



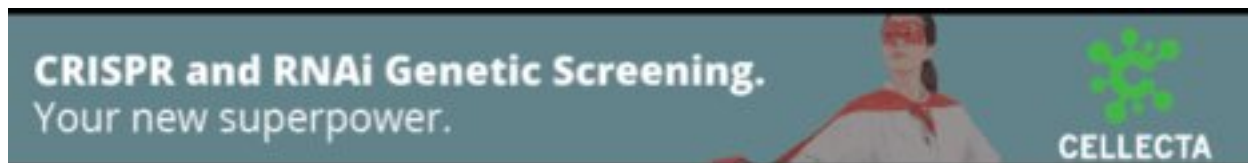
TransRate: reference free quality assessment of de novo transcriptome assemblies

Richard Smith-Unna, Chris Boursnell, Rob Patro, et al.

Genome Res. published online June 1, 2016

Access the most recent version at doi:[10.1101/gr.196469.115](https://doi.org/10.1101/gr.196469.115)

P<P	Published online June 1, 2016 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Title**

2 TransRate: reference free quality assessment of *de-novo* transcriptome assemblies

3 **Authors**

4 Richard Smith-Unna¹, Chris Bournnell¹, Rob Patro², Julian M Hibberd¹, and Steven Kelly^{3*}

5 **Affiliations**

6 1) Department of Plant Sciences, University of Cambridge, Downing Street, CB2 3EA, UK

7 2) Department of Computer Science, Stony Brook University, Stony Brook, NY 11794-4400, USA

8 3) Department of Plant Sciences, University of Oxford, South Parks Road, OX1 3RB, UK

9 **Corresponding Author**

10 * email: steven.kelly@plants.ox.ac.uk, phone: 0044 (0) 1865 275123

11 **Running title:**

12 Transcriptome assembly accuracy evaluation

13 **Keywords:**

14 *de novo* assembly evaluation, transcriptome assembly

15 **Abstract**

16 TransRate is a tool for reference-free quality assessment of *de novo* transcriptome assemblies. Using
17 only the sequenced reads and the assembly as input, we show multiple common artifacts of *de novo*
18 transcriptome assembly can be readily detected. These include chimeras, structural errors, incomplete
19 assembly and base errors. TransRate evaluates these errors to produce a diagnostic quality score for
20 each contig and these contig scores are integrated to evaluate whole assemblies. Thus TransRate
21 can be used for *do novo* assembly filtering and optimisation as well as comparison of assemblies
22 generated using different methods from the same input reads. Applying the method to a dataset of 155
23 published *de novo* transcriptome assemblies we deconstruct the contribution that assembly method,
24 read length, read quantity, and read quality make to the accuracy of *de novo* transcriptome
25 assemblies and reveal that variance in the quality of the input data explains 43% of the variance in the

26 quality of published *de novo* transcriptome assemblies. As TransRate is reference-free it is suitable for
27 assessment of assemblies of all types of RNA, including assemblies of long non-coding RNA, rRNA,
28 mRNA, and mixed RNA samples.

29 **Introduction**

30 High-throughput sequencing of RNA has revolutionized our ability to assess the genetic and
31 quantitative basis of many complex biological traits. For organisms that have sequenced and
32 annotated genomes, short reads can be directly mapped to these resources and quantitative
33 estimates of gene expression (as well as splice-variants and mutations) can be determined using a
34 variety of different methods. In the absence of an appropriate reference genome, *de novo*
35 transcriptome assembly must be performed. These assemblies provide the primary data for gene
36 discovery and evolutionary analyses, and facilitate quantitative assessment of differential gene
37 expression. Given the importance of these applications to comparative biological research several
38 algorithms have been developed to produce *de novo* transcriptome assemblies from raw sequence
39 data. Popular amongst these algorithms are Trinity (Grabherr et al., 2011), Oases (Schulz et al.,
40 2012), Trans-ABYSS (Robertson et al., 2010), IDBA-tran (Peng et al., 2013) and SOAPdenovo-Trans
41 (Xie et al., 2014), each of which takes a different approach to the problem of reconstituting a
42 transcriptome from short sequence reads. Furthermore, they all provide considerable flexibility with
43 multiple parameters and heuristics that can be modified to allow the user to tailor their assembly
44 settings for variations in RNA-seq library construction, coverage depth and differences between
45 organisms. These large parameter spaces mean that the same read data can generate substantially
46 different assemblies both within and between assembly methods. Likewise, altering parameter
47 combinations can result in the assembly of contigs with varying properties such that disparate
48 conclusions relating to transcript content and expression level can be reached from the same input
49 data.

50 In addition to the considerable algorithmic flexibility, the data being assembled can be generated from
51 multiple different RNA types. These can range from specifically amplified sub-populations of particular

52 types of RNA, to total RNA encompassing all RNA types within the cell. Given the wide range of input
53 data and assembly methods there is a need to be able to evaluate the quality of any *de novo*
54 transcriptome in the absence of a known reference and identify the set of parameters, or assembly
55 methods that best reconstruct the transcriptome from which the raw read data was generated.
56 Moreover, there is a need to be able to identify within a given assembly the set of contigs that are well-
57 assembled from those that are not, so that incorrect data do not influence downstream biological
58 interpretation.

59 Algorithms to assess the outputs of DNA-directed (e.g. genome and meta-genome) assembly have
60 been developed. These range in complexity from descriptive metrics (Gurevich et al., 2013), to explicit
61 modelling of the sequencing and assembly process to provide a likelihood-based measure of
62 assembly quality (Clark et al., 2013; Rahman and Pachter, 2013). However, the assumptions used for
63 evaluation of DNA-directed assembly such as uniformity of coverage (except in repetitive regions) and
64 assembled contig length are not appropriate for the assembly of transcriptomes due to the
65 exponentially distributed coverage of different transcripts and log-normally distributed transcript
66 lengths. Therefore alternative criteria that are tailored for the biological properties of transcriptomes
67 need to be used for the assessment of *de novo* assembled transcriptomes.

68 To date the majority of *de novo* transcriptome assessment methods have exploited comparative
69 approaches in which the assembled transcriptome is compared to a known reference dataset (Lowe et
70 al., 2014; O'Neil and Emrich, 2013). These comparative methods provide insight into the complement
71 of known proteins that are represented within a *de novo* assembly but do not reveal the extent to
72 which the contigs representing those proteins are assembled correctly. Furthermore, due to the
73 inherent limitations of such comparative analyses they only assess the *de novo* transcriptome on the
74 subset of contigs that represent conserved proteins. Highly divergent transcripts, novel transcripts,
75 and non-coding transcripts are not assessed by these methods and thus the assessment measures do
76 not consider all of the data. Moreover, *de novo* assembly of non-coding RNA or specific sub-
77 populations of RNAs are poorly evaluated by these comparative methods. To date only a single

78 reference-free transcriptome assembly evaluation tool has been produced, RSEM-eval (Li et al.,
79 2014). RSEM-eval provides an assembly likelihood given the read data, allowing the comparison of
80 assemblies generated from the same input data. Although RSEM-eval quantifies the relative
81 contribution that each contig makes to an overall assembly score, it does not provide descriptive
82 statistics about the quality of contigs within an assembly.

