliced alignment against a reference *genome* or unspliced alignment against a reference transcriptome (a database of all isoforms). A main challenge in spliced alignment against a

reference genome is the proper alignment of reads that span a SJ, especially when these junctions are not annotated *a priori*. Meanwhile, the main challenge in unspliced alignment to a transcriptome is the redundant sequence among related isoforms, which often leads to a high multi-mapping rate.

200

Spliced alignment to a reference genome

A popular solution for handling alignments of RNA-seq data is to use a splice-aware aligner. Early RNA-seq aligners, e.g., TopHat [52], make use of DNA-seq aligners, such as Bowtie [53], by first building a catalog of putative SJs to which the reads can be directly aligned.

205

More recent splice-aware alignment tools [54–64] account for read splicing directly. They also can utilize the locations of known SJs and discover previously unannotated SJs. When a read partially aligns, the annotated SJ database is consulted to check if the alignment ends prematurely as the result of the read spanning a known splice site. In this case, compatible downstream splice sites can be considered as candidate loci to align the remaining portion of the read. Even if no annotated splice site exists at the point where the alignment ends, the tool can interrogate the terminal nucleotides in the partial alignment to see if they are compatible with known canonical (or user provided) donor or acceptor sites, providing evidence that the partial alignment stops as the result of a splicing event.

210

215

One of the primary difficulties in aligning reads across SJs is that only a small portion of the read spans into one of the exons. Splice-aware aligners including STAR [55], HISAT(2) [56], Subread [57] and GMAP [58], attempt to deal with such cases by using evidence from reads that confidently align across SJs. In such strategies, new SJs are added to the index when they display “high-confidence” evidence, i.e., when multiple reads with sufficient anchoring sequence span the SJ. Subsequently, the trusted SJs are used to help align reads that start or end near the junction boundaries.

220

Unspliced alignment to a reference transcriptome

In organisms where transcriptomes are well-characterized, an alternative to splice-aware genome alignment is direct transcriptome alignment, which consists of aligning against a set of known transcripts. Since the transcript sequences are already spliced, reads should align contiguously and many of the computationally expensive steps and heuristics can be avoided. Moreover, when no reasonable quality reference genome is available for reference-based transcript assembly (e.g., when a transcriptome has been assembled *de novo*), alignment directly to the assembled transcripts is the only available option. However, transcriptome alignment induces a high degree of multi-mapping and dealing with this becomes a primary computational challenge. For example, if a gene has 3 distinct isoforms, a constitutive exon of this gene will appear 3 times in the transcriptome reference (e.g., the blue exon in Fig 4). Additionally, mapping only to annotated transcripts gives no capability to find novel splicing or expression patterns (e.g., novel exons) and it becomes difficult to assess retained introns or partial splicing; of course, it is possible to augment the transcriptome with unspliced variants.

225

230

235 The choice of genome versus transcriptome alignment is largely driven by the desired target application, and the constraints of downstream analyses.

Gene- and transcript-level quantification from RNA-seq data

240 One of the main uses of RNA-seq is to assess gene- and transcript-level abundances. Accurate abundance estimation is crucial to common downstream applications, including assessing all the notions of DE. Most commonly, abundances are estimated at the level of “genes”, but recently transcript-level abundances have also become more widely used, and there are inherent tradeoffs in choosing between the two levels of resolution.

245 Gene-level quantification consists of assigning fragments (reads or read pairs) to “genes”, where gene is often taken to represent the amalgamation of all transcripts produced from a specific strand at a specific locus [65], which typically share some exons or parts of exons. The total expression of a gene is the sum of expression of its isoforms. Any fragment arising from any isoform of a gene is assigned to the underlying gene. There are typically two paths that can be taken to obtain gene-level quantifications: direct fragment *overlap* counting of gene features, and transcript-level quantification followed by aggregation to the gene level.

250 Direct fragment counting of gene features is done by first mapping RNA-seq reads to the genome with a splice-aware aligner, and then using a tool like featureCounts [66], HTSeq [67], or the built-in capability of STAR [55], to assess how many fragments overlap each gene; the same approach can be used to quantify other disjoint genomic features, such as non-overlapping exonic segments. Even in this basic pipeline, there are many variations of how certain conditions should be handled. For example, should a fragment reside completely within a feature to be counted? If a fragment maps to multiple features, should it be discarded, counted toward each feature, or somehow partially allocated? Of course, direct fragment counting approaches exhibit desirable features: they are conceptually simple and are typically quite fast. Conversely, they suffer from various disadvantages: they have no principled way of handling multi-mapping reads (e.g., arising from paralogous genes), and they are oblivious to potentially important compositional changes that are not reflected directly in gene-level read counts (e.g., isoform switching). Additionally, since such methods assess the frequency of reads overlapping a gene, they must grapple with the concept of gene definition. For example, should a gene be considered to be the union of exons of all transcripts of the gene, or the intersection? Should intronic reads be included? Though the notion of the gene is a useful abstraction, transcripts are 260 assayed in RNA-seq, and so present a conceptually cleaner target for quantification. 265

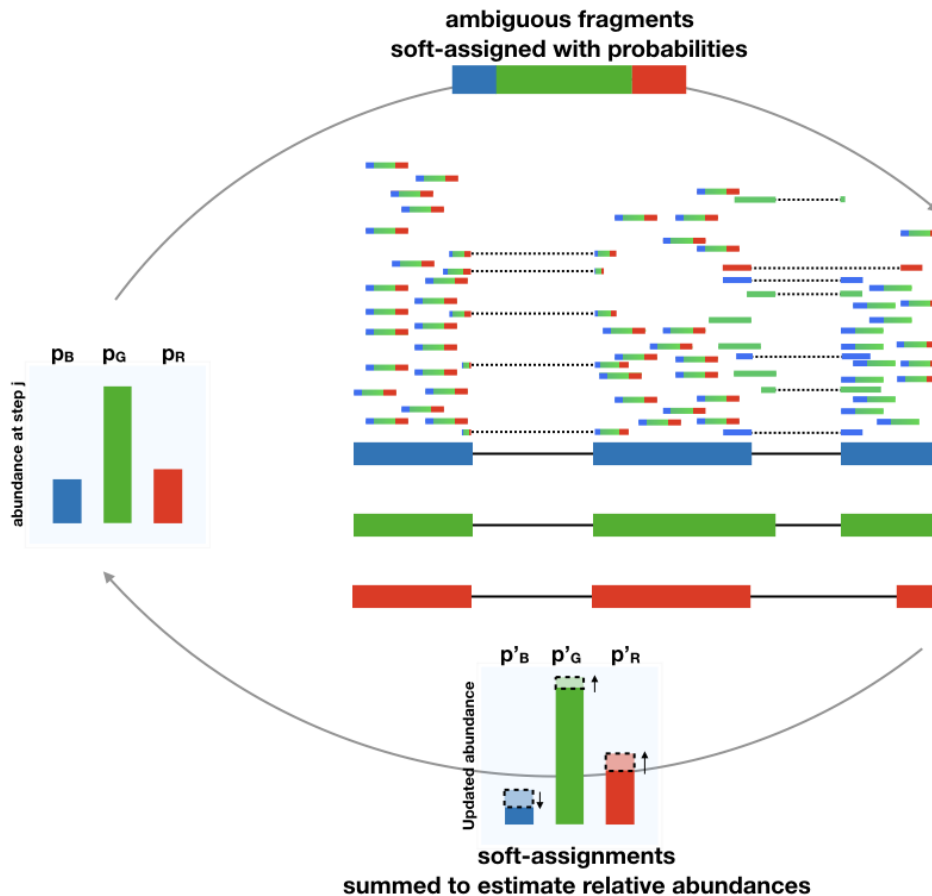


Figure 5: Illustration of alignment of various reads to a gene with 3 isoforms (B - blue; G - green; R - red). In this example, we wish to estimate the abundances of these isoforms, but the majority of reads have ambiguous origins and need to be probabilistically assigned to the transcripts (relative probabilities for each read is shown by the magnitude of the three colors). Some reads are consistent only with the B and G transcripts (colored blue and green, respectively) and a small number of reads uniquely align to a single transcript (single color). In the expectation-maximization (or related) algorithm, given the current abundance estimates, fragments are probabilistically assigned to transcripts, and then estimated abundances are updated by summarizing the (proportional) allocations over all fragments; transcript abundance estimates are determined after iterating the procedure until convergence.

Transcript-level quantification consists of the assignment of fragments to specific transcripts, which is fundamentally more challenging but has a number of advantages: it admits a clear interpretation since transcripts are what the cell expresses; it allows for improved biological resolution and allows decoding of potentially important biological changes, such as isoform switching; it is the most appropriate level to model and correct for technical biases [68–71]; it provides a proper model for handling reads that multi-map, as failing to do so can lead to systematically poor quantification for genes in gene families [72]; solving the transcript-level abundance estimation problem implies a principled solution to aggregating to gene-level estimates [73–75]. Conversely, transcript-level quantification is not without disadvantages: alternative splicing implies that many fragments are ambiguous in their origin, and they must be assigned probabilistically, necessitating the adoption of a model, which may fail to adequately

capture reality; this read ambiguity translates to additional uncertainty in the estimated transcript abundances.

Transcript quantification

280 Methods for transcript quantification are based primarily on defining a generative model of RNA-seq reads, and then trying to perform inference on this model to obtain the relevant quantities (i.e., transcript abundances); see schematic in Figure 5. There has been a tremendous amount of research aimed at solving the problem of quantifying transcript-level abundance from high-throughput sequencing data and here we describe a few major highlights
285 along the arc this research has taken.

Initial probabilistic frameworks for transcript identification and abundance estimation using EST (expressed sequence tags) data were already being developed before the wave of Illumina-based sequencing [76], but Jiang et al. [77] were among the first to attempt isoform-level abundance estimation using RNA-seq data. They define counts over exons and
290 exon junctions as arising according to a Poisson model, and view transcripts as vectors of inclusion and exclusion of these exons and junctions. By expressing the likelihood of the model parameters given the observed data, they pose a statistical model that admits efficient inference, for which they both obtain the point estimate by gradient ascent, and provide estimates of the posterior distributions of the parameters via importance sampling. This work
295 represents one of the first proper statistical formulations of the problem. However, the approach does not account for fragments that map to multiple genes, and requires annotations of transcripts in terms of the gene-transcript relationship as well as the exon and junction read inclusion matrix.

Li et al. [73,78] proposed one of the most widely-adopted generative models for transcript
300 quantification, RSEM. They define a fragment-level model of RNA-seq experiments in terms of sampling molecules from an underlying population, proportional to the product of their abundance and length, and then generate fragments from the sampled molecules. Primary quantities of interest, including the nucleotide fractions (the fraction of all sequenced nucleotides deriving from each transcript) and the transcript fractions (the fraction of all transcripts in the
305 initial population that consists of each transcript species), are estimated, and can be directly converted into popular abundance units, such as transcripts per million (TPM) or estimated counts. Notably, they propose computing the maximum likelihood (ML) estimates using an expectation-maximization (EM) algorithm (see Figure 5), and introduce a modified Gibbs sampling procedure to allow estimating credible intervals for the abundance estimates [73]. The
310 model is quite general: it works at the fragment level; it can account for numerous protocol-related aspects, including single-end and paired-end sequencing, directional vs. unstranded protocols, various coverage biases, etc. Further, the model relies only on knowing the transcript sequences and not the relationships to genes or annotations of exons and SJs. Thus, it can be easily applied to both well-characterized or newly-assembled transcriptomes.

315 One drawback of adopting a fragment-level model, however, is that each EM iteration scales in
the total number of alignments, which is indeed large in most RNA-seq experiments.

320 Instead of trying to model each fragment individually, MMSeq focuses on modeling sufficient
statistics [79,80]. Reads are categorized into equivalence classes, where two reads are
equivalent if they align to the same set of transcripts. The approach works both within and
325 across genes, and does not require the shared regions giving rise to the equivalence classes to
correspond to any known annotation (e.g., exon or SJ). MMSeq uses an EM method that works
directly over these equivalence classes, allowing efficient inference of transcript-level
abundance in this model. In addition to this ML approach, they also introduce a Gibbs sampling
330 procedure that allows estimating transcript abundances using summary statistics from samples
of the estimated posterior, which also allows the assessment of uncertainty in the transcript-level
abundance estimation and for assessing groups of transcripts with correlated posterior
estimates. It is worth noting that the underlying likelihood function of the equivalence
class-based model is not equivalent to that of the fragment-level model in RSEM, although
335 subsequent work explored other factorizations of the full fragment-level likelihood that either
preserved equality with the RSEM model while speeding up inference [81], or sacrificed equality
in an attempt to balance efficiency and fidelity [82,83]. eXpress demonstrated how
fragment-level inference could be made much more efficient by modifying the inferential
algorithm itself (i.e., online-EM), rather than the factorization of the underlying likelihood function
[84] [85].

335 Cufflinks is widely known both as a reference-guided transcript assembly algorithm and a
quantification tool [86]. Quantification is either restricted to a reference annotation, or allows new
transcripts to be identified via alignments; transcript abundances are estimated via an EM
algorithm to determine the ML estimates given the observed data. While we do not focus on
assembly methods here, it is worth mentioning that, given the close relationship between
340 transcript identification (assembly) and quantification, numerous approaches attempt to solve
both problems together, either stagewise or jointly [82,87–94].

A model similar to RSEM that jointly performs quantification and DE, together with fully
Bayesian inference, was introduced in BitSeq [95]. BitSeq focused on sampling from the
posterior distributions of transcript abundances, given the fragment alignments, giving accurate
345 estimates [96] and useful information about posterior uncertainty and posterior correlation,
which is used in the DE step [95]. To combat the heavy computational requirements, Hensman
et al. introduced a variational Bayesian (VB) approximation [97] that can be efficiently optimized.
A VB approach to the transcript abundance estimation problem was introduced in TIGAR [98],
where the VB EM algorithm was shown to outperform the standard EM algorithm. However,
350 Hensman et al. [97] introduced a novel optimization procedure called VBNG (Variational
Bayesian Natural Gradient), which is a gradient ascent algorithm that takes into account the
information geometry [99] of the underlying problem. They also suggest that EM-based methods
tend to find solutions near the boundary of the parameter space, and their quantifications are
less robust than either fully Bayesian or variational Bayesian estimates [97].

355 Many of the mentioned approaches, among others, simplify the model or improve the efficiency
of the inferential procedure, but they all rely on full alignments of each read, which can be a
computationally intensive and time consuming process. Recently, a number of new methods
bypass the alignment step, and instead adopt *lightweight* models for quantification. Sailfish
defines the transcript abundance likelihood in terms of the constituent *k-mers* of the underlying
360 transcriptome and their abundance in the read data [100]. Since the *k-mers* are completely
known in advance, the relevant equivalence classes can be pre-computed, which reduces the
inferential problem to one of simply counting *k-mers* and performing inference via an EM
algorithm, e.g., the SQUAREM algorithm [101]. This approach increases the speed of
abundance estimation by over an order of magnitude compared to full alignment approaches.
365 Building on the idea of *k-mer*-based abundance estimation, RNA-skim takes the approach of
Sailfish even further, identifying sets of distinctive *k-mers*, known as “sigmers” [102]. Transcripts
are clustered into groups, and sigmers are identified as *k-mers* that are unique to (and indicative
of) each cluster. Quantification is then performed by counting the sigmers in the read data,
instead of all *k-mers*, and the EM algorithm is used to estimate transcript abundances from
370 sigmer equivalence class counts. While very fast, these *k-mer* based approaches do not retain
the coherence of the *k-mers* along a read, which can reduce specificity, and they cannot easily
estimate certain aspects of the generative model, like the fragment length distribution.
Addressing these shortcomings, kallisto relies on the use of pseudo-alignments to directly
compute the sufficient statistics of the equivalence-class-based model of transcript abundance
375 estimation [103]. This approach uses *k-mers* to identify the transcripts with which fragments are
compatible, but does not treat the *k-mers* independently. The pseudo-alignments can be
computed in such a way that equivalence class counts are generated without the need to
consider or compute individual fragment to transcript alignments, and this can often be achieved
by querying only a small number of the *k-mers* present in a fragment, making the approach very
380 efficient and allowing accurate estimation in the equivalence-class based model using an EM
algorithm. Salmon is another lightweight quantification approach that avoids full alignments,
although they can still be used as input [104]; it uses a two-phase algorithm for transcript
abundance estimation: an online phase using a stochastic collapsed VB inference algorithm
[105], where abundances and auxiliary parameters (e.g., GC bias parameters,
385 sequence-specific bias parameters, fragment length distribution) are estimated; updates are
then made using mini-batches of mappings. Salmon uses a lightweight mapping algorithm to
compute the likely transcripts, positions and orientations of origin of each fragment, and adopts
a *fragment-level* GC bias modeling approach [71], which reduces mis-identification of expressed
isoforms when read coverage is not uniform along the transcripts due to GC-content. In the
390 offline phase, a factorized likelihood function is optimized until parameter convergence. The
granularity of the likelihood factorization used by Salmon can be adjusted [83] in a way that
allows one to tradeoff between the fragment model of RSEM and the count-based model of
MMSeq. In the offline phase, the factorized likelihood is optimized using a VB EM algorithm [98]
or a traditional EM algorithm. Combining the ideas of efficiently determining fragment-transcript
395 compatibility with the sigmer concept of RNA-skim, Fleximer uses a new matching algorithm that
makes use of sets of sigmers to determine the likely loci of origin of reads, instead of treating

each sigmer independently [106]. A generalized suffix tree is used to organize the reference sequences and a “segment graph” that demonstrates how segments of sequence are shared among reference transcripts is used to select an informative and robust set of sigmers for
400 quantification. Reads are mapped against the reference by matching them to sigmers using a pre-computed automaton. This process produces a set of transcript equivalence classes, along with a corresponding count for each, which is sometimes referred to as transcript compatibility counts (TCC); these are used in conjunction with an EM algorithm to estimate transcript abundances.

405 Due to their vastly improved speed, ease of use, and reduced computational requirements, “alignment-free” approaches have become popular for assessing transcript and gene-level abundance using RNA-seq data. Recent benchmarks [107–110] suggest that, in addition to being fast, such methods are capable of producing accurate abundance estimates — at least to the extent that simulation-based studies, sometimes adopting the assumed generative models
410 of the quantification approaches, can be relied upon to assess such accuracy. However, work still remains to be done with respect to fast, accurate, and robust transcript-level quantification. For example, it is likely the case that the underlying models can be further improved to account for complexities in the fragmentation patterns of molecules prior to sequencing [111], to better balance robustness and sample-specific accuracy [112], or to address as-yet-uncharacterized
415 biases. Also, most of these approaches (lightweight and otherwise) assume that the annotation of transcripts to be quantified is complete. The accuracy of quantification can suffer when this is not the case, though it is possible to computationally flag transcripts whose estimates are unreliable [113].

Basics of differential expression

420 Following alignment and quantification, assessing *differential expression* (DE) from the estimated feature abundances is often the next challenge. We will first present a general context and describe the statistical frameworks and overall workflow. The starting point is a count table with rows representing features (e.g., genes) and columns representing samples (i.e., experimental units). The goal of DE is to formulate and test a statistical hypothesis for each
425 feature. Depending on the experimental design, the context and the research question, more complex analyses are often required. As such, we elaborate on further variations of the overall workflow in the ‘Variants of differential expression’ section.

The general workflow involves the following steps (see Figure 6): filtering and normalization (preprocessing), specification of the statistical model and estimation of model parameters, statistical inference on the relevant parameters and adjustment for multiple testing. We
430 introduce this general workflow from the perspective of classical models for count regression. Then, we discuss various notable deviations, including alternative estimation and inference frameworks and additional strategies to ensure robustness.

435 It is worth first considering the magnitude of the inferential problem. Typically, only a limited
 number of replicates are available (e.g., 3-5 replicates per condition). One can ponder on the
 achievable statistical power from such small sample sizes, even for a single feature, with the
 real interest lying in thousands of features simultaneously. This parallel inference challenge is
 common to various genome-scale experiments and the statistical community has contributed
 strategies to at least improve the overall performance; from this, a few themes have emerged.
 440 For example, in estimating parameters for a given feature, it is often beneficial to consider the
 information coming from the other features in the dataset [114]. In general, genomics data is ripe
 for the use of empirical Bayes methods to moderate estimates, where priors for a feature are
 derived from a suitable set of other features measured in the dataset. In addition, it has become
 clear that moderating *variance parameters* is critical and indeed much of the success of earlier
 445 parallel inference frameworks (e.g., for microarrays) can be attributed to variance moderation,
 whether this occurs in an *ad hoc* strategy [115] or in hierarchical models [116]. Other “tricks”,
 such as regularization of regression parameters or considerations for robustness, provide
 additional performance benefits. Taken together, the challenges associated with vast parallel
 inference can be greatly eased by adopting one or more of these strategies.

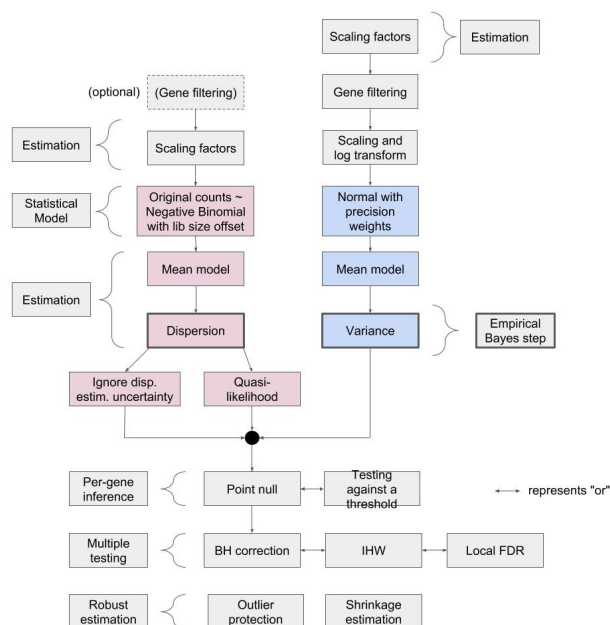


Figure 6: Schematic overview of a DE analysis for RNA-seq data. The red boxes correspond to pipelines for count-based models (e.g., edgeR, DESeq2), while the blue boxes correspond to a linear model based pipeline (e.g., limma-voom).

450

Preprocessing: filtering and normalization

The vast number of features in a typical RNA-seq experiment leads to a large multiple testing burden. However, many features are largely uninformative, e.g., features with low expression provide little evidence for DE. Therefore, filtering strategies are employed that predominantly remove uninformative features and reduce the multiple testing burden. Bourgon et al. [117] showed that filtering is valid if it is independent of the DE test statistic; thus, filtering on residual variance is invalid, while filtering on expression strength, as is commonly done, is valid.

455

The observed counts of the features cannot be directly compared across samples, since there are differences in sequencing depth across libraries. Several methods have been developed to “normalize” counts to facilitate across-sample comparisons, although in most count-based models, the counts themselves are not modified and instead scaling factors accompany the analysis. Initial attempts focused on a simple correction for sequencing depth, using the total sum of counts for each sample (i.e., the library size) as a scaling factor [3,118]. However, variation in library preparation or RNA composition between samples also contribute to between-sample variability, and should be accounted for [119]. In addition, a few highly expressed genes can largely drive the sampling of fragments, thus leading to inaccurate scaling of the counts. A popular approach is to calculate a size factor [119,120] for each sample. This can be considered to be a robust global fold change between the current sample and a (pseudo-)reference sample derived from all samples. DESeq’s median-of-ratio and edgeR’s trimmed mean of M-values (TMM, where M-values denote empirical fold changes between two samples) method are the most popular scaling approaches [121]. Both procedures assume that the majority of genes are indeed not DE, and adopt robust summarization methods to calculate the size factors (effective library sizes) in order to reduce the impact of DE genes (TMM: a trimmed weighted mean; DESeq: median of the log-expression ratios). More advanced normalization methods have since emerged, to address other technical artifacts such as GC content and transcript length effects, and to accommodate within- and between-lane normalisation, e.g., CQN: [122]; EDASEQ: [123]. Moreover, methods based on external spike-in features have been introduced to address normalization for applications where many features are DE or where the basic assumptions of conventional normalization methods are violated [124–126]. Recently, a normalization technique has been proposed for RNA-seq data with large differences between conditions that assumes similar distributions in biological replicates, while accommodating for differences between conditions [127].

The normalization size factors are built into the DE analysis workflow as *offsets* in the statistical models (see below). Notably, size factors are treated as fixed and known, while they are actually random variables that have been estimated from the data [128], and it is unclear how ignoring their associated uncertainty affects the downstream DE analysis.

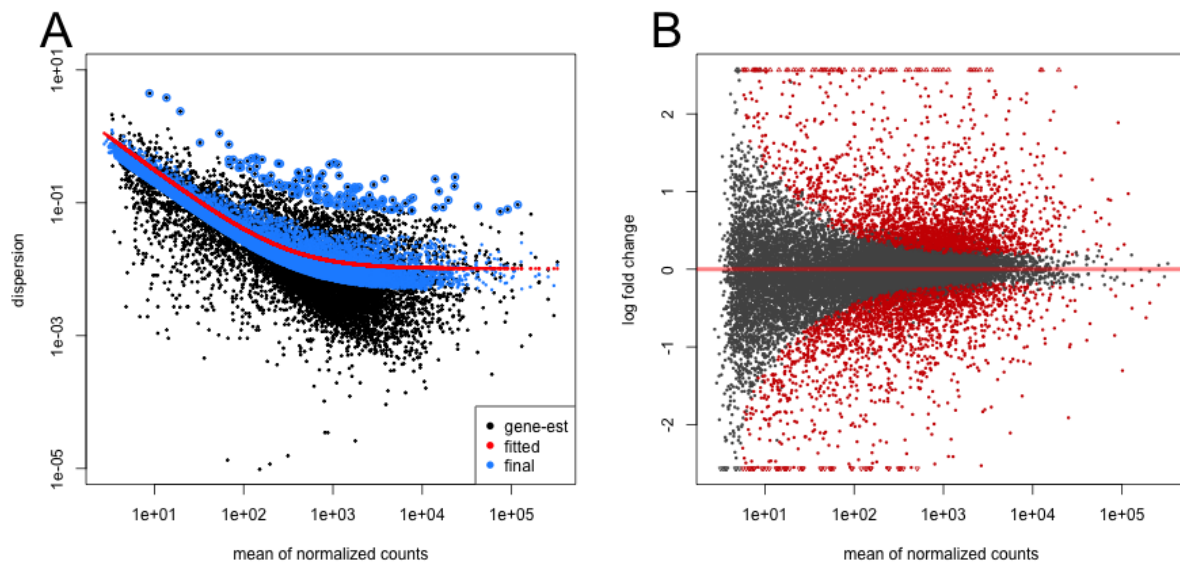


Figure 7: (A) A mean-dispersion plot of the RNA-seq experiment from [129], as processed in [130]. The dispersion smoothly decreases for genes with higher expression and eventually reaches an asymptote, which can be considered as the biological variability that is present in the dataset. (B) MA-plot of the same RNA-seq experiment. The y-axis shows the log-fold-change (M) and the x-axis shows the mean of normalized counts (A). The variability on the fold changes is higher for lowly expressed genes, which is intrinsic to count data. Red points denote DE detection according to an FDR threshold of 0.05.