83 Here we present TransRate, a novel method for evaluation of the accuracy and completeness of *de*
84 *novo* transcriptome assemblies. TransRate assesses these features through two novel reference-free
85 statistics: the TransRate contig score and the TransRate assembly score. The TransRate contig score
86 provides a quantitative measure of the accuracy of assembly for each individual contig and The
87 TransRate assembly score provides a quantitative measure of the accuracy and completeness of the
88 assembly.

89 **Results**

90 **Problem definition and approach**

91 The aim of *de novo* transcriptome assembly is to accurately reconstruct the complete set of transcripts
92 that are represented in the read data in the absence of a reference genome. There are several
93 contributing factors that negatively affect the accuracy of this reconstruction process. These factors
94 include error in the sequencing process, incomplete coverage of transcripts (due to insufficient
95 sequencing depth), and real biological variability (such as variation in exon/intron retention, variation in
96 exon boundary usage, and variation in nucleotide sequence between alleles). Moreover, assembly
97 errors can originate from algorithmic simplifications (such as representing the information contained in
98 the reads as shorter words) and allowances (e.g. permitting assembly of fragments containing mis-
99 matches) that are used to mitigate the computational complexity of the assembly problem. Together,
100 these factors cause several common assembly artifacts including hybrid assembly of gene families,
101 transcript fusion (chimerism), spurious insertions in contigs, and structural abnormalities such as
102 incompleteness, fragmentation and local mis-assembly of contigs (Figure 1).

103 TransRate is focused on a clear problem definition; to assess the accuracy and completeness of a *de*
104 *novo* assembled transcriptome using only the input reads. TransRate proceeds by mapping the reads
105 to the assembled contigs, proportionally assigning multi-mapping reads in a probabilistic manner to
106 their contig of origin, analyzing the alignments, calculating contig level metrics (Table 1), integrating
107 these contig level metrics to provide a contig score, and then combining the completeness of the
108 assembly with the score of each contig to produce an overall assembly score (Figure 2). TransRate
109 also provides an abundance weighted assembly score which weights each constituent contig score by
110 the relative abundance level of each contig.

111 **Contig assessment criteria**

112 To calculate the TransRate contig score a correctly assembled contig is assumed to have the
113 following four properties. 1) The identity of the nucleotides in the contig will accurately represent the
114 nucleotides of the true transcript. 2) The number of nucleotides in the contig (i.e. the assembled
115 transcript length) will accurately represent the number in the true transcript. 3) The order of the
116 nucleotides in the contig will accurately represent the order in the true transcript. 4) The contig will
117 represent a single transcript. We propose that each of these four statements can be approximated
118 through analysis of the reads that map to the assembled contigs and are encapsulated by the four
119 metrics presented in Table 1. For a detailed description of these metrics and how they are calculated
120 see the TransRate contig score section in the Methods.

121 To determine whether these four contig level metrics were discrete, and thus captured different
122 properties of each assembled contig, their performance was evaluated on a range of assemblies
123 generated using different algorithms from multiple different species (Figure 3A). For each contig level
124 metric the distributions of observed scores was broadly similar irrespective of species or assembly
125 algorithm (Figure 3A). One notable exception to this observation is that the distribution of $s(C_{cov})$
126 (Table 1) observed for rice and mouse contigs generated using SOAPdenovo-Trans was markedly
127 different to that observed for Oases and Trinity for the same species. This reveals that the contigs
128 generated using SOAPdenovo-Trans on this rice data contained fewer regions that had zero coverage
129 after read mapping.

130 Visual inspection of the global behavior of the contig level metrics suggested that the four scores could
131 be classified into two groups based on the density function of the observed score values. Both $s(C_{ord})$
132 and $s(C_{seg})$ (Table 1) produced approximately uniform distributions spanning the entire score range
133 (Figure 3A), whereas $s(C_{cov})$ and $s(C_{nuc})$ (Table 1) produced distributions whose density increased
134 towards higher values (Figure 3A). To determine if these visually similar distributions were correlated,
135 and thus measured features of the assembled contigs that were inter-dependent, we analyzed the
136 pairwise Spearman's rank correlation between the score components. This revealed that the metrics
137 were poorly correlated (Figure 3B) and thus each provided discrete assessment of the assembled
138 contigs to which they were applied.

139 Manual inspection of reference-based results for the 30 lowest-scoring contigs according to each
140 score component was consistent with the individual score components capturing their target properties
141 (Supplemental Fig S1). The Bayesian segmentation of coverage depth, $s(C_{seg})$, was also evaluated by
142 inspection of coverage depth profiles (Supplemental Fig S2) and simulation of artificial transcript
143 chimeras. The latter was done by *in silico* fusion of randomly selected transcripts from the yeast
144 transcriptome and assessment of $s(C_{seg})$ scores as a function of the difference in abundance between
145 the fused transcripts (Supplemental Fig S3). Here, the segmentation method was unable to distinguish
146 chimeras between transcripts whose abundance differed by < 2 fold (Supplemental Fig S3). The
147 individual score components are provided in the TransRate program output so that end-users can gain
148 insight into the common sources of error in their assembly.

149 **Evaluation of the TransRate contig score**

150 As the contig-level metrics provided discrete evaluation of assembled contigs, we sought to determine
151 if the geometric mean of these metrics (see Methods equation [1]) was informative of the accuracy of
152 assembly. To assess this, 4 million read pairs were simulated from each of the four test species (rice,
153 mouse, human, and yeast, see Independence of score components) and assembled using
154 SOAPdenovo-Trans with default settings. Simulated reads were used here so that the true set of
155 transcripts was known and hence the accuracy of the assembled contigs could be assessed. The

156 resultant assemblies were subjected to TransRate assessment, and the utility of the TransRate contig
157 scores was assessed by comparing them to a conventional measure of contig accuracy calculated by
158 alignment of the assembled contigs to the transcripts used to simulate the reads (see Calculation of
159 contig accuracy). Comparison of these measures revealed that there was a strong monotonic
160 relationship between contig accuracy and TransRate contig score (Figure 4A). Across all simulated
161 datasets, the TransRate contig score exhibited a Spearman's rank correlation with contig accuracy of
162 $\rho = 0.71$ (Figure 4A, Supplemental Table S1). For comparison we also applied RSEM-eval to the
163 same dataset (Figure 4B). Here, the contig impact score from RSEM-eval, which measures the
164 relative contribution of every contig to the assembly score, also showed a positive correlation with
165 contig accuracy, however the Spearman's rank correlation with accuracy was lower than that
166 observed for TransRate ($\rho = 0.36$, Supplemental Table S1). Non-parametric correlation measures
167 were used here to enable unbiased comparison of TransRate and RSEM-eval scores, as their score
168 distributions differ in type, location, scale and shape.

169 Analysis of the interrelationship between contig scores and contig accuracy revealed that both
170 assessment methods exhibited minimum value inflation (Figure 4A & B). Though some of these
171 minimum value contigs comprise accurately assembled transcript sequences, they are assigned
172 minimum score values as they fail to acquire mapped reads during the read-mapping process. This
173 occurs due to the presence of contigs within the assembly that better represent the true contig than
174 the contig in question and thus preferentially obtain all of the mapped reads during the probabilistic
175 read assignment stage. This phenomenon commonly occurs when the contig in question is a substring
176 of longer contig in the assembly. As these contigs are redundant and they would be quantified as "not
177 expressed" in downstream expression analyses of the assemblies, both TransRate and RSEM-eval
178 are justified in the assignment of minimum value scores to these contigs. In the absence of these
179 minimum value contigs the Spearman's correlation coefficients for both TransRate and RSEM-eval are
180 $\rho = 0.70$ and $\rho = 0.77$ respectively.

181 **Application of TransRate for relative evaluation of *de novo* assemblies from the same**
182 **read data**

183 Given that the TransRate contig score is strongly related to contig accuracy, we sought to develop an
184 assembly-level score that summarised the information captured by assessment the individual contigs
185 (Figure 4A). Here, the geometric mean of all contig scores was selected such that each contig
186 contributed equally to the final assembly assessment (see Methods equation [2]). Analysis of the
187 TransRate contig score distributions for assemblies generated using different assembly algorithms
188 from different species revealed that most assemblers produced contigs that obtained a wide range of
189 scores (Figure 5A). Some distributions also appeared to be multi-modal with overlapping populations
190 of low and high scoring contigs (Figure 5A).

191 Comparison of the geometric mean of the contig scores revealed that on different datasets, different
192 assemblers tended to produce more accurate assemblies (Figure 5B). On average, Oases (version
193 0.2.06 with Velvet version 1.2.07) produced the highest mean contig scores for mouse and rice, while
194 Trinity (version Trinity-r2013-02-25) produced the highest mean contig scores for human and yeast
195 (Figure 5B). The percentage of the input that could be mapped to these assemblies ranged from 65-
196 85% and thus significant amounts of read data failed to be assembled by each method (Figure 5C). To
197 provide a single assembly assessment score that combined the proportion of read data contained
198 within the assembly and the mean accuracy of the constituent contigs we took the product of the
199 geometric mean contig score and the proportion of reads mapping to the assembly (Figure 5D). This
200 assembly score places equal importance on the accuracy of each of the assembled contigs and the
201 proportion of the input read data that is captured by the *de novo* assembly. In an ideal scenario where
202 all of the input reads map back to the assembled contigs with no disagreement between the reads and
203 the assembly the assembly score will be 1. Errors in the sequencing or assembly process that cause
204 reads to be omitted from the assembly or reads to disagree with the assembled contigs will cause the
205 assembly score to tend towards 0.

206 TransRate also provides an abundance-weighted contig score (see Methods equation [3]) where
207 transcripts with assembly errors are penalized in proportion to their abundance. That is, highly

208 abundant transcripts with errors are penalized more heavily than low abundance transcripts with the
209 same errors. Using these abundance-weighted contigs scores an abundance-weighted assembly
210 score can also be evaluated (see Methods equation [4]). The results from using these abundance-
211 weighted scores exhibit the same trend as for the TransRate contig and assembly scores
212 (Supplemental Fig S4). However, the additional penalty due to abundance weighting causes the
213 overall scores to be much lower (Supplemental Fig S4). Caution should be exercised by the user
214 when using the abundance-weighted contig scores as they are not comparable between contigs. That
215 is, a highly abundant transcript with an assembly error will have a lower score than a transcript with
216 the same error that is expressed to a lower level.

217 **Further comparison of *de novo* assemblies using BLAST and TransRate**

218 To demonstrate additional ways in which TransRate can be combined with BLAST based assessment
219 of *de novo* transcriptome assemblies, the *de novo* assemblies was annotated using reciprocal best
220 BLAST (bi-directional best BLAST hit) against the appropriate Ensembl reference dataset for that
221 species. The TransRate scores for these contigs were compared and the proportion of transcripts that
222 had the highest TransRate score for each assembly was recorded (Figure 5E). No one method
223 consistently outperformed the others, rather the different assemblers produced the best assembly for
224 >25% of transcripts (Figure 5E). Analysis of the total number of reference transcripts that were
225 assembled by the different methods revealed that, though there was significant agreement between
226 the methods, each method uniquely assembled a large number of *bona fide* transcripts not assembled
227 by the other methods (Figure 5F). Taken together these analyses lend support to the idea that
228 combining contigs from multiple assembly methods is an effective way to increase the completeness
229 of a *de novo* assembled transcriptome.

230 **Filtration of contigs using TransRate contig scores**

231 As shown in Figure 4A, 4B & 5A, many contigs within a given assembly can achieve low or minimum
232 value scores and thus users may desire to remove them from the assembly. While TransRate allows
233 the user to specify any contig score cut-off between 0 and 1 for filtration of assembled contigs, it also
234 provides an alternative option whereby a specific contig score cut-off can be learned for any given

235 assembly. To do this TransRate uses a global optimisation method to find the contig score cut-off
236 value such that the TransRate assembly score function is maximised (Supplemental Fig S5). This
237 automated cut-off method is consistent with the problem definition and overall aim of TransRate (to
238 assess the accuracy and completeness of a *de novo* assembled transcriptome using only the input
239 reads) as it automatically selects the subset of contigs that maximises both accuracy and
240 completeness. It should be noted that filtering contigs in this way may remove some accurately
241 assembled low abundance transcripts that have incomplete coverage.

242 To provide an example of the results obtained from the application of the automated TransRate contig
243 filtering the 10 assemblies analyzed in figure 5 above were subject to filtering. Those *de novo*
244 assembled contigs that contained regions with >95% identity to predicted genes in the genomes of the
245 source species were selected for further analysis. On average 20% of genes that had contigs
246 matching at least part of a predicted gene were filtered out by TransRate (Supplemental Fig S6). Of
247 the genes whose entire length was encompassed in a single transcript ~12% were discarded by
248 TransRate (Supplemental Fig S6). Although TransRate has identified these transcripts as poorly
249 assembled, and caution should be exercised against using abundance level estimates for these
250 contigs, they may contain regions that have utility in certain analyses (e.g. phylogenetic analysis).

251 **Comparative analysis of 155 published assemblies provides a reference for calibration** 252 **and relative assessment of assembly quality**

253 To provide a reference distribution of TransRate assembly scores that end-users can use to assess
254 the relative merit of their own assemblies, TransRate was applied to a set of 155 published *de novo*
255 assembled transcriptomes (Supplemental Table S2). All assembled transcriptomes were downloaded
256 from the NCBI Transcriptome Shotgun Archive (<http://www.ncbi.nlm.nih.gov/genbank/tsa>) and were
257 chosen for analysis if they met the following criteria: 1) The assembly program was listed; 2) The
258 reads were Illumina paired-end reads; 3) The published assembly contained at least 5,000 contigs.
259 TransRate assembly scores for this set of published assemblies ranged from 0.001 to 0.52 (Figure 6A,
260 red line). Each assembly was also subject to automated assembly score optimisation producing
261 optimised assembly scores that ranged from 0.001 to 0.6 (Figure 6A, teal line). Although some

262 assembly scores showed little or no change following removal of low scoring transcripts, most
263 improved when contigs below the learned cut-off were discarded (Figure 6B).