Modeling and Estimation

Because of the typically small sample size, DE tools mainly implement parametric methods [120,131–134]. Initially, count data were log-transformed and linear models were used for DE analysis [4]. However, log-transformed counts suffer from *heteroscedasticity* (a systematic mean-variance trend) intrinsic to count data, rendering analysis with the standard linear model, which assumes homoscedasticity, suboptimal. In addition, fitting continuous models to (transformed) count data introduces a further approximation. Therefore discrete count distributions gained more traction in the initial frameworks.

Gene expression variability across technical replicates (i.e., resequencing the same sample), so-called *shot noise*, has been shown to approximately follow a Poisson distribution [118], for which the variance is equal to the mean. Biological replication introduces additional between-sample variability and analysis frameworks therefore resorted to one of the natural extensions, the gamma-Poisson or negative binomial (NB) distribution, which has an additional *dispersion* parameter and a quadratic mean-variance relationship:

$$Y_{fi} \sim NB(\mu_{fi}, \phi_f)$$

$$Var(Y_{fi}) = \mu_{fi} + \phi_f \mu_{fi}^2,$$

where Y_{fi} denotes the read count of feature f in sample i , ϕ_f is the dispersion for feature f and $\mu_{fi} = s_i \theta_{fi}$ represents the average expression, which is driven by the true (relative) mRNA concentration in the sample, θ_{fi} , multiplied by a normalization scaling factor, s_i ; there also
 505 exists a characteristic dispersion-mean trend in RNA-seq datasets (Figure 7A). Initial implementations focused on two-group comparisons [120,135] and were later extended to the generalized linear model (GLM) framework, an extension of classical linear models to non-Gaussian responses [136]. GLMs allow for the inclusion of multiple treatments or covariates, thus broadening the applicability. The NB GLM model can be formulated as:

$$\begin{aligned}
 510 \quad Y_{fi} &\sim NB(\mu_{fi}, \phi_f) \\
 &\log(\mu_{fi}) = \eta_{fi} \\
 &\eta_{fi} = X_i \beta_f + \log(s_i)
 \end{aligned}$$

where η_{fi} is the linear predictor, X_i denotes the design matrix, β_f represents the regression parameters and $\log(s_i)$ are scaling (normalization) offsets. Regardless of the model, the
 515 parameters θ_{fi} or, equivalently, (a linear contrast of) β_f , would represent the parameter(s) of interest for inference.

Reliable estimation of the dispersion parameter ϕ_f is non-trivial due to limited sample sizes. Traditional ML estimators for the dispersion are negatively biased [137], since they do not account for the fact that the mean is also estimated from the data. Early implementations
 520 estimated a single common dispersion parameter for all features [137], with the rationale to obtain a stable estimate by borrowing strength over all genes. However, the common dispersion assumption is unrealistic and relaxed estimation schemes were proposed, such as moderation toward a common dispersion [135], or estimation in strata of similar expression strength [120]. For example, DESeq adopts a method of moments (MM) estimator and assumes the dispersion
 525 to be a smooth function of the mean. In order to avoid too liberal inference, the dispersion is then set as the *maximum* between the smooth fit and the gene-wise MM estimate; however, while robust to outliers, this method tends to overestimate the variance and is therefore conservative [138,139]. Later approaches resorted to an approximate conditional inference scheme, the Cox-Reid adjusted profile likelihood (APL) [140], to correct for the bias in the ML
 530 estimator [136]. Again, stable estimation is provided by leveraging information across genes (Figure 8). In particular, edgeR uses a maximized weighted APL to tradeoff between gene-specific and shared dispersion estimators upon estimating the dispersion-mean trend across all genes (similar to DESeq). The weighted likelihood:

$$535 \quad APL_f(\phi_f) + G_0 APL_{sf}(\phi_f)$$

consists of the APL for a specific feature f (first component), and a shared likelihood (second component), which can be interpreted as a prior from a Bayesian perspective, thus representing an approximate empirical Bayes solution [135]. The weight given to the prior likelihood (G_0) can also be estimated from the data [141]. Analogously, Dispersion Shrinkage for Sequencing (DSS) and DESeq2 model the $\log(\phi_f)$ as a Gaussian random variable, and Bayes formula is applied to

540 generate a posterior mode for each gene [131,142]. Hyperparameters for the (Gaussian) prior are inferred from the data, using either the MM or Cox-Reid estimator across all genes.

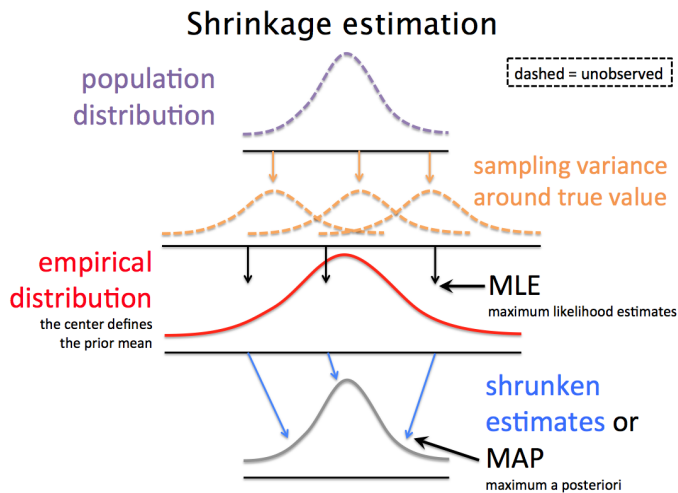


Figure 8: Empirical Bayes. In an RNA-seq experiment, the observed differences in gene expression across groups of samples with respect to within-group variance is assessed. The unobserved population distribution for the true within-group variance of each gene is depicted with a purple curve. Variances are estimated from limited sample size experiments, and so there is sampling variance in our estimate of the variance (orange curves). A ML estimate or a bias-corrected estimator for expression variance can be used (black arrow). Thousands of genes are typically observed and estimates are made for each, providing an empirical distribution of ML estimates across all genes (red curve). This empirical distribution of ML estimates can be used to determine a prior distribution for empirical Bayes analysis; the posterior distribution for the variance of each gene is calculated using Bayes formula. The maximum a posteriori (MAP), or posterior mode, represents a "shrunken estimate" (blue arrow), where the amount of shrinkage is determined by the shape of the likelihood and the width of the prior distribution.

Once dispersion estimates are available, the parameters of the mean model, β_f , can be estimated using standard algorithms for GLMs.

Statistical inference

545 After fitting a GLM to each feature, the statistical inference involves testing the null hypothesis H_0 that there is no DE between conditions, e.g., that the log-fold-change is zero ($LFC = 0$), against the alternative H_1 that the LFC differs from zero ($LFC \neq 0$). In the GLM framework, the null hypothesis can be represented either as a single regression parameter or a linear combination of parameters (contrasts), which is defined by a vector or matrix L so that $H_0: LFC = L\beta_f = 0$. Indeed, a regression parameter in a NB-GLM with canonical link function can be interpreted as a log-fold-change between groups and thus provides a measure of effect size.

550 There are multiple hypothesis tests available for GLMs with known (asymptotic) distribution under the null. Likelihood ratio tests (LRTs) compare the likelihood of a full model, upon estimating all parameters without constraints, with the likelihood of a reduced model, where one or some of the parameters are constrained according to H_0 . LRT statistics are asymptotically chi-squared distributed under H_0 and this type of test is implemented in both edgeR and DESeq2. By default, however, DESeq2 adopts a Wald test. Wald tests are attractive from a computational point of view, since they only require fitting the full model and calculating the

555

560 variance-covariance matrix of the regression coefficients. The Wald test statistic for a single
model parameter or a single contrast $W = \widehat{LFC} / \text{se}(\widehat{LFC})$ asymptotically follows a standard
normal distribution under H_0 , where $\text{se}(\cdot)$ indicates the standard error and \widehat{LFC} is the ML
estimate of LFC. From ML theory, it is known that LRTs have better properties (e.g., invariance
to transformation) than Wald tests in GLMs [143]; however, RNA-seq tools moderate dispersion
estimates and do not re-estimate them under H_0 , so it is unclear whether these benefits carry
565 over to RNA-seq data analysis in practice.

Multiple testing

The p-values obtained from the statistical inference must be corrected for multiple testing to
avoid excess false positives. While it is possible to control the probability of returning at least
one false positive in the list of detections by adopting family wise error rate (FWER) corrections,
570 this stringent form of correction is overly conservative. Indeed, when screening many thousands
of features, one is typically willing to tolerate a certain proportion of false positives in order to
obtain a larger number of true positives. The false discovery rate (FDR), which gained
significant popularity, controls the expected fraction of false positives in the detected set of
features, i.e., $\text{FDR} = E[V/\max(R, 1)]$, where V is the number of false positive rejections and R the
575 total number of detections. The FDR was introduced in the seminal paper of Benjamini &
Hochberg [144] and has become common practice in high-dimensional data analyses because
of its simplicity and solid theoretical justification. Indeed, it can be shown that the FDR is justified
under a range of dependency structures between the genes [145] and can be approached from
both frequentist and Bayesian perspectives.

580

Variations to the general workflow

There is a large and growing number of alternatives to the basic framework mentioned above:
different inference based on the same models, alternative models, more robust approaches,
different testing regimes, variations on multiple testing corrections and so on. In this section, we
summarize some of the many developments.

585 Alternative models (inference frameworks). NB count models, which underpin many DE tools,
assume a quadratic mean-variance relationship. Inference, however, may benefit from a more
flexible variance structure and, for this, other models have been proposed. One strategy uses
quasi-likelihood (QL), which requires that mean and variance are specified to be able to make
inference on the mean model parameters [146]. The QL method adopts the same mean model
590 structure as the NB, but introduces an additional overdispersion parameter, such that $\text{var}(Y_{fi}) =$
 $\psi_f (\mu_{fi} + \phi_f \mu_{fi}^2)$; ψ_f is estimated using a moderated MM estimator. QL naturally allows
(asymptotic) hypothesis tests based on t- and F-statistics, thus accommodating the uncertainty
in the estimation of the additional QL dispersion parameter. Another variation is the use of a
more flexible distribution, such as the negative binomial power distribution, which adds an
595 additional parameter [147] to the NB. Within the NB framework itself, Bayesian methods have
also been developed. A fully Bayesian approach has the benefit that various aspects of the

posterior can be reported (e.g., credible intervals) and the degree of parameter shrinkage naturally depends on the amount of information available for that gene (a trade-off between expression magnitude, dispersion and residual degrees of freedom). One of the early methods was ShrinkBayes, a fully Bayesian approach that included multiple mixture priors (e.g., Gaussian) [148,149] and where fitting was accomplished using Integrated nested Laplace approximations (INLA) [150], which avoids the Markov Chain Monte Carlo sampling. Another alternative is to remain within computationally and inferentially efficient Gaussian linear models, after suitably transforming the (normalized) count data. For example, limma-voom models log-transformed normalized counts using a linear model while adjusting for heteroskedasticity via weighted regression, where the observation weights are computed from the observed variance-mean relationship [151]. In this case, moderated t- and F-statistics are used for inference. Finally, non-parametric methods have been developed, which are more robust to outliers and do not require distributional assumptions. For example, SAMSeq [152], adopts the Wilcoxon test to assess DE between groups, and uses resampling procedures to adjust for differences in sequencing depth.

Robust LFC estimation. The standard NB workflow typically makes use of APL NB likelihood for parameter estimation, combined with empirical Bayes procedures to borrow strength across features when estimating the dispersion parameter. There are two related challenges: i) ratios of smaller counts result in more variable LFCs (Figure 7B); ii) estimation of LFC can be sensitive to outliers. This makes a ranking of genes according to LFC difficult, since lowly expressed or outlier-affected genes are likely to dominate the top list. In order to derive more robust LFC estimates, several approaches have been adopted. First, the use of “prior counts” in the numerator and denominator of the LFC; effective shrinkage is accomplished by augmenting each count with a carefully chosen value, although the optimal value may vary across datasets. Second, edgeR-robust [139], for instance, adopts an M-estimation approach by iteratively downweighting outlying observations within the GLM fitting procedure, dampening the effect of outliers on both mean and variance estimates. Alternatively, outliers can be identified and removed and/or imputed by taking advantage of the remaining data for a feature [131]. Lastly, priors can be imposed on the LFC parameters. For example, DESeq2 includes a zero-centered Gaussian prior in the NB GLM and provides the posterior mode of LFC as output [131]. The width of the prior is set conservatively, using a weighted upper quantile of the observed log-fold-changes. New alternative shrinkage estimators in DESeq2 incorporate priors with heavier tails that introduce less bias, using either a mixture of Normals [153] or a Cauchy distribution [154].