264 It has been suggested that the transcriptomes from certain groups of organisms may be more difficult
265 to assemble than others (Martin and Wang, 2011). To investigate whether TransRate assembly scores
266 varied for different taxa the results were analyzed according to their major phylogenetic groups (Figure
267 6C). For clades with more than 10 representative assemblies no association between assembly
268 quality and taxonomic group was found (Figure 6C).

269 To determine if any assembler consistently produced higher TransRate assembly scores on end-user
270 datasets, the performance of methods that had at least 10 assemblies was compared (Figure 6D). In
271 this test Trinity, Oases, and SOAPdenovo-Trans all produced assemblies that spanned similar score
272 ranges, with the highest mean score exhibited by Trinity (Figure 6D). In contrast, Newbler, Agalma
273 and Trans-ABYSS assemblies produced lower TransRate scores (Figure 6D). However, caution
274 should be exercised when interpreting these results as the user-modifiable settings and post-assembly
275 processing steps were not reported for these published assemblies. Thus the extent to which the
276 TransRate assembly scores were influenced by changes in user-modifiable assembly parameters or
277 post-assembly processing is unknown.

278 Given that neither assembly method nor taxonomic group produced a major effect on the TransRate
279 score of an assembly we sought to determine whether the quality of the input read data was
280 responsible for some of the variation in TransRate assembly scores. The read data for each assembly
281 was analyzed using FastQC and the resulting read-level metrics compared to the TransRate assembly
282 scores of the assemblies generated using those reads. This revealed that neither the read length nor
283 the percentage GC of the read dataset exhibited any correlation with TransRate assembly score
284 (Figure 6E & F). However, significant associations were observed for both read quality ($r^2 = 0.27$,
285 Figure 6G) and the level of read-duplication in the dataset ($r^2 = 0.1$, Figure 6H). In Illumina sequencing,
286 low read qualities are predominantly caused by errors in the sequencing process, common sources
287 include over-clustering of the flow cell and phasing. In contrast, increases in read-duplication is

288 caused by errors in the sample preparation stage. It occurs during the PCR amplification stage of the
289 read library preparation and is generally caused by either conducting the library preparation from too
290 little starting material, or by having a large variance in the fragment size such that smaller fragments
291 become over-represented during the limited cycle PCR. While there is little correlation between the
292 number of sequenced reads and the TransRate score of the assembled transcriptome (Figure 6I)
293 there is a clear association between the relative coverage implied by those reads and the TransRate
294 score ($r^2 = 0.16$, Figure 6J). In summary, the quality of the sequence reads, the number of reads per
295 gene and the quality of the input cDNA library (in order of relative contribution) explain 43% of the
296 variance in *de novo* assembly quality. Thus, the quality of the input data is more important in
297 determining the quality of a *de novo* assembly than the choice of assembly method that is used.

298 **Discussion**

299 Here we present TransRate a novel method for reference free assessment and filtering of *de novo*
300 assembled transcriptomes. Our method is focused on a clear definition of an optimal *de novo*
301 assembled transcriptome, that it should be a complete and accurate representation of the transcripts
302 encompassed in the raw read data. TransRate avoids conflating assessment of *de novo* assembly
303 quality with other criteria (such as coverage of expected reference transcript subsets) that are not
304 equivalent to correct or complete assembly of the input reads. Moreover, the method is not biased by
305 expression level of the transcripts and each transcript is weighted equally in the overall transcriptome
306 assessment (unless the alternative abundance weighted metric is used). As the majority of published
307 *de novo* assembled transcriptomes use Illumina paired-end sequencing, our analysis of the efficacy of
308 TransRate is focused on this data type. However, the method is suitable for the analysis of other types
309 of sequencing and thus is not restricted to use in the analysis of Illumina data.

310 TransRate is specifically designed to provide detailed insight into the quality of any *de novo*
311 assembled transcriptome and each of its constituent contigs such that comparative analysis between
312 assembly methods and post-assembly filtering of good and bad contigs can be performed. As
313 TransRate does not use reference datasets in the evaluation of assemblies it is equally suitable for the

314 assessment of assemblies of all types of RNA, including long non-coding RNA, mRNA, ribosomal
315 RNA and mixed RNA samples. Moreover, given multiple assemblies generated using the same input
316 reads, TransRate can also be used to determine the assembly that best represents the input read
317 data. Thus TransRate could be used to help improve the performance of multiple different *de novo*
318 transcriptome assembly algorithms. TransRate can also be used to filter out low scoring contigs,
319 however caution should be exercised here as application of filtering may result in removal of
320 transcripts that have some utility. For example transcripts with very low coverage are more likely to
321 have low contig scores because of fragmentation and encapsulated bases in gapped regions, these
322 transcripts while incompletely assembled may have utility in pathway reconstruction, quantitative
323 expression analysis or phylogenetic analysis. Similarly, transcripts with low $s(C_{seg})$ scores are likely to
324 represent chimeric transcripts. Here, although the transcript itself may be incorrectly assembled the
325 component segments of the transcript may themselves be correctly assembled and of utility if
326 separated. To help users identify and diagnose likely assembly errors affecting low scoring contigs
327 TransRate provides each of the separate contig scores (in addition to the overall contig score). This
328 information can be used to help resolve assembly errors on a contig by contig basis. Further
329 investigation by systematically exploring a large range of read mapping parameters across a large
330 range of read mapping algorithms and assembly tools may yield new ways to improve the
331 performance of TransRate. This may improve the $s(C_{ord})$ and $s(C_{cov})$ measures that are affected by
332 read coherency (Myers, 2005), which may, in turn, suggest how the assemblies could be improved.

333 To help end users to interpret the TransRate scores that they obtain for their own assemblies and
334 place them in context of previously published assemblies, we provide a meta-analysis of 155
335 published *de novo* assemblies. Here, a user generated *de novo* assembly with a TransRate score of
336 0.22 (optimised score of 0.35) would be better than 50% of published *de novo* assembled
337 transcriptomes that have been deposited in the NCBI TSA. Through detailed analysis of these 155
338 published assemblies we reveal that the quality of the input read data is the major factor determining

339 the quality of any *de novo* transcriptome assembly, explaining more of the variance in quality between
340 assemblies than quantity of read data or assembly method that is used.

341 **Methods**

342 **Algorithm overview**

343 TransRate is a reference-free qualitative assessment tool for the analysis of *de novo* transcriptome
344 assemblies. TransRate requires one or more transcriptome assemblies and the reads used to
345 generate those assemblies. TransRate aligns the reads to the assembly, processes those read
346 alignments, and calculates contig scores using the full set of processed read alignments. TransRate
347 classifies contigs into those that are well assembled and those that are poorly assembled, by learning
348 a score cutoff from the data that maximises the overall assembly score.

349 **Read alignment**

350 Reads are aligned to a given assembly using SNAP v1.0.0 (Zaharia et al., 2011). Alignments are
351 reported up to a maximum edit distance of 30. Up to 10 multiple alignments are reported per read
352 where available (-omax 10), up to a maximum edit distance of 5 from the best-scoring alignment (-om
353 5). Exploration within an edit distance of 5 from each alignment is allowed for the calculation of MAPQ
354 scores (-D 5). BAM-format alignments produced by SNAP are processed by Salmon (Patro et al.,
355 2015). Separate strategies are employed for abundance estimation and posterior read assignment.
356 For abundance estimation, each mapped read is fractionally assigned to each potential contig of origin
357 using Salmon (Patro et al., 2015) in a process that is analogous to the proportional assignment of the
358 EM procedure used in RSEM (Li et al., 2010). For contig score evaluation a different approach was
359 taken in which a single assignment was produced for each read. Here each read was assigned
360 entirely to a single contig, but the probability of assignment for multi-mapping reads was sampled from
361 the distribution of relative transcript abundances. Thus during contig evaluation each read is given an
362 all-or-nothing assignment with assignments sampled proportional to the estimated abundances.

363 **Simulation of chimeric transcripts**

364 The complete set of transcripts (n = 5917) for the *Saccharomyces cerevisiae* genome were
365 downloaded from <http://www.yeastgenome.org/>. The transcripts were quantified and mRNA

366 abundances recorded using Salmon and the same set of reads used in the *de novo* assembly
367 evaluation described in the section entitled “Analysis of assemblies generated from real reads”. To
368 simulate transcript chimeras, 1000 transcripts were selected at random without replacement from the
369 complete set of transcripts. Pairs of transcripts ($n = 500$) were fused *in silico* by concatenation of two
370 of the randomly selected full length transcript sequences head-to-tail. These 500 transcript chimeras
371 were placed back into the reference transcriptome file (replacing both of their constituent transcripts)
372 such that the transcriptome submitted to TransRate contained the 500 chimeric transcripts and the
373 4917 transcripts that were not chimeric ($n = 5417$). The transcriptome was subject to assessment with
374 TransRate using the same set of RNA-seq reads. This processes was repeated 20 times to obtain the
375 results for the analysis of 10,000 chimeras. The $s(C_{seg})$ score for each transcript chimera was
376 compared to the difference in the relative abundance of the constituent transcripts in the chimera.

377 **TransRate contig scores**

378 TransRate outputs scores for every contig. Here, an assembly consists of a set of contigs C derived
379 from a set of reads \hat{R} . Reads are aligned and assigned to contigs such that R_i is the set of reads
380 assigned to C_i . We propose that a correctly assembled contig derived from a *de novo* transcriptome
381 assembly will have the following four intuitive properties.

- 382 1. **The identity of the nucleotides in the contig will accurately represent the nucleotides of**
383 **the true transcript** $s(C_{nuc})$. This score measures the extent to which the nucleotides in the
384 mapped reads are the same as those in the assembled contig. If the mapped reads do not
385 support the nucleotides of the contig then this likely because: A) The non-supportive reads should
386 map to a different contig or a contig that is not represented in the assembly (a similar gene family
387 variant, alternative allele, or other similarly encoded gene), or B) the assembled sequence is
388 incorrect. In the case of the former, a missing contig (i.e. one that is not assembled) will
389 negatively affect the score of the contig to which its reads incorrectly map. Though the contig to
390 which they map may be correctly assembled, the negative score for this contig can be justified as
391 the incorrectly mapped reads will render the abundance estimate of the assembled contig invalid.

392 In the case of the latter, disagreement between the reads and the contig must be due to mis-
393 assembly. To ensure that stochastic read errors that result in disagreement between a read and a
394 contig do not affect the overall score for that contig support for an alternative nucleotide
395 sequence needs to be provided by multiple reads, (see below).

396 2. **The number of nucleotides in the contig will accurately represent the number in the true**
397 **transcript** $s(C_{cov})$. This score measures the proportion of nucleotides in the contig that have zero
398 coverage and thus have no supporting read data. If there are nucleotides in the contig that are
399 not covered by any reads (regardless of the agreement between the reads and the sequence of
400 the contig) then this should negatively impact on the contig score.

401 3. **The order of the nucleotides in the contig will accurately represent the order in the true**
402 **transcript** $s(C_{ord})$. This score measures the extent to which the order of the bases in contig are
403 correct by analyzing the pairing information in the mapped reads. Here, if the orientation of the
404 mapped reads does not conform to an expected mapping estimated from an analysis of a sub-
405 sample of mapped read pairs then these incorrectly mapping reads will negatively affect the
406 contig score. Similarly, if the contig could have been extended, i.e. there are read-pairs that map
407 such that one read is present near a terminus of the contig and its pair is not mapped and would
408 be expected to map beyond the scope of the contig, then such cases indicate that the contig
409 does not use all of the available reads and thus is incompletely assembled. This metric is
410 informative for the identification partially assembled transcripts.

411 4. **The contig will represent a single transcript** $s(C_{seg})$. This score measures the probability that
412 the coverage depth of the transcript is univariate, i.e. that it represents an assembly of a single
413 transcript and not a hybrid/chimeric assembly of multiple transcripts expressed at different
414 expression levels. Here the per-nucleotide coverage depth of the contig must be best modelled
415 by a single Dirichlet distribution (described below). If the contig is better modelled by the product
416 of two or more Dirichlet distributions then this indicates that two or more contigs with different
417 transcript abundances have been erroneously assembled together.

418 The TransRate contig score is the product of the scores for each of these properties using the aligned
419 reads as evidence. These four properties is evaluated as follows.

420 **Calculation of $s(C_{nuc})$**

421 The alignment edit distance is used to quantify the extent to which the contig sequence is correct. The
422 alignment edit distance is the number of changes that must be made to the sequence of a read in
423 order for it to perfectly match the contig sequence. Here the edit distance of an aligned read $r_{ij} \in R_i$ is
424 denoted as $e_{r_{ij}}$ and the set of reads that cover nucleotide k ($k \in [1, n]$) as ρ_k . The maximum possible
425 edit distance for an alignment is limited by the read alignment algorithm (described in the Read
426 alignment section above) and is denoted as e . The support for the contig provided by the reads is then
427 evaluated as $1 - \frac{e_{r_{ij}}}{e}$ for each $r_i \in \rho_k$, and the mean of all support values is used to calculate $s(C_{nuc})$.

428 **Calculation of $s(C_{cov})$**

429 This score is evaluated as the fraction of nucleotides in the contig that receive at least one mapped
430 read irrespective of the agreement between the read and the contig.

431 **Calculation of $s(C_{ord})$**

432 The pairing information of the mapped reads is used to evaluate this score. To determine the
433 parameters of the read library preparation a randomly selected sub-sample of 1% of all mapped read
434 pairs are analyzed. From these alignments the orientation of the paired end reads is determined and
435 the mean and standard deviation of the fragment size is inferred. All read pair alignments are then
436 classified according to whether they are plausible given the estimated parameters of the library
437 preparation and assuming that the assembled contig is correct. A read pair is considered correct if the
438 following criteria are met: (a) both reads in the pair align to the same contig, (b) the relative orientation
439 of the reads in the pair is consistent with the inferred library preparation parameters, (c) the relative
440 position of the reads is consistent with the mean and standard deviation of the inferred fragment size.
441 $s(C_{ord})$ is then evaluated as the proportion of all mapped read pairs that are correct.