Accounting for unobserved effects. As mentioned, (G)LMs can adjust for known confounders. However, genomic data can also be affected by unknown, hence unobserved, confounders. This problem is widespread in publicly available data, which typically do not contain sufficient metadata on potential batch effects caused by lab, protocol, date, etc. Batch correction methods can leverage the parallel structure of high throughput transcriptomic data to identify unknown and unobserved systematic effects. SVA [155,156] and RUV [125], for instance, estimate *surrogate variables* through singular value decomposition on control features or on a matrix of

640 model residuals so as to avoid that the phenotypic effect of interest is captured by the surrogates. RUV also has the option to exploit information in replicate samples. The estimated surrogate variables can subsequently be included as predictors in the statistical model to adjust for the batch effects.

645 Statistical inference by testing against a threshold. The standard approach for detecting DE in RNA-seq involves a simple null hypothesis: $H_0: LFC=0$. However, statistical significance does not guarantee that the fold changes are large enough to be biologically relevant. Analysts often produce candidate gene lists by applying a threshold on the magnitude of the LFC, but the statistical properties of this approach are unclear. The FDR is a *set* property and has no interpretation when the set, post-FDR calculation, is altered [157]. To address these practical and theoretical concerns, several tools have adopted tests relative to a LFC threshold, a procedure initially proposed for microarray data [158]. This results in a composite null hypothesis, such as $H_0: |LFC|<a$. Implementations differ; for example, DESeq2 replaces the composite null with a simple null hypothesis at the boundary of the parameter space [131]; edgeR uses a modified likelihood ratio test or a quasi-likelihood F-test against a threshold [159].

655 Small-sample inference. The null distributions for Wald or LRT statistics for count models are only valid asymptotically and the number of replicates is often too low for these approximations to be fully effective, which may lead to inflated FDR. Initial implementations provided exact tests [137], but these can only be applied in simple designs. Another strategy is “small sample asymptotics”, essentially making use of higher-order approximations that are still compatible with the GLM framework [160].

660 Multiple testing. While the FDR achieves a more reasonable sensitivity-specificity tradeoff than family wise error rate correction approaches, other developments beyond simple filtering aim to further reduce the multiple testing burden. Storey's q-value, for instance, estimates the proportion of true null hypotheses from the data to increase power [161], while others adopt a data-driven *weighting* of the p-values in the FDR correction [162]. Although the FDR is deeply rooted in statistical theory, it is not guaranteed that methods will control error rates at the nominal level in real applications. NB methods, for instance, rely on the asymptotic theory, which might not hold for applications with low sample sizes. A study has suggested that co-regulation of genes induces intergene correlations, which can alter the null distribution of the statistical test [163]; local FDR approaches were introduced that empirically estimate the null distribution [164]. Other developments address issues in testing many hypotheses for every gene (e.g., multifactorial designs). The conventional approach is to control the FDR on each hypothesis, but this does not allow for straightforward prioritisation, since genes typically have a different ranking for each hypothesis. Stage-wise testing procedures can be interpreted as a generalisation of analysis of variance with post-hoc tests towards a high throughput context [165,166], thus allowing a natural ordering of the genes according to an “omnibus” (all effects of interest) test, while providing FDR control at the gene level.

675

Variants of differential expression

The previous section introduced count-based DE in general terms: each row of a count matrix is submitted to a statistical model (often by first estimating moderated variance parameters over the whole dataset) and hypothesis tests of interest are conducted, with an adjustment for multiple testing. In this section, we unravel a set of additional approaches to interrogate RNA-seq data in terms of DE.

Although it may be obvious that DE is of interest, this can manifest or be defined in multiple ways (see Schematic Figure 9). It is important to remember that while one may want to cast inferences to the gene level, measurements are made at the fragment level. We use the term differential gene expression (DGE) to refer to hypothesis testing related to the *total outcome* of an annotated gene, either by comparing accumulated transcript-per-million (TPM) estimates or by comparing raw counts while including an adjustment for average transcript length via offsets [74]. If the expression of transcripts is the feature of interest (independent of other transcripts), differential transcript expression (DTE) analyses can be conducted. Alternatively, one could be interested in whether *at least one transcript* from a gene is DE. This requires statistical testing at the transcript level and then *aggregation* to the gene level. Yet another strategy is to consider whether the relative abundance (i.e., proportions) of transcripts for a specific genomic locus changes between conditions, which is commonly referred to as differential transcript usage (DTU) or, more generally, differential splicing (DS). A surrogate for DTU, differential exon usage (DEU), is conducted on exon-level quantifications; in this case, the goal is to identify exons that deviate from proportional expression to separate differential usage from DE. Yet another alternative is to quantify and test differences at the *event level*, where reads supporting (or not supporting) an event (e.g., inclusion of a cassette exon) are summarized and compared [167].

There is certainly a question of which analysis path to choose. Conceptually, pure DTE points to all kinds of DE and while casting a wide net of potentially interesting genes might seem appealing, there are some considerations to be made. For example, if a given transcript is DE, often the question becomes: what happens to the expression of the other transcripts for this gene? Are all transcripts changing in the same direction? If so, it may be better in terms of sensitivity to detect an aggregated output (i.e., DGE). Or, transcript-level expression can be represented as a gene-wise multivariate outcome and isoform switches considered collectively, i.e., by assessing DTU, which is not affected, in either direction, by DGE. DTU implies DTE while the opposite is not necessary true. Generally speaking, we favor the strategy of two clear but orthogonal analyses (DGE and DTU), over a catch-all DTE analysis [74], but this will ultimately be application dependent and scientists should clearly define their question of interest in advance.

Differential transcript expression

Modeling transcript-level count data for DE presents some additional challenges due to increased variability and resolution compared to gene-level analyses. For example, transcript-level abundance estimates are considerably more variable than gene-level counts due to ambiguous assignment of fragments to isoforms [74]. Thus, transcript quantifications inferred by popular tools such as RSEM, Salmon or kallisto, carry a higher degree of uncertainty, which should be accounted for in the downstream DE analysis.

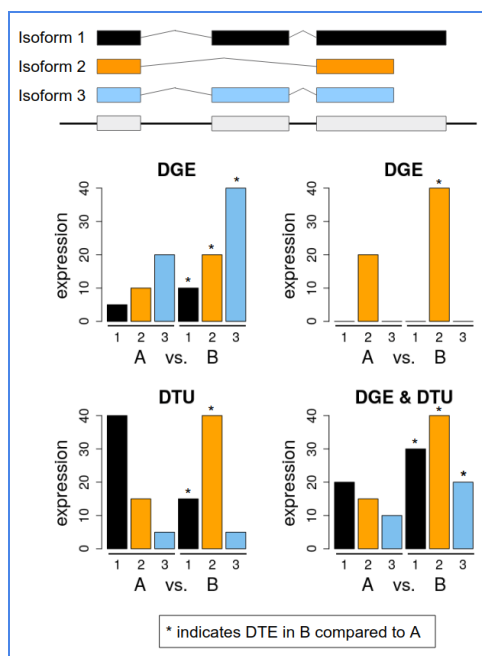


Figure 9: Schematic illustration of some examples of DGE, DTE and DTU for a gene with three isoforms (Isoform 1, 2 and 3) in a two group comparison (A vs B). Bars marked with an asterisk indicate DTE in group B relative to group A.

Transcript quantifications still have many of the properties of count data (e.g., mean-variance relationships) and thus could be used as inputs to the frameworks mentioned above. However, quantifications are estimates that may obscure inference when plugging them into count-based RNA-seq tools. Cufflinks was one of the first methods to use estimated abundances and their corresponding standard errors to perform DTE (and also DGE) analyses; the method quantifies transcript abundances via a likelihood model and EM algorithm and tests of DE are performed by applying the delta method on the abundance parameters [75,86]. Bayesian approaches for identifying DTE based on estimated counts, e.g., ranking via Bayes factors, include EBSeq [134], which uses an empirical Bayesian hierarchical model, and MMSeq [79,168], which fits a linear mixed model to data via Markov Chain Monte Carlo techniques. Similarly, BitSeq (and later cjBitSeq) introduced a generative model that couples both quantification and DE using fully Bayesian inference [95,169]. Most recently, with the advent of ultrafast transcript quantification algorithms, sleuth uses bootstrap samples of each sample of reads to determine the so-called

“inferential variance” and integrates this into the DE calculation through a variance components model on the log-transformed scale [170].

Differential transcript usage (differential splicing)

735 One of the first statistical models for DTU, cuffdiff, calculates the square root of the Jensen-Shannon divergence on estimated *transcript proportions*, and uses the delta method to estimate the variance of this metric under the null hypothesis of no change in proportions [86]. Another conceptually distinct approach formulates a Poisson mixed-effects model on exon- and junction-level quantifications and searches for exon-condition interactions that represent differential usage [171]. Such *departure-from-parallelism* modeling was introduced in earlier
740 analyses of probe-level microarray data for DTU [172]; on RNA-seq data, this approach was further formalized with DEXSeq [173], where a NB model on exon-level counts is formulated. Exon by exon, DEXSeq tests whether an improvement in fit is achieved by adding a single exon-condition interaction, which represents the differential usage of that exon across conditions. A comparison study showed that DEXSeq has a good performance in well-annotated
745 transcriptomes and that filtering of lowly expressed transcripts improves error control [174]; in addition, DEXSeq also works well with transcript quantifications as input [175].

In a similar vein, DRIMSeq [176] and LeafCutter [177] employ the Dirichlet-Multinomial (DM) distribution to perform the same inference task, but treat the output of a gene’s expression as a multivariate response; Bayesian inference for the DM model has also been considered in
750 BayesDRIMSeq [178]. Several tools neglect the uncertainty in estimated transcript-level counts and this is perhaps the reason for inflated FDRs [175]. To address this, RATs uses bootstrapped (transcript-level) quantifications to infer DTU via a G-test of independence, based on the multinomial distribution, on the two groups’ isoform counts [179]. Bayesian methods, such as cjBitSeq [169], instead of considering estimated counts and their uncertainties, focus on the
755 group of transcripts that each read is compatible with (i.e., equivalence classes). In this way, quantification is not required because the DS tools treat the transcript allocation of reads as an unknown latent variable.

Event-level analyses based on percent-spliced-in

760 Some methods perform differential analyses based on percent-spliced-in (PSI) values. PSIs can be computed either for specific events (retained intron, cassette exon, etc.) or at the transcript level, and indicate the fraction of RNA-seq reads supporting the event, obtained as the ratio between the number of reads including the event and the total number of reads including and excluding the event. The difference of the PSIs between conditions is then used to assess DS, performed separately for each event (or transcript). Some of the main DS tools based on PSIs
765 include: rMATS [180], which uses a LRT, and SUPPA2 [181], whose test is based on comparing the observed difference in PSIs across conditions to the empirical cumulative density function of the within-replicates differences of PSIs of splice junctions from similarly expressed transcripts.

770 Event-level analysis, similar to DEXSeq's exon-level approach, separately focuses on each splicing event and results could be aggregated to the gene level by considering the most significant event- or transcript-level test, appropriately adjusted for multiple testing [173,181,182].

Multi stage testing

775 As mentioned, DS analyses can be approached at the gene-, transcript- or event-specific level. While gene-level tests often have higher sensitivity, testing each individual transcript provides increased resolution. However, neither gene- nor transcript-level tests guarantee FDR control on the full set. Stage-wise testing procedures [165,166], instead, first screen for significant genes, and only consider significant transcripts from those genes. This procedure gives gene-level FDR control and allows researchers to leverage the power from gene-level tests while allowing them to interpret results at the transcript level [175]. The same procedure can be applied replacing transcript-level tests with exon- or event-specific tests.

780

Single Cell Transcriptome Sequencing

785 One of the emerging data types in transcriptomics is single cell RNA-seq (scRNA-seq), whereby the expressed content of individual cells is prepared and sequenced. In this case, experimental design is again of critical importance to avoid confounding [183]. Experimentally, capture and reverse transcription efficiency become important, given that the number of mRNA transcripts per mammalian cell is estimated to vary between 50k and 300k [184].

790 Two main experimental approaches are used: plate-based, where cells are sorted into individual wells for lysis and library preparation; or, droplet-based, where each cell is absorbed (together with reagents) and processed within an oil droplet [185]. Several variations of these protocols are now available, increasing the number of cells assayed, but ultimately only a small fraction of the expressed RNAs (cDNAs), often the most highly expressed transcripts, are captured. The features that distinguish scRNA-seq from bulk RNA-seq data include: i) generally low depth of sampling for each cell (due to cost, but also due to lower diversity of cDNA fragments); ii) so-called "dropout" where a cell expresses a transcript but it is unobserved; and, iii) higher levels of biological (since no averaging) and technical (e.g., more amplification) variation.

795

800 Nonetheless, researchers are able to distinguish cell "identities", where identity represents the combined effects of cell "type" (permanent feature) and cell "state" (transient feature) [186]. The Human Cell Atlas, amongst other projects, opens the door for exploring spatial context [187], developmental patterns [188], immune responses [189], response to therapy [190], and an increasing range of basic science and clinical investigations [191–193].

Although many computational aspects of scRNA-seq data are beyond the scope of this review (e.g., dimensionality reduction techniques, ordering cells into lineages), one connected application area that has already received considerable attention is DE analysis. In the simplest

805 setting, cells are first partitioned into different classes (e.g., assumed to correspond to different cell types) via clustering, with the subsequent aim of finding markers for each cluster, e.g., to annotate cell types. To perform this task, a statistical model uses *cells* as experimental units, as opposed to samples in bulk analyses; thus it is worth considering the population to which the conclusions extrapolate to.

810 To date, several methods have been developed to decipher DE between cell types, many of which have been comparatively assessed in recent benchmarks [194,195]. Many of these single-cell-specific methods are extensions or variations of existing bulk approaches. For example, SCDE formulated the RPM (read-per-million) data for a given gene across cells as a mixture of Poisson and negative binomial components; using a Bayesian approach, probabilities of observing a given fold change are converted into empirical p-values [196]. MAST uses a hurdle model on $\log(\text{TPM}+1)$ data, where a logistic regression is used to model whether a gene is expressed or not and conditional on expression, a Gaussian linear model is used. Inferences for the two sets of regression parameters are done in a Bayesian framework that also provides regularization [197]. Again, extending existing approaches, Van den Berge *et al.* [198] proposed a zero-inflated NB (ZINB) model; model fitting is done within the ZINB-WaVE framework [199], 815 estimating cell- and gene-specific posterior probabilities for counts to belong to the NB count component of the ZINB mixture model. These probabilities are used as observation weights in the downstream estimation of regression parameters in the classical NB framework.

820 Nonetheless, many DE methods focus on assessing changes in the mean parameters. But since cell subsets are being compared, we may not expect to have simple shifts in the mean. 825 Instead, it may be informative to detect and understand changes in expression *variability* across conditions [200]. Alternatively, full distributions (instead of means or variances) can be compared, as was proposed in a Bayesian framework in scDD [201], highlighting not only DE but also differential proportions (change in the relative usage of low and high expression), differential modality (change in the number and place of the mode of expression) or some combination thereof. 830

In many applications of single-cell DE analysis, the “sample sizes” (numbers of cells) are generally larger than those commonly used within the optimized frameworks built for bulk RNA-seq data, and thus it seems that the distributional assumptions play less of a role for effective inference. Indeed, a recent comparison highlighted decent performance of t-tests and 835 Wilcoxon rank-sum non-parametric tests in comparing single cell subsets [194].

Beyond comparing cell types, which may or may not involve multiple experimental units (e.g., patients), it will be of increasing interest to compare expression levels of genes *across* biological replicates and conditions. For example, it may be of interest to understand cell-type-specific immune responses following a stimulus. A recent study compared multiple patients across 840 stimulated and unstimulated conditions by first computationally separating immune cell types [189]; to do this, cells from a given cell type were *aggregated* into a “pseudo-bulk” RNA-seq dataset and DE was performed using standard tools.

Long Read Transcriptome Sequencing

845 The short read length of Illumina-based RNA-seq complicates unambiguous placement of reads to the genome, especially in repeat regions [202], and adds difficulties to the assembly, identification and quantification of expressed isoforms [203–205]. In contrast, so-called *third-generation*, or long-read, sequencing technologies, led by Pacific Biosciences (PacBio) [206] and Oxford Nanopore Technologies (ONT), are able to generate reads that are much longer. By sequencing single molecules, they can also avoid the need for PCR amplification, 850 hence reducing coverage biases [207,208]. Currently, long-read technologies incur a higher average cost and a higher error rate than short read sequencing [203,209]. However, this is a rapidly developing field and improvements in error rates and throughput are to be expected.

The strategies employed by PacBio and ONT to generate long sequencing reads of single molecules differ in many ways. PacBio, with its RSII and Sequel instruments, uses SMRT 855 (single molecule real time) sequencing [210], where the reactions take place inside so-called zero-mode waveguides (ZMWs) [211]. At the bottom of each ZMW, there is a single DNA polymerase molecule. As the polymerase processes a DNA fragment, the incorporation of each nucleotide leads to a fluorescent signal, which is detected by the ZMW and converted to a base call. A specific characteristic of the PacBio library preparation system is the creation of 860 SMRTbell templates [212], which are obtained by ligating SMRTbell hairpin adapters. The result is a circular construct, where the two strands of the template are separated by adapters with known sequences. As the construct is processed by the polymerase in the ZMW, the original template can be passed multiple times. Since the sequencing errors are largely random [213], the base-level error rate can be considerably reduced by forming a consensus over these 865 passes.

ONT, in contrast, uses a different sequencing strategy, based on protein nanopores placed in a polymer membrane [214] for its MinION and PromethION sequencers. A current is passed through the nanopores, and as the template molecule is passed through the pore by a motor protein, each combination of bases induces a change in the current. Analyzing the exact nature 870 of this change allows the identification of the template sequence. By adding a hairpin sequence to the end of the double-stranded cDNA fragment before denaturing it into a single-stranded molecule and passing it through the nanopore, both the template sequence and its complement are included in a single read and can be combined at the base-calling step to generate a higher-quality, so called 2D, read [215]. In contrast to PacBio, ONT also offers direct sequencing 875 of RNA [216]. Advantages of this include that the reverse transcription step is avoided, which may reduce biases, and that RNA modifications can be directly observed, since they also change the current passing through the nanopore in characteristic ways [217]. However, at present, the required amount of starting material is considerably higher than for cDNA protocols.

880 Applications to cDNA (RNA) include both transcriptome-wide sequencing and characterization
of specific genes via targeted sequencing [15,203,218–222], as well as performance
evaluations based on synthetic transcript catalogs (ERCC with 92 sequences, or SIRV, with 68).
Long-read transcript sequencing (LRTS) offers the potential that every read represents a
full-length transcript. If this was indeed true, *de novo* (reference-free) identification of the full set
885 of observed isoforms would be straightforward, and only require grouping together reads
expected to differ only by sequencing errors (which, depending on the error rate and isoform
similarity, may of course still not be trivial). This is, however, not currently the case, both due to
fragmentation and degradation of template molecules during library preparation and because of
early termination of the sequencing, which leads to ambiguities in transcript identification [223].
In particular, this means that it is not easy to determine whether truncated variants are present.

890 Transcript identification from LRTS can be either reference-based or reference-free. The latter
typically involves clustering reads based on similarity, followed by polishing of the consensus
sequence within each cluster [15,224–227]. Since LRTS is still a young field, methods and tools
for reference-based alignment are still emerging, but so far include a mix of established tools,
such as GMAP [58] and new innovations, such as minimap2 [228]. A recent study comparing
895 PacBio, ONT and Illumina data [203], showed that the long-read technologies were indeed
much better at correctly identifying expressed SIRV transcripts than *de novo* assembly of short
reads.

The rapid technological developments in LRTS also mean that the read generation process,
e.g., biases affecting the ability of observing a given read, is still largely unknown. In addition,
900 read lengths are extremely variable, error rates are relatively high, and throughput is still
relatively low. In particular for the PacBio RSII instrument, the selection of transcript molecules
is biased towards short sequences [223]. Thus, samples are typically size-fractionated before
sequencing, which distorts the abundance estimates. Taken together, these and other aspects
make accurate transcript quantification from LRTS more difficult, and new models and tools will
905 be needed. Encouragingly, a recent study showed that by combining LRTS and Illumina data,
more accurate quantifications for the artificial SIRV transcripts could be achieved [203].

Since abundance estimation for long reads returns values in the form of read (or transcript)
counts, it is plausible that the DE machinery developed for short-read data can be applied in a
similar way. The quality of the DE calls will be directly dependent on the accuracy of the
910 abundance estimates. However, the current low depth of sequencing compared to short-read
data sets will ultimately lead to low power to detect DE features.

Summary

915 *"I'm a scientist and I know what constitutes proof. But the reason I call myself by my childhood
name is to remind myself that a scientist must also be absolutely like a child. If [they] see a
thing, [they] must say that [they] see it, whether it was what [they] thought [they] were going to*

see or not. See first, think later, then test. But always see first. Otherwise you will only see what you were expecting. Most scientists forget that.”

— adapted from *The Ultimate Hitchhiker's Guide to the Galaxy* by Douglas Adams

920 In this review, we gave an overview of the data science of gene expression analysis, with a
focus on methods to estimate transcript-level abundance and statistical tools for assessing DE.
Notably, RNA-seq data is often an intermediate discovery step where the detected molecular
changes represent candidates for further follow up. Nonetheless, the analysis of RNA-seq data
for gene expression is already very mature, due to a deep understanding of the biases present,
925 (estimated) count tables and to a refined understanding of how well tools perform via the many
benchmarks available.

Ultimately, the success of RNA-seq lies in its wide range of applications and it is likely that
Illumina-based short-fragment RNA-seq will continue to be the workhorse for the field for many
years. With increasing fidelity of single-cell protocols, many tools are emerging to deal with the
930 additional complexities of single cell measurements and these will be further refined in the
coming years. Similarly, with the decreasing costs and lower error rates of long-read
technologies, it may be possible to characterize alternative transcription quantitatively with
full-length transcript sequencing, thus considerably reducing read-to-transcript ambiguity;
however, much still needs to be learned about the biases present.

935

References

1. Lister R, O'Malley RC, Tonti-Filippini J, Gregory BD, Berry CC, Millar AH, et al. Highly integrated single-base resolution maps of the epigenome in Arabidopsis. *Cell*. 2008;133: 523–536.
2. Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, et al. The transcriptional landscape of the yeast genome defined by RNA sequencing. *Science*. 2008;320: 1344–1349.
3. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods*. 2008;5: 621–628.
4. Cloonan N, Forrest ARR, Kolle G, Gardiner BBA, Faulkner GJ, Brown MK, et al. Stem cell transcriptome profiling via massive-scale mRNA sequencing. *Nat Methods*. 2008;5: 613–619.
5. Wilhelm BT, Marguerat S, Watt S, Schubert F, Wood V, Goodhead I, et al. Dynamic repertoire of a eukaryotic transcriptome surveyed at single-nucleotide resolution. *Nature*. 2008;453: 1239–1243.
6. Palazzo AF, Lee ES. Non-coding RNA: what is functional and what is junk? *Front Genet*. 2015;6: 2.
7. Zhao W, He X, Hoadley K a., Parker JS, Hayes DN, Perou CM. Comparison of RNA-Seq by poly (A) capture, ribosomal RNA depletion, and DNA microarray for expression profiling. *BMC Genomics*. 2014;15: 419.
8. Bentley DR, Balasubramanian S, Swerdlow HP, Smith GP, Milton J, Brown CG, et al.