442 **Calculation of $s(C_{seg})$**

443 The per-nucleotide read coverage data is used to evaluate this score. To evaluate the probability that
444 the contig originates from a single transcript (i.e. it is not chimeric) a Bayesian segmentation analysis

445 of the per-nucleotide coverage depth is performed. For a correctly assembled contig it is assumed that
 446 the distribution of per-nucleotide coverage values in that contig is best described by a single Dirichlet
 447 distribution. i.e. all nucleotides in the same transcript should have the same expression level and thus
 448 should be best modelled as a stochastic sample from a single distribution. In contrast, a contig that is
 449 a chimera derived from concatenation of two or more transcripts will have per-nucleotide coverage
 450 values that are best described by two or more different Dirichlet distributions. The probability that the
 451 distribution of per-nucleotide read coverage values comes from a single Dirichlet distribution is
 452 evaluated using a Bayesian segmentation algorithm previously developed for analysis of changes in
 453 nucleotide composition (Liu and Lawrence, 1999). To facilitate the use of this method, the per-
 454 nucleotide coverage along the contig is encoded as a sequence of symbols in an unordered alphabet
 455 by taking \log_2 of the read depth rounded to the nearest integer. As the probability will be a value
 456 between 0 and 1, this probability is used directly as $s(C_{seg})$.

457 **TransRate assembly score**

458 The aim of the TransRate assembly score is to provide insight into the accuracy and completeness of
 459 any given assembly. Thus the assembly score weights equally a summary statistic of the TransRate
 460 contig scores and the proportion of the input reads that are contained within this assembly. We note
 461 here that alternative methods for summarizing contig scores that weight contig scores by their
 462 expression level would produce different results. However, such schemes would not be consistent with
 463 the problem definition and aim of TransRate: to assess the accuracy and completeness of a *de novo*
 464 assembled transcriptome using only the input reads. This score assumes that an ideal assembly will
 465 contain a set of contigs that represent unique and complete transcripts to which all of the reads used
 466 to assemble those transcripts can be mapped. The TransRate assembly score (T) is evaluated as the
 467 geometric mean of the mean contig score and the proportion of read pairs that map to the assembly
 468 such that

$$469 \quad [1] \quad T = \sqrt{(\prod_{c=1}^n s(C))^{\frac{1}{n}} R_{valid}}$$

470 Where

471
$$[2] s(C) = s(C_{nuc})s(C_{cov})s(C_{ord})s(C_{seg})$$

472 **The abundance-weighted TransRate score**

473 An abundance-weighted contig and assembly score are also provided by TransRate. The contig score
474 is evaluated as

475
$$[3] s_w(C) = s(C)^{1+\log_n(TPM+1)}$$

476 where $s(C)$ is as defined in equation [2] and TPM is the transcripts per million transcripts value
477 assigned to that contig by Salmon. Under this framework highly abundant transcripts that have
478 assembly errors are penalized more heavily than low abundance transcripts with the same errors. The
479 abundance-weighted assembly score (T_w) is thus evaluated as

480
$$[4] T_w = \sqrt{(\prod_{c=1}^n s_w(C))^{\frac{1}{n}} R_{valid}}$$

481 **Analysis of assemblies generated from real reads**

482 To demonstrate the utility TransRate contig and assembly scores using real data, TransRate was
483 applied to publicly available benchmark assemblies from two previous analyses (Davidson and
484 Oshlack, 2014; Xie et al., 2014). One set comprised different assemblies generated for rice (*Oryza*
485 *sativa*) and mouse (*Mus musculus*) using the Oases, Trinity, and SOAPdenovo-Trans assemblers (Xie
486 et al., 2014). The other set comprised assemblies for human (*Homo sapiens*) and yeast
487 (*Saccharomyces cerevisiae*) that had been assembled with Oases and Trinity (Davidson and Oshlack,
488 2014). These assemblies were chosen as they have previously been independently used in
489 benchmark comparisons and each of the species has a completed annotated reference genome
490 available. In all cases, TransRate was run with the published reads and the published assembly as
491 input.

492 **Independence of score components**

493 Correlation between the contig score components was measured for the assemblies from real data.

494 To prevent larger assemblies from biasing the results, 5,000 contigs were sampled at random from
495 each assembly. These contigs were used to calculate a Spearman's rank correlation coefficient using

496 R version 3.1.1 (R Core Team, 2014). The correlation between any two score components was taken
 497 as the mean of the correlation across all datasets.

498 **Identification of reconstructed reference transcripts**

499 The full set of coding and non-coding transcripts for each species were downloaded from Ensembl
 500 Genomes version 25 (<ftp://ftp.ensemblgenomes.org/pub/release-25/>). Assembled contigs were then
 501 identified by BLAST searching the reference dataset for the corresponding species using bidirectional
 502 blastn local alignment with an e-value cutoff of 10^{-5} (BLAST+ version 2.2.29 (Camacho et al., 2009)).
 503 Only reciprocal best hits were retained for further analysis.

504 **Assembly from simulated read data**

505 For each species, a total of 10 million mRNA molecules were simulated from the full set of annotated
 506 mRNAs from the Ensembl reference with exponentially distributed expression values using the flux
 507 simulator v1.2.1 (Griebel et al., 2012). mRNA molecules were uniform-randomly fragmented and then
 508 size-selected to a mean of 400 nucleotides and standard deviation of 50 nucleotides. From the
 509 resulting fragments, 4 million pairs of 100bp reads were simulated using the default error profile
 510 included in flux-simulator. An assembly was generated from these simulated reads using
 511 SOAPdenovo-Trans with default parameters.

512 **Calculation of contig accuracy**

513 Accuracy was calculated by comparing contigs assembled from simulated data to the set of transcripts
 514 from which the read data were simulated. Reciprocal best BLAST hits were identified and the
 515 accuracy of each contig assembled from simulated read data was evaluated as the contig F-score
 516 where

$$517 \quad [5] \text{ Contig precision} = \frac{\text{Number of correct nucleotides in contig}}{\text{Number of nucleotides in contig}}$$

$$518 \quad [6] \text{ Contig recall} = \frac{\text{Number of correct nucleotides in contig}}{\text{Number of nucleotides in reference transcript}}$$

$$519 \quad [7] \text{ Contig } F - \text{score} = 2 \left(\frac{(\text{contig precision})(\text{contig recall})}{(\text{contig precision} + \text{contig recall})} \right)$$

520 Spearman's rank correlation coefficient between the contig F-score and the TransRate contig score
521 was calculated using R version 3.1.1. The same contigs were also subject to analysis using RSEM-
522 eval and the relationship between contig impact score and contig F-score analyzed using the same
523 method.

524 **Constructing a benchmark dataset of TransRate scores**

525 A survey of the range of assembly scores for published *de novo* transcriptome assemblies was
526 conducted by analyzing a sub-set of transcriptome assemblies from the Transcriptome Shotgun
527 Archive (<http://www.ncbi.nlm.nih.gov/genbank/tsa>). *De novo* assembled transcriptomes were used in
528 this analysis only if paired-end reads were provided, the assembler and species were named in the
529 metadata, and the assembly contained at least 5,000 contigs (TransRate has no minimum or
530 maximum contig requirements but a minimum number of 5,000 was imposed to ensure sufficient raw
531 data was available for analysis). For each of these test datasets, the assembly and reads were
532 downloaded. TransRate was run on all assemblies and FastQC version 2.3 (Andrews, 2010) was used
533 to evaluate the quality of the read datasets.