955

- Accurate whole human genome sequencing using reversible terminator chemistry. *Nature*. 2008;456: 53–59.
9. Ju J, Kim DH, Bi L, Meng Q, Bai X, Li Z, et al. Four-color DNA sequencing by synthesis using cleavable fluorescent nucleotide reversible terminators. *Proc Natl Acad Sci U S A*. 2006;103: 19635–19640.
- 960 10. Collado-Torres L, Nellore A, Kammers K, Ellis SE, Taub MA, Hansen KD, et al. Reproducible RNA-seq analysis using recount2. *Nat Biotechnol*. 2017;35: 319–321.
11. Levin JZ, Yassour M, Adiconis X, Nusbaum C, Thompson DA, Friedman N, et al. Comprehensive comparative analysis of strand-specific RNA sequencing methods. *Nat Methods*. Nature Publishing Group; 2010;7: 709–715.
- 965 12. Zhao S, Zhang Y, Gordon W, Quan J, Xi H, Du S, et al. Comparison of stranded and non-stranded RNA-seq transcriptome profiling and investigation of gene overlap. *BMC Genomics*. 2015;16: 675.
13. Parkhomchuk D, Borodina T, Amstislavskiy V, Banaru M, Hallen L, Krobitsch S, et al. Transcriptome analysis by strand-specific sequencing of complementary DNA. *Nucleic Acids Res*. 2009;37: e123.
- 970 14. Mamanova L, Turner DJ. Low-bias, strand-specific transcriptome Illumina sequencing by on-flowcell reverse transcription (FRT-seq). *Nat Protoc*. 2011;6: 1736–1747.
15. Wang B, Tseng E, Regulski M, Clark TA, Hon T, Jiao Y, et al. Unveiling the complexity of the maize transcriptome by single-molecule long-read sequencing. *Nat Commun*. Nature Publishing Group; 2016;7: 11708.
- 975 16. Lazic SE. *Experimental Design for Laboratory Biologists: Maximising Information and Improving Reproducibility*. 1 edition. Cambridge University Press; 2017.
17. Hart SN, Therneau TM, Zhang Y, Poland G a., Kocher J-P. Calculating sample size estimates for RNA sequencing data. *J Comput Biol*. 2013;20: 970–978.
- 980 18. Guo Y, Zhao S, Li CI, Sheng Q, Shyr Y. RNAseqPS: A web tool for estimating sample size and power for RNAseq experiment. *Cancer Inform*. 2014;13: 1–5.
19. Zhao S, Li C-I, Guo Y, Sheng Q, Shyr Y. RnaSeqSampleSize: real data based sample size estimation for RNA sequencing. *BMC Bioinformatics*. 2018;19: 191.
- 985 20. Busby M a., Stewart C, Miller C a., Grzeda KR, Marth GT. Scotty: a web tool for designing RNA-Seq experiments to measure differential gene expression. *Bioinformatics*. 2013;29: 656–657.
21. Poplawski A, Binder H. Feasibility of sample size calculation for RNA-seq studies. *Brief Bioinform*. 2018;19: 713–720.
- 990 22. Oshlack A, Wakefield MJ. Transcript length bias in RNA-seq data confounds systems biology. *Biol Direct*. 2009;4: 14.
23. Liu Y, Zhou J, White KP. RNA-seq differential expression studies: more sequence or more replication? *Bioinformatics*. 2014;30: 301–304.
24. Schurch NJ, Schofield P, Gierliński M, Cole C, Sherstnev A, Singh V, et al. How many biological replicates are needed in an RNA-seq experiment and which differential expression tool should you use? *RNA*. 2016;22: 839–851.
- 995 25. Mercer TR, Clark MB, Crawford J, Brunck ME, Gerhardt DJ, Taft RJ, et al. Targeted sequencing for gene discovery and quantification using RNA CaptureSeq. *Nat Protoc*. 2014;9: 989–1009.
- 1000 26. Cabanski CR, Magrini V, Griffith M, Griffith OL, McGrath S, Zhang J, et al. cDNA hybrid capture improves transcriptome analysis on low-input and archived samples. *J Mol Diagn*. 2014;16: 440–451.
27. Irimia M, Weatheritt RJ, Ellis JD, Parikshak NN, Gonatopoulos-Pournatzis T, Babor M, et al.

- 1005 A highly conserved program of neuronal microexons is misregulated in autistic brains. *Cell*. 2014;159: 1511–1523.
28. Eom T, Zhang C, Wang H, Lay K, Fak J, Noebels JL, et al. NOVA-dependent regulation of cryptic NMD exons controls synaptic protein levels after seizure. *Elife*. 2013;2: e00178.
29. Fratta P, Sivakumar P, Humphrey J, Lo K, Ricketts T, Oliveira H, et al. Mice with endogenous TDP-43 mutations exhibit gain of splicing function and characteristics of amyotrophic lateral sclerosis. *EMBO J*. 2018;37. doi:10.15252/embj.201798684
- 1010 30. Salzman J, Gawad C, Wang PL, Lacayo N, Brown PO. Circular RNAs are the predominant transcript isoform from hundreds of human genes in diverse cell types. *PLoS One*. 2012;7: e30733.
31. Kim T-K, Hemberg M, Gray JM. Enhancer RNAs: a class of long noncoding RNAs synthesized at enhancers. *Cold Spring Harb Perspect Biol*. 2015;7: a018622.
- 1015 32. Parker BC, Zhang W. Fusion genes in solid tumors: an emerging target for cancer diagnosis and treatment. *Chin J Cancer*. 2013;32: 594–603.
33. Frye M, Jaffrey SR, Pan T, Rechavi G, Suzuki T. RNA modifications: what have we learned and where are we headed? *Nat Rev Genet*. 2016;17: 365–372.
- 1020 34. Perou CM, Sørlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, et al. Molecular portraits of human breast tumours. *Nature*. 2000;406: 747–752.
35. Climente-González H, Porta-Pardo E, Godzik A, Eyraes E. The Functional Impact of Alternative Splicing in Cancer. *Cell Rep*. 2017;20: 2215–2226.
36. Cieślak M, Chinnaiyan AM. Cancer transcriptome profiling at the juncture of clinical translation. *Nat Rev Genet*. 2017; doi:10.1038/nrg.2017.96
- 1025 37. Pedersen G, Kanigan T. Clinical RNA sequencing in oncology: where are we? *Per Med*. 2016;13: 209–213.
38. Ståhl PL, Salmén F, Vickovic S, Lundmark A, Navarro JF, Magnusson J, et al. Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science*. 2016;353: 78–82.
- 1030 39. Westermann AJ, Gorski SA, Vogel J. Dual RNA-seq of pathogen and host. *Nat Rev Microbiol*. Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved.; 2012;10: 618.
40. Piskol R, Ramaswami G, Li JB. Reliable identification of genomic variants from RNA-seq data. *Am J Hum Genet*. 2013;93: 641–651.
- 1035 41. Park E, Williams B, Wold BJ, Mortazavi A. RNA editing in the human ENCODE RNA-seq data. *Genome Res*. 2012;22: 1626–1633.
42. Uszczyńska-Ratajczak B, Lagarde J, Frankish A, Guigó R, Johnson R. Towards a complete map of the human long non-coding RNA transcriptome. *Nat Rev Genet*. 2018;19: 535–548.
- 1040 43. Bashiardes S, Zilberman-Schapira G, Elinav E. Use of Metatranscriptomics in Microbiome Research. *Bioinform Biol Insights*. 2016;10: 19–25.
44. Racle J, de Jonge K, Baumgaertner P, Speiser DE, Gfeller D. Simultaneous enumeration of cancer and immune cell types from bulk tumor gene expression data. *Elife*. 2017;6. doi:10.7554/eLife.26476
- 1045 45. Martin JA, Wang Z. Next-generation transcriptome assembly. *Nat Rev Genet*. 2011;12: 671–682.
46. Castel SE, Levy-Moonshine A, Mohammadi P, Banks E, Lappalainen T. Tools and best practices for data processing in allelic expression analysis. *Genome Biol*. 2015;16: 195.
47. Sun W, Hu Y. eQTL Mapping Using RNA-seq Data. *Stat Biosci*. 2013;5: 198–219.
- 1050 48. Alamancos GP, Agirre E, Eyraes E. Methods to study splicing from high-throughput RNA sequencing data. *Methods Mol Biol*. 2014;1126: 357–397.

49. van Dam S, Vösa U, van der Graaf A, Franke L, de Magalhães JP. Gene co-expression analysis for functional classification and gene-disease predictions. *Brief Bioinform.* 2018;19: 575–592.
- 1055 50. Khatri P, Sirota M, Butte AJ. Ten years of pathway analysis: current approaches and outstanding challenges. *PLoS Comput Biol.* 2012;8: e1002375.
51. de Leeuw CA, Neale BM, Heskes T, Posthuma D. The statistical properties of gene-set analysis. *Nat Rev Genet.* 2016;17: 353–364.
- 1060 52. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics.* Oxford University Press; 2009;25: 1105–1111.
53. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* 2009;10: R25.
54. Engström PG, Steijger T, Sipos B, Grant GR, Kahles A, Rättsch G, et al. Systematic evaluation of spliced alignment programs for RNA-seq data. *Nat Methods.* 2013;10: 1185–1191.
- 1065 55. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29: 15–21.
56. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory requirements. *Nat Methods.* 2015;12: 357–360.
- 1070 57. Liao Y, Smyth GK, Shi W. The Subread aligner: fast, accurate and scalable read mapping by seed-and-vote. *Nucleic Acids Res.* 2013;41: e108.
58. Wu TD, Watanabe CK. GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics.* 2005;21: 1859–1875.
- 1075 59. Lin H-N, Hsu W-L. DART - a fast and accurate RNA-seq mapper with a partitioning strategy. *Bioinformatics.* 2017; doi:10.1093/bioinformatics/btx558
60. Sedlazeck FJ, Rescheneder P, von Haeseler A. NextGenMap: fast and accurate read mapping in highly polymorphic genomes. *Bioinformatics.* 2013;29: 2790–2791.
61. Medina I, Tárraga J, Martínez H, Barrachina S, Castillo MI, Paschall J, et al. Highly sensitive and ultrafast read mapping for RNA-seq analysis. *DNA Res.* 2016;23: 93–100.
- 1080 62. Baruzzo G, Hayer KE, Kim EJ, Di Camillo B, FitzGerald GA, Grant GR. Simulation-based comprehensive benchmarking of RNA-seq aligners. *Nat Methods.* 2017;14: 135–139.
63. Bonfert T, Kirner E, Csaba G, Zimmer R, Friedel CC. ContextMap 2: fast and accurate context-based RNA-seq mapping. *BMC Bioinformatics.* 2015;16: 122.
- 1085 64. Wu TD, Nacu S. Fast and SNP-tolerant detection of complex variants and splicing in short reads. *Bioinformatics.* 2010;26: 873–881.
65. Aken BL, Ayling S, Barrell D, Clarke L, Curwen V, Fairley S, et al. The Ensembl gene annotation system. *Database.* 2016;2016. doi:10.1093/database/baw093
66. Liao Y, Smyth GK, Shi W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics.* 2014;30: 923–930.
- 1090 67. Anders S, Pyl PT, Huber W. HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics.* 2015;31: 166–169.
68. Roberts A, Trapnell C, Donaghey J, Rinn JL, Pachter L. Improving RNA-Seq expression estimates by correcting for fragment bias. *Genome Biol.* 2011;12: R22.
69. Hansen KD, Brenner SE, Dudoit S. Biases in Illumina transcriptome sequencing caused by random hexamer priming. *Nucleic Acids Res.* 2010;38: e131.
- 1095 70. Liu X, Shi X, Chen C, Zhang L. Improving RNA-Seq expression estimation by modeling isoform- and exon-specific read sequencing rate. *BMC Bioinformatics.* 2015;16: 332.
71. Love MI, Hogenesch JB, Irizarry RA. Modeling of RNA-seq fragment sequence bias reduces systematic errors in transcript abundance estimation. *Nat Biotechnol.* 2016;34:

- 1100 1287–1291.
72. Robert C, Watson M. Errors in RNA-Seq quantification affect genes of relevance to human disease. *Genome Biol.* 2015;16: 177.
73. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics.* 2011;12: 323.
- 1105 74. Sonesson C, Love MI, Robinson MD. Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Res.* 2015;4: 1521.
75. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc.* 2012;7: 562–578.
- 1110 76. Xing Y, Yu T, Wu YN, Roy M, Kim J, Lee C. An expectation-maximization algorithm for probabilistic reconstructions of full-length isoforms from splice graphs. *Nucleic Acids Res.* 2006;34: 3150–3160.
77. Jiang H, Wong WH. Statistical inferences for isoform expression in RNA-Seq. *Bioinformatics.* 2009;25: 1026–1032.
- 1115 78. Li B, Ruotti V, Stewart RM, Thomson JA, Dewey CN. RNA-Seq gene expression estimation with read mapping uncertainty. *Bioinformatics.* 2010;26: 493–500.
79. Turro E, Su S-Y, Gonçalves Â, Coin LJM, Richardson S, Lewin A. Haplotype and isoform specific expression estimation using multi-mapping RNA-seq reads. *Genome Biol.* 2011;12: R13.
- 1120 80. Richard H, Schulz MH, Sultan M, Nürnberger A, Schrinner S, Balzereit D, et al. Prediction of alternative isoforms from exon expression levels in RNA-Seq experiments. *Nucleic Acids Res.* 2010;38: e112.
81. Nicolae M, Mangul S, Măndoiu II, Zelikovsky A. Estimation of alternative splicing isoform frequencies from RNA-Seq data. *Algorithms Mol Biol.* 2011;6: 9.
- 1125 82. Mezlini AM, Smith EJM, Fiume M, Buske O, Savich GL, Shah S, et al. iReckon: simultaneous isoform discovery and abundance estimation from RNA-seq data. *Genome Res.* 2013;23: 519–529.
83. Zakeri M, Srivastava A, Almodaresi F, Patro R. Improved data-driven likelihood factorizations for transcript abundance estimation. *Bioinformatics.* 2017;33: i142–i151.
- 1130 84. Roberts A, Pachter L. Streaming fragment assignment for real-time analysis of sequencing experiments. *Nat Methods.* 2013;10: 71–73.
85. Cappé O, Moulines E. On-line expectation--maximization algorithm for latent data models. *J R Stat Soc Series B Stat Methodol.* Wiley Online Library; 2009;71: 593–613.
- 1135 86. Trapnell C, Williams B a., Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol.* 2010;28: 511–515.
87. Li W, Feng J, Jiang T. IsoLasso: a LASSO regression approach to RNA-Seq based transcriptome assembly. *J Comput Biol.* 2011;18: 1693–1707.
- 1140 88. Canzar S, Andreotti S, Weese D, Reinert K, Klau GW. CIDANE: comprehensive isoform discovery and abundance estimation. *Genome Biol.* 2016;17: 16.
89. Maretty L, Sibbesen JA, Krogh A. Bayesian transcriptome assembly. *Genome Biol.* 2014;15: 501.
90. Shi X, Wang X, Wang T-L, Hilakivi-Clarke L, Clarke R, Xuan J. Sparselso: a novel Bayesian approach to identify alternatively spliced isoforms from RNA-seq data. *Bioinformatics.* 2018;34: 56–63.
- 1145 91. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol.*

- 2015;33: 290–295.
- 1150 92. Tomescu AI, Kuosmanen A, Rizzi R, Mäkinen V. A novel min-cost flow method for estimating transcript expression with RNA-Seq. *BMC Bioinformatics*. 2013;14 Suppl 5: S15.
93. Bernard E, Jacob L, Mairal J, Vert J-P. Efficient RNA isoform identification and quantification from RNA-Seq data with network flows. *Bioinformatics*. 2014;30: 2447–2455.
94. Shao M, Kingsford C. Accurate assembly of transcripts through phase-preserving graph decomposition. *Nat Biotechnol*. 2017;35: 1167–1169.
- 1155 95. Glaus P, Honkela A, Rattray M. Identifying differentially expressed transcripts from RNA-seq data with biological variation. *Bioinformatics*. 2012;28: 1721–1728.
96. SEQC/MAQC-III Consortium. A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium. *Nat Biotechnol*. 2014;32: 903–914.
- 1160 97. Hensman J, Papastamoulis P, Glaus P, Honkela A, Rattray M. Fast and accurate approximate inference of transcript expression from RNA-seq data. *Bioinformatics*. 2015;31: 3881–3889.
98. Nariai N, Hirose O, Kojima K, Nagasaki M. TIGAR: transcript isoform abundance estimation method with gapped alignment of RNA-Seq data by variational Bayesian inference.
- 1165 99. Amari S-I, Nagaoka H. *Methods of information geometry*, volume 191 of *Translations of Mathematical Monographs*. American Mathematical Society. 2000;13.
100. Patro R, Mount SM, Kingsford C. Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nat Biotechnol*. 2014;32: 462–464.
- 1170 101. Varadhan R, Roland C. Squared Extrapolation Methods (SQUAREM): A New Class of Simple and Efficient Numerical Schemes for Accelerating the Convergence of the EM Algorithm. *bepress*; 2004; Available: <https://biostats.bepress.com/jhubiostat/paper63/>
102. Zhang Z, Wang W. RNA-Skim: a rapid method for RNA-Seq quantification at transcript level. *Bioinformatics*. 2014;30: i283–i292.
- 1175 103. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol*. 2016;34: 525–527.
104. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods*. 2017;14: 417–419.
- 1180 105. Foulds J, Boyles L, DuBois C, Smyth P, Welling M. Stochastic Collapsed Variational Bayesian Inference for Latent Dirichlet Allocation. *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM; 2013. pp. 446–454.
106. Ju CJ-T, Li R, Wu Z, Jiang J-Y, Yang Z, Wang W. Fleximer: Accurate Quantification of RNA-Seq via Variable-Length k-mers. *Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*. ACM; 2017. pp. 263–272.
- 1185 107. Kanitz A, Gypas F, Gruber AJ, Gruber AR, Martin G, Zavolan M. Comparative assessment of methods for the computational inference of transcript isoform abundance from RNA-seq data. *Genome Biol*. 2015;16: 150.
- 1190 108. Germain P-L, Vitriolo A, Adamo A, Laise P, Das V, Testa G. RNAontheBENCH: computational and empirical resources for benchmarking RNAseq quantification and differential expression methods. *Nucleic Acids Res*. 2016;44: 5054–5067.
109. Teng M, Love MI, Davis CA, Djebali S, Dobin A, Graveley BR, et al. A benchmark for RNA-seq quantification pipelines. *Genome Biol*. 2016;17: 74.
- 1195 110. Zhang C, Zhang B, Lin L-L, Zhao S. Evaluation and comparison of computational tools

- for RNA-seq isoform quantification. *BMC Genomics*. 2017;18: 583.
111. Prakash C, Haeseler AV. An Enumerative Combinatorics Model for Fragmentation Patterns in RNA Sequencing Provides Insights into Nonuniformity of the Expected Fragment Starting-Point and Coverage Profile. *J Comput Biol*. 2017;24: 200–212.
- 1200 112. Jones DC, Kuppusamy KT, Palpant NJ, Peng X, Murry CE, Ruohola-Baker H, et al. Isolator: accurate and stable analysis of isoform-level expression in RNA-Seq experiments [Internet]. *bioRxiv*. 2016. p. 088765. doi:10.1101/088765
113. Sonesson C, Love MI, Patro R, Hussain S, Malhotra D, Robinson MD. A junction coverage compatibility score to quantify the reliability of transcript abundance estimates and annotation catalogs [Internet]. *bioRxiv*. 2018. p. 378539. doi:10.1101/378539
- 1205 114. Efron B, Hastie T. *Computer Age Statistical Inference: Algorithms, Evidence, and Data Science*. 1st ed. New York, NY, USA: Cambridge University Press; 2016.
115. Tusher VG, Tibshirani R, Chu G. Significance analysis of microarrays applied to the ionizing radiation response. *Proc Natl Acad Sci U S A*. 2001;98: 5116–5121.
- 1210 116. Smyth GK. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Stat Appl Genet Mol Biol*. 2004;3: Article3.
117. Bourgon R, Gentleman R, Huber W. Independent filtering increases detection power for high-throughput experiments. *Proc Natl Acad Sci U S A*. 2010;107: 9546–9551.
118. Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y. RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays. *Genome Res*. 2008;18: 1509–1517.
- 1215 119. Robinson MD, Oshlack A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol*. 2010;11: R25.
120. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol*. 2010;11: R106.
- 1220 121. Dillies M-A, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, et al. A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Brief Bioinform*. 2013;14: 671–683.
122. Hansen KD, Irizarry R a., Wu Z. Removing technical variability in RNA-seq data using conditional quantile normalization. *Biostatistics*. 2012;13: 204–216.
- 1225 123. Risso D, Schwartz K, Sherlock G, Dudoit S. GC-content normalization for RNA-Seq data. *BMC Bioinformatics*. 2011;12: 480.
124. Lovén J, Orlando DA, Sigova AA, Lin CY, Rahl PB, Burge CB, et al. Revisiting Global Gene Expression Analysis. *Cell*. 2012;151: 476–482.
- 1230 125. Risso D, Ngai J, Speed TP, Dudoit S. Normalization of RNA-seq data using factor analysis of control genes or samples. *Nat Biotechnol*. Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved.; 2014;32: 896.
126. Taruttis F, Feist M, Schwarzfischer P, Gronwald W, Kube D, Spang R, et al. External calibration with *Drosophila* whole-cell spike-ins delivers absolute mRNA fold changes from human RNA-Seq and qPCR data. *Biotechniques*. 2017;62: 53–61.
- 1235 127. Hicks SC, Okrah K, Paulson JN, Quackenbush J, Irizarry RA, Bravo HC. Smooth quantile normalization. *Biostatistics*. 2018;19: 185–198.
128. Vallejos CA, Risso D, Scialdone A, Dudoit S, Marioni JC. Normalizing single-cell RNA sequencing data: challenges and opportunities. *Nat Methods*. 2017; doi:10.1038/nmeth.4292
- 1240 129. Himes BE, Jiang X, Wagner P, Hu R, Wang Q, Klanderman B, et al. RNA-Seq transcriptome profiling identifies CRISPLD2 as a glucocorticoid responsive gene that modulates cytokine function in airway smooth muscle cells. *PLoS One*. 2014;9: e99625.

- 1245 130. Love MI, Anders S, Kim V, Huber W. RNA-Seq workflow: gene-level exploratory analysis and differential expression. *F1000Res*. 2016;4. doi:10.12688/f1000research.7035.2
131. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*. 2014;15: 550.
132. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010;26: 139–140.
- 1250 133. Hardcastle TJ, Kelly K a. baySeq: empirical Bayesian methods for identifying differential expression in sequence count data. *BMC Bioinformatics*. 2010;11: 422.
134. Leng N, Dawson J a., Thomson J a., Ruotti V, Rissman AI, Smits BMG, et al. EBSeq: an empirical Bayes hierarchical model for inference in RNA-seq experiments. *Bioinformatics*. 2013;29: 1035–1043.
- 1255 135. Robinson MD, Smyth GK. Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics*. 2007;23: 2881–2887.
136. McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res*. 2012;40: 4288–4297.
- 1260 137. Robinson MD, Smyth GK. Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics*. 2008;9: 321–332.
138. Soneson C, Delorenzi M. A comparison of methods for differential expression analysis of RNA-seq data. *BMC Bioinformatics*. 2013;14: 91.
139. Zhou X, Lindsay H, Robinson MD. Robustly detecting differential expression in RNA sequencing data using observation weights. *Nucleic Acids Res*. 2014;42: e91.
- 1265 140. Cox DR, Reid N. Parameter orthogonality and approximate conditional inference. *J R Stat Soc Series B Stat Methodol*. 1987;49: 1–39.
141. Chen Y, Lun ATL, Smyth GK. Differential Expression Analysis of Complex RNA-seq Experiments Using edgeR. In: Datta S, Nettleton D, editors. *Statistical Analysis of Next Generation Sequencing Data*. Springer International Publishing; 2014. pp. 51–74.
- 1270 142. Wu H, Wang C, Wu Z. A new shrinkage estimator for dispersion improves differential expression detection in RNA-seq data. *Biostatistics*. 2013;14: 232–243.
143. Nelder JA, Wedderburn RWM. *Generalized Linear Models*. *J R Stat Soc Ser A*. [Royal Statistical Society, Wiley]; 1972;135: 370–384.
- 1275 144. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Series B Stat Methodol*. 1995;57: 289–300.
145. Benjamini Y, Yekutieli D. The control of the false discovery rate in multiple testing under dependency. *Ann Stat*. 2001;29: 1165–1188.
- 1280 146. Lund SP, Nettleton D, McCarthy DJ, Smyth GK. Detecting differential expression in RNA-sequence data using quasi-likelihood with shrunken dispersion estimates. *Stat Appl Genet Mol Biol*. 2012;11. doi:10.1515/1544-6115.1826
147. Di Y, Schafer DW, Cumbie JS, Chang JH. The NBP Negative Binomial Model for Assessing Differential Gene Expression from RNA-Seq. *Stat Appl Genet Mol Biol*. 2011;10. doi:10.2202/1544-6115.1637
- 1285 148. Van De Wiel M a., Leday GGR, Pardo L, Rue H, Van Der Vaart AW, Van Wieringen WN. Bayesian analysis of RNA sequencing data by estimating multiple shrinkage priors. *Biostatistics*. 2013;14: 113–128.
149. van de Wiel MA, Neerincx M, Buffart TE, Sie D, Verheul HMW. ShrinkBayes: a versatile R-package for analysis of count-based sequencing data in complex study designs. *BMC Bioinformatics*. 2014;15: 116.
- 1290 150. Rue H, Martino S, Chopin N. Approximate Bayesian inference for latent Gaussian

- models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2009;71: 319–392.
- 1295 151. 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.
152. Li J, Tibshirani R. Finding consistent patterns: a nonparametric approach for identifying
differential expression in RNA-Seq data. *Stat Methods Med Res*. 2013;22: 519–536.
153. Stephens M. False discovery rates: a new deal. *Biostatistics*. 2016;
doi:10.1093/biostatistics/kxw041
- 1300 154. Zhu A, Ibrahim JG, Love MI. Heavy-tailed prior distributions for sequence count data:
removing the noise and preserving large differences [Internet]. *bioRxiv*. 2018. p. 303255.
doi:10.1101/303255
155. Leek JT, Storey JD. Capturing Heterogeneity in Gene Expression Studies by Surrogate
Variable Analysis. Gibson G, editor. *PLoS Genet*. Public Library of Science; 2007;3: 12.
- 1305 156. Leek JT. svaseq: removing batch effects and other unwanted noise from sequencing
data. *Nucleic Acids Res*. 2014;42. doi:10.1093/nar/gku864
157. Finner H, Roters M. On the False Discovery Rate and Expected Type I Errors.
Biometrical Journal. 2001;43: 985–1005.
158. McCarthy DJ, Smyth GK. Testing significance relative to a fold-change threshold is a
TREAT. *Bioinformatics*. 2009;25: 765–771.
- 1310 159. Chen Y, Lun ATL, Smyth GK. From reads to genes to pathways: differential expression
analysis of RNA-Seq experiments using Rsubread and the edgeR quasi-likelihood pipeline.
F1000Res. 2016;5: 1438.
160. Di Y, Emerson SC, Schafer DW, Kimbrel JA, Chang JH. Higher order asymptotics for
negative binomial regression inferences from RNA-sequencing data. *Stat Appl Genet Mol
Biol*. 2013;12: 49–70.
- 1315 161. Storey JD. The positive false discovery rate: A Bayesian interpretation and the q-value.
Ann Stat. 2003;31: 2013–2035.
162. Ignatiadis N, Klaus B, Zaugg JB, Huber W. Data-driven hypothesis weighting increases
detection power in genome-scale multiple testing. *Nat Methods*. 2016;13: 577–580.
- 1320 163. Efron B. Large-Scale Simultaneous Hypothesis Testing. *J Am Stat Assoc*. Taylor &
Francis; 2004;99: 96–104.
164. Efron B. Size, Power and False Discovery Rates. *Ann Stat*. Institute of Mathematical
Statistics; 2007;35: 1351–1377.
- 1325 165. Van den Berge K, Sonesson C, Robinson MD, Clement L. stageR: a general stage-wise
method for controlling the gene-level false discovery rate in differential expression and
differential transcript usage. *Genome Biol*. 2017;18: 151.
166. Heller R, Manduchi E, Grant GR, Ewens WJ. A flexible two-stage procedure for
identifying gene sets that are differentially expressed. *Bioinformatics*. 2009;25: 1019–1025.
- 1330 167. Kakaradov B, Xiong HY, Lee LJ, Jovic N, Frey BJ. Challenges in estimating percent
inclusion of alternatively spliced junctions from RNA-seq data. *BMC Bioinformatics*. 2012;13
Suppl 6: S11.
168. Turro E, Astle WJ, Tavar?? S. Flexible analysis of RNA-seq data using mixed effects
models. *Bioinformatics*. 2014;30: 180–188.
- 1335 169. Papastamoulis P, Rattray M. A Bayesian model selection approach for identifying
differentially expressed transcripts from RNA sequencing data. *J R Stat Soc C*. 2018;67:
3–23.
170. Pimentel H, Bray NL, Puente S, Melsted P, Pachter L. Differential analysis of RNA-seq
incorporating quantification uncertainty. *Nat Methods*. Nature Publishing Group, a division