534 **Software availability**

535 TransRate is written in Ruby and C++, and makes use of the BioRuby (Goto et al., 2010) and
536 BAMtools (Barnett et al., 2011) libraries. The source code is available in a compressed archive in the
537 supplemental material (Supplemental File S1) and at <http://github.com/Blahah/transrate> and is
538 released under the open source MIT license. Binary downloads and full documentation are available
539 at <http://hibberdlab.com/transrate>. The software is operated via a command line interface and can be
540 used on OSX and Linux. TransRate can also be used programmatically as a Ruby gem.

541 **Disclosure declaration**

542 The authors declare no competing interests.

543 **Acknowledgements**

544 The authors thank the TransRate user community for testing, bug reports and feedback. In particular,
545 we thank Matt MacManes for meticulous testing. RSU was funded by The Millennium Seed Bank
546 (Royal Botanical Gardens, Kew). SK is a Royal Society University Research Fellow. CB was

547 supported by a grant to JMH from the Biotechnology and Biological Sciences Research Council, the
548 Department for International Development and (through a grant to BBSRC) the Bill & Melinda Gates
549 Foundation, under the Sustainable Crop Production Research for International Development
550 programme, a joint initiative with the Department of Biotechnology of the Government of India's
551 Ministry of Science and Technology. This project has received funding from the European Union's
552 Horizon 2020 research and innovation program under grant agreement no 637765.

553

554 **Tables**555 **Table 1**

Score component	Description
$s(C_{nuc})$	The proportion of nucleotides in the mapped reads that are the same as those in the assembled contig
$s(C_{cov})$	The proportion of nucleotides in the contig that have have no supporting read data
$s(C_{ord})$	The extent to which the order of the bases in contig are correct by analyzing the pairing information in the mapped reads
$s(C_{seg})$	The probability that the coverage depth of the transcript is univariate

556

557 **Figure legends**558 **Figure 1 Common errors in *de-novo* transcriptome assembly, and how they can be detected**559 **using read mapping data.** *Family collapse* occurs when multiple members of a gene family are

560 assembled into a single hybrid contig. This error can be detected by measuring the extent that the

561 nucleotides in the contig are supported by the mapped reads. *Chimerism* occurs when two or more

562 transcripts (that may or may not be related) are concatenated together in a single contig during

563 assembly. This can be detected when the expression levels of the transcripts differ, leading to a

564 change-point in the read coverage along the contig. *Unsupported insertions* can be detected as bases565 in a contig that are unsupported by the read evidence. *Incompleteness* can be detected when reads or566 fragments align off the end of the contig. *Fragmentation* is caused by low coverage and is detectable567 when read pairs bridge two contigs. *Local misassembly* encompasses various structural errors that

568 can occur during assembly, such as inversions, usually as a result of assembler heuristics. These are

569 detectable when both members of a read pairs align to single contig, but in manner inconsistent with

570 the sequencing protocol. *Redundancy* occurs when a single transcript is represented by multiple

571 overlapping contigs in an assembly. This is detectable when reads align to multiple contigs but the
 572 assignment process assigns them all to the contig that best represents the original transcript.

573 **Figure 2. The TransRate workflow.** (1) TransRate takes as input one or more *de novo* transcriptome
 574 assemblies and the paired-end reads used to generate them. (2) The reads are aligned to the contigs.
 575 (3) Multi-mapping reads are proportionally assigned to contigs based on the posterior probability that
 576 each contig was the true origin of the read. (4) The alignments are evaluated using four discrete score
 577 components. (5) The four score components are integrated to generate the TransRate contig score.
 578 (6) The TransRate assembly score is calculated from analysis of all contig scores.

579 **Figure 3. Distribution and interrelationship of contig score components.** (A) Distribution of contig
 580 score components in ten different assemblies spanning four species and three different assemblers.
 581 $s(C_{nuc})$ is the fraction of nucleotides in a contig whose sequence identity agrees with the aligned
 582 reads. $s(C_{cov})$ is the fraction of nucleotides in a contig that have one or more mapped reads. $s(C_{ord})$ is
 583 the fraction of reads that map to the contig in the correct orientation. $s(C_{seg})$ is the probability that the
 584 read coverage along the length of the contig is best explained by a single Dirichlet distribution, as
 585 opposed to two or more distributions. (B) The Spearman's rank correlation coefficient between the
 586 contig score components, averaged across all species and assemblers.

587 **Figure 4. TransRate contig score is related to assembly accuracy.** Contigs from assemblies of
 588 simulated reads from four species (rice, mouse, yeast, and human) were evaluated using TransRate
 589 and RSEM-eval. Reciprocal best-BLAST against the true set of transcripts was used to determine the
 590 F-score, or reference-based accuracy, of the assembled contig. Each point is a contig in an assembly,
 591 with all four assemblies on the same plot. A) Comparison between TransRate contig score and contig
 592 F-score. B) Comparison between RSEM-eval contig impact score and contig F-score, with contig
 593 impact scores below 0 set to the smallest positive value in the data to enable plotting.

594 **Figure 5. Calculation of TransRate assembly scores.** A) Distribution of TransRate contig scores for
 595 the 10 representative assemblies from real data. B) Geometric mean of TransRate contig scores for all
 596 assemblies. C) Proportion of reads that map to each assembly. D) Final TransRate assembly scores
 597 for the 10 representative assemblies. E) The proportion of reference transcripts that are best

598 assembled by individual assembly methods. F) The number of reference transcripts (identified by
599 reciprocal best BLAST) that are assembled by each assembler.