- 1340 of Macmillan Publishers Limited. All Rights Reserved.; 2017;14: 687.
171. Blekhman R, Marioni JC, Zumbo P, Stephens M, Gilad Y. Sex-specific and lineage-specific alternative splicing in primates. *Genome Res.* 2010;20: 180–189.
172. Purdom E, Simpson KM, Robinson MD, Conboy JG, Lapuk AV, Speed TP. FIRMA: a method for detection of alternative splicing from exon array data. *Bioinformatics.* 2008;24: 1707–1714.
- 1345 173. Anders S, Reyes A, Huber W. Detecting differential usage of exons from RNA-seq data. *Genome Res.* 2012;22: 2008–2017.
174. Sonesson 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.
- 1350 175. Love MI, Sonesson C, Patro R. Swimming downstream: statistical analysis of differential transcript usage following Salmon quantification. *F1000Res.* 2018;7. doi:10.12688/f1000research.15398.2
176. Nowicka M, Robinson MD. DRIMSeq: a Dirichlet-multinomial framework for multivariate count outcomes in genomics. *F1000Res.* 2016;5: 1356.
- 1355 177. Li YI, Knowles DA, Humphrey J, Barbeira AN, Dickinson SP, Im HK, et al. Annotation-free quantification of RNA splicing using LeafCutter. *Nat Genet.* 2017; doi:10.1038/s41588-017-0004-9
178. Papastamoulis P, Rattray M. Bayesian estimation of differential transcript usage from RNA-seq data. *Stat Appl Genet Mol Biol.* 2017;16: 367–386.
- 1360 179. Froussios K, Mourão K, Simpson GG, Barton GJ, Schurch NJ. Identifying differential isoform abundance with RATs: a universal tool and a warning [Internet]. *bioRxiv.* 2017. p. 132761. doi:10.1101/132761
180. Shen S, Park JW, Lu Z-X, Lin L, Henry MD, Wu YN, et al. rMATS: Robust and flexible detection of differential alternative splicing from replicate RNA-Seq data. *Proceedings of the National Academy of Sciences.* 2014;111: E5593–E5601.
- 1365 181. Trincado JL, Entizne JC, Hysenaj G, Singh B, Skalic M, Elliott DJ, et al. SUPPA2: fast, accurate, and uncertainty-aware differential splicing analysis across multiple conditions. *Genome Biol.* 2018;19: 40.
- 1370 182. Yi L, Pimentel H, Bray NL, Pachter L. Gene-level differential analysis at transcript-level resolution. *Genome Biol.* 2018;19: 53.
183. Hicks SC, Townes FW, Teng M, Irizarry RA. Missing data and technical variability in single-cell RNA-sequencing experiments. *Biostatistics.* 2017; doi:10.1093/biostatistics/kxx053
- 1375 184. Marinov GK, Williams BA, McCue K, Schroth GP, Gertz J, Myers RM, et al. From single-cell to cell-pool transcriptomes: stochasticity in gene expression and RNA splicing. *Genome Res.* 2014;24: 496–510.
185. Svensson V, Vento-Tormo R, Teichmann SA. Exponential scaling of single-cell RNA-seq in the past decade. *Nat Protoc.* 2018;13: 599–604.
- 1380 186. Wagner A, Regev A, Yosef N. Revealing the vectors of cellular identity with single-cell genomics. *Nat Biotechnol.* 2016;34: 1145–1160.
187. Moor AE, Itzkovitz S. Spatial transcriptomics: paving the way for tissue-level systems biology. *Curr Opin Biotechnol.* 2017;46: 126–133.
188. Kumar P, Tan Y, Cahan P. Understanding development and stem cells using single cell-based analyses of gene expression. *Development.* 2017;144: 17–32.
- 1385 189. Kang HM, Subramaniam M, Targ S, Nguyen M, Maliskova L, McCarthy E, et al. Multiplexed droplet single-cell RNA-sequencing using natural genetic variation. *Nat*

- Biotechnol. 2018;36: 89–94.
- 1390 190. Paulson KG, Voillet V, McAfee MS, Hunter DS, Wagener FD, Perdicchio M, et al. Acquired cancer resistance to combination immunotherapy from transcriptional loss of class I HLA. *Nat Commun.* 2018;9: 3868.
191. Giladi A, Amit I. Single-Cell Genomics: A Stepping Stone for Future Immunology Discoveries. *Cell.* 2018;172: 14–21.
- 1395 192. Trapnell C. Defining cell types and states with single-cell genomics. *Genome Res.* 2015;25: 1491–1498.
193. Sandberg R. Entering the era of single-cell transcriptomics in biology and medicine. *Nat Methods.* 2014;11: 22–24.
194. Soneson C, Robinson MD. Bias, robustness and scalability in single-cell differential expression analysis. *Nat Methods.* 2018;15: 255–261.
- 1400 195. Jaakkola MK, Seyednasrollah F, Mehmood A, Elo LL. Comparison of methods to detect differentially expressed genes between single-cell populations. *Brief Bioinform.* 2017;18: 735–743.
196. Kharchenko PV, Silberstein L, Scadden DT. Bayesian approach to single-cell differential expression analysis. *Nat Methods.* 2014;11: 740–742.
- 1405 197. Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, et al. MAST: A flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA-seq data. *Genome Biol. Genome Biology;* 2015;16: 278.
198. Van den Berge K, Perraudeau F, Soneson C, Love MI, Risso D, Vert J-P, et al. Observation weights unlock bulk RNA-seq tools for zero inflation and single-cell applications. *Genome Biol.* 2018;19: 24.
- 1410 199. Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert J-P. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat Commun.* 2018;9: 284.
200. Eling N, Richard AC, Richardson S, Marioni JC, Vallejos CA. Correcting the Mean-Variance Dependency for Differential Variability Testing Using Single-Cell RNA Sequencing Data. *Cell Syst.* 2018; doi:10.1016/j.cels.2018.06.011
- 1415 201. Korthauer KD, Chu L-F, Newton MA, Li Y, Thomson J, Stewart R, et al. A statistical approach for identifying differential distributions in single-cell RNA-seq experiments. *Genome Biol.* 2016;17: 222.
202. Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing: computational challenges and solutions. *Nat Rev Genet.* 2012;13: 36–46.
- 1420 203. Weirather JL, de Cesare M, Wang Y, Piazza P, Sebastiano V, Wang X-J, et al. Comprehensive comparison of Pacific Biosciences and Oxford Nanopore Technologies and their applications to transcriptome analysis. *F1000Res.* 2017;6: 100.
204. Steijger T, Abril JF, Engström PG, Kokocinski F, The RGASP Consortium, Hubbard TJ, et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Methods. The Author(s);* 2013;10: 1177.
- 1425 205. Tilgner H, Grubert F, Sharon D, Snyder MP. Defining a personal, allele-specific, and single-molecule long-read transcriptome. *Proc Natl Acad Sci U S A. National Academy of Sciences;* 2014;111: 9869–9874.
- 1430 206. Gonzalez-Garay ML. Introduction to Isoform Sequencing Using Pacific Biosciences Technology (Iso-Seq). In: Wu J, editor. *Transcriptomics and Gene Regulation.* Dordrecht: Springer Netherlands; 2016. pp. 141–160.
207. Laver T, Harrison J, O'Neill PA, Moore K, Farbos A, Paszkiewicz K, et al. Assessing the performance of the Oxford Nanopore Technologies MinION. *Biomol Detect Quantif.* 2015;3: 1–8.
- 1435

208. Ross MG, Russ C, Costello M, Hollinger A, Lennon NJ, Hegarty R, et al. Characterizing and measuring bias in sequence data. *Genome Biol.* 2013;14: R51.
209. Teng JLL, Yeung ML, Chan E, Jia L, Lin CH, Huang Y, et al. PacBio But Not Illumina Technology Can Achieve Fast, Accurate and Complete Closure of the High GC, Complex Burkholderia pseudomallei Two-Chromosome Genome. *Front Microbiol.* 2017;8: 1448.
- 1440 210. Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, et al. Real-time DNA sequencing from single polymerase molecules. *Science.* 2009;323: 133–138.
211. Levene MJ, Korfach J, Turner SW, Foquet M, Craighead HG, Webb WW. Zero-mode waveguides for single-molecule analysis at high concentrations. *Science.* 2003;299: 682–686.
- 1445 212. Travers KJ, Chin C-S, Rank DR, Eid JS, Turner SW. A flexible and efficient template format for circular consensus sequencing and SNP detection. *Nucleic Acids Res.* 2010;38: e159.
213. Carneiro MO, Russ C, Ross MG, Gabriel SB, Nusbaum C, DePristo M a. Pacific biosciences sequencing technology for genotyping and variation discovery in human data. *BMC Genomics.* *BMC Genomics;* 2012;13: 375.
- 1450 214. Wang Y, Yang Q, Wang Z. The evolution of nanopore sequencing. *Front Genet.* 2014;5: 449.
215. Quick J, Quinlan AR, Loman NJ. A reference bacterial genome dataset generated on the MinION™ portable single-molecule nanopore sequencer. *Gigascience.* 2014;3: 22.
- 1455 216. Garalde DR, Snell EA, Jachimowicz D, Sipos B, Lloyd JH, Bruce M, et al. Highly parallel direct RNA sequencing on an array of nanopores. *Nat Methods.* 2018; doi:10.1038/nmeth.4577
217. Smith AM, Jain M, Mulrone L, Garalde DR, Akeson M. Reading canonical and modified nucleotides in 16S ribosomal RNA using nanopore direct RNA sequencing [Internet]. *bioRxiv.* 2017. p. 132274. doi:10.1101/132274
- 1460 218. Sharon D, Tilgner H, Grubert F, Snyder M. A single-molecule long-read survey of the human transcriptome. *Nat Biotechnol.* 2013;31: 1009–1014.
219. Oikonomopoulos S, Wang YC, Djambazian H, Badescu D, Ragoussis J. Benchmarking of the Oxford Nanopore MinION sequencing for quantitative and qualitative assessment of cDNA populations. *Sci Rep.* 2016;6: 31602.
- 1465 220. Abdel-Ghany SE, Hamilton M, Jacobi JL, Ngam P, Devitt N, Schilkey F, et al. A survey of the sorghum transcriptome using single-molecule long reads. *Nat Commun.* 2016;7: 11706.
221. Au KF, Sebastiano V, Afshar PT, Durruthy JD, Lee L, Williams B a., et al. Characterization of the human ESC transcriptome by hybrid sequencing. *Proc Natl Acad Sci U S A.* 2013;110: E4821–30.
- 1470 222. Treutlein B, Gokce O, Quake SR, Südhof TC. Cartography of neurexin alternative splicing mapped by single-molecule long-read mRNA sequencing. *Proc Natl Acad Sci U S A.* 2014;111: E1291–9.
- 1475 223. Rhoads A, Au KF. PacBio Sequencing and Its Applications. *Genomics Proteomics Bioinformatics.* 2015;13: 278–289.
224. Marchet C, Lecompte L, Da Silva C, Cruaud C, Aury JM, Nicolas J, et al. De novo Clustering Nanopore Long Reads of Transcriptomics Data by Gene [Internet]. *bioRxiv.* 2017. p. 170035. doi:10.1101/170035
- 1480 225. Workman RE, Myrka AM, Wong GW, Tseng E, Welch KC Jr, Timp W. Single-molecule, full-length transcript sequencing provides insight into the extreme metabolism of the ruby-throated hummingbird *Archilochus colubris*. *Gigascience.* 2018;7: 1–12.
226. An D, Cao HX, Li C, Humbeck K, Wang W. Isoform Sequencing and State-of-Art

- 1485 Applications for Unravelling Complexity of Plant Transcriptomes. *Genes* . 2018;9.
doi:10.3390/genes9010043
227. Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, et al. Widespread Polycistronic Transcripts in Fungi Revealed by Single-Molecule mRNA Sequencing. *PLoS One*. 2015;10: e0132628.
- 1490 228. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;
doi:10.1093/bioinformatics/bty191