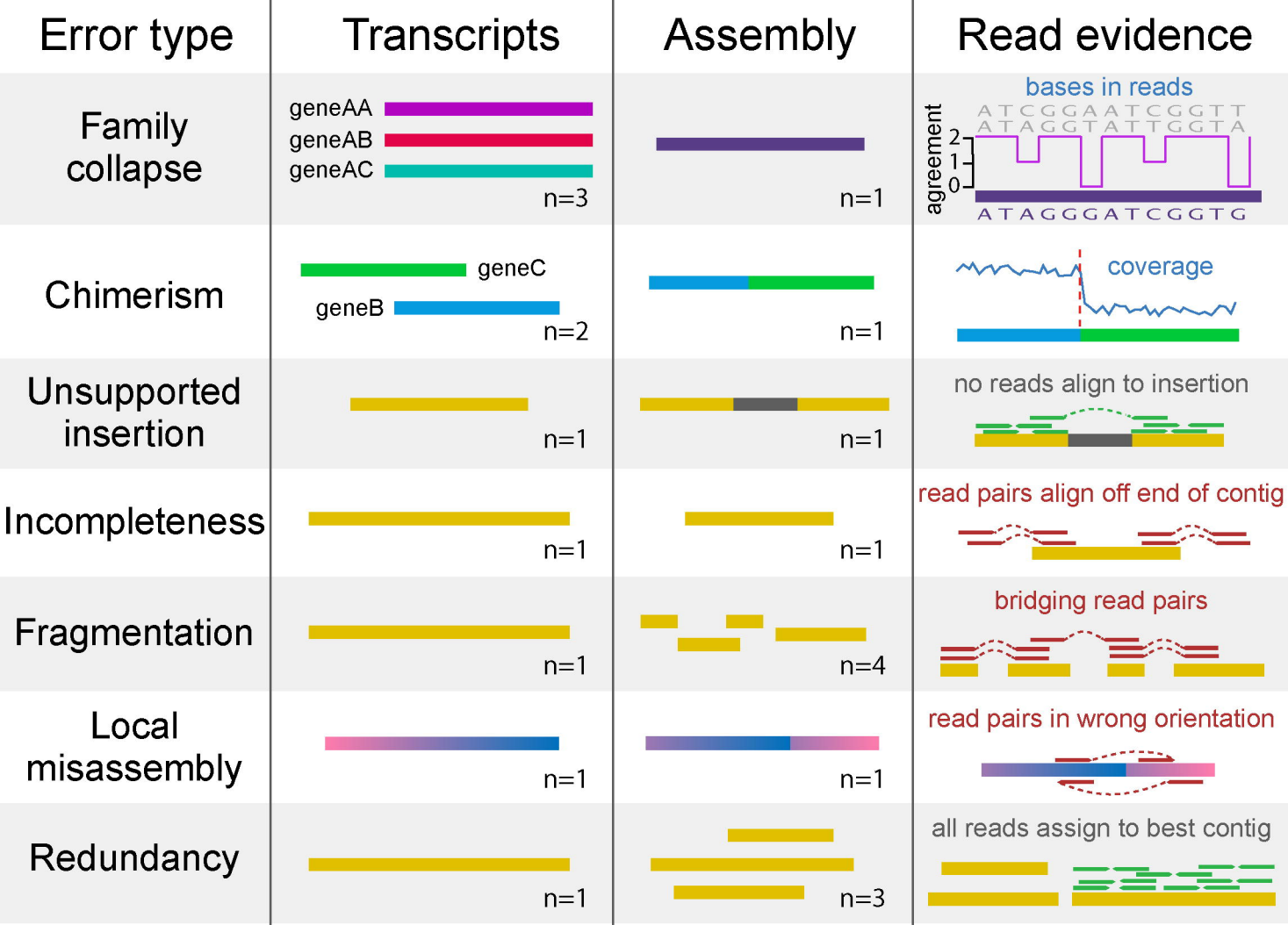
600 **Figure 6. Application of TransRate to 155 published assemblies from the NCBI Transcriptome**

601 **Shotgun Archive.** 155 assemblies from the Transcriptome Shotgun Archive were analysed using
602 TransRate. The quality of the reads used to generate the assemblies were also analysed using
603 FastQC. A) Cumulative distribution of TransRate raw and optimised assembly scores for each of the
604 155 assemblies. B) Comparison between raw and optimised assembly score. C) Distribution of
605 TransRate optimised assembly scores partitioned by taxonomic group. C) Distribution of TransRate
606 optimised assembly scores partitioned by assembly method. (E-J) TransRate optimised assembly
607 scores compared to various summary statistics of the input reads: E) read length, F) read GC%, G)
608 mean read per-base Phred score, H) percent of reads that were PCR duplicates, I) number of read
609 pairs, and J) read bases per assembled base.

610 **References**

- 611 Andrews, S., 2010. FastQC: a quality control tool for high throughput sequence data [WWW
612 Document]. URL <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- 613 Barnett, D.W., Garrison, E.K., Quinlan, A.R., Strömberg, M.P., Marth, G.T., 2011. BAMTools: a C++
614 API and toolkit for analyzing and managing BAM files. *Bioinformatics* 27, 1691–1692.
615 doi:10.1093/bioinformatics/btr174
- 616 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., Madden, T.L., 2009.
617 BLAST+: architecture and applications. *BMC Bioinformatics* 10, 421. doi:10.1186/1471-2105-
618 10-421
- 619 Clark, S.C., Egan, R., Frazier, P.I., Wang, Z., 2013. ALE: a generic assembly likelihood evaluation
620 framework for assessing the accuracy of genome and metagenome assemblies. *Bioinformatics*
621 29, 435–443. doi:10.1093/bioinformatics/bts723
- 622 Davidson, N.M., Oshlack, A., 2014. Corset: enabling differential gene expression analysis for de novo
623 assembled transcriptomes. *Genome Biol.* 15, 410. doi:10.1186/s13059-014-0410-6
- 624 Goto, N., Prins, P., Nakao, M., Bonnal, R., Aerts, J., Katayama, T., 2010. BioRuby: bioinformatics
625 software for the Ruby programming language. *Bioinformatics* 26, 2617–2619.
626 doi:10.1093/bioinformatics/btq475
- 627 Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L.,
628 Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, A., Rhind, N., Palma,
629 F. di, Birren, B.W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., Regev, A., 2011. Full-length
630 transcriptome assembly from RNA-Seq data without a reference genome. *Nat. Biotechnol.* 29,
631 644–652. doi:10.1038/nbt.1883
- 632 Griebel, T., Zacher, B., Ribeca, P., Raineri, E., Lacroix, V., Guigó, R., Sammeth, M., 2012. Modelling
633 and simulating generic RNA-Seq experiments with the flux simulator. *Nucleic Acids Res.* 40,
634 10073–10083. doi:10.1093/nar/gks666
- 635 Gurevich, A., Saveliev, V., Vyahhi, N., Tesler, G., 2013. QUASt: quality assessment tool for genome
636 assemblies. *Bioinformatics* 29, 1072–1075. doi:10.1093/bioinformatics/btt086

- 637 Li, B., Fillmore, N., Bai, Y., Collins, M., Thomson, J.A., Stewart, R., Dewey, C.N., 2014. Evaluation of
638 de novo transcriptome assemblies from RNA-Seq data. *Genome Biol.* 15, 553.
639 doi:10.1186/s13059-014-0553-5
- 640 Li, B., Ruotti, V., Stewart, R.M., Thomson, J.A., Dewey, C.N., 2010. RNA-Seq gene expression
641 estimation with read mapping uncertainty. *Bioinformatics* 26, 493–500.
642 doi:10.1093/bioinformatics/btp692
- 643 Liu, J.S., Lawrence, C.E., 1999. Bayesian inference on biopolymer models. *Bioinformatics* 15, 38–52.
644 doi:10.1093/bioinformatics/15.1.38
- 645 Lowe, E.K., Swalla, B.J., Brown, C.T., 2014. Evaluating a lightweight transcriptome assembly pipeline
646 on two closely related ascidian species. *PeerJ Prepr.* 505. doi:10.7287/peerj.preprints.505v1
- 647 Martin, J.A., Wang, Z., 2011. Next-generation transcriptome assembly. *Nat. Rev. Genet.* 12, 671–682.
648 doi:10.1038/nrg3068
- 649 Myers, E.W., 2005. The fragment assembly string graph. *Bioinformatics* 21, ii79–ii85.
650 doi:10.1093/bioinformatics/bti1114
- 651 O’Neil, S.T., Emrich, S.J., 2013. Assessing De Novo transcriptome assembly metrics for consistency
652 and utility. *BMC Genomics* 14, 465. doi:10.1186/1471-2164-14-465
- 653 Patro, R., Duggal, G., Kingsford, C., 2015. Accurate, fast, and model-aware transcript expression
654 quantification with Salmon. *bioRxiv* 021592. doi:10.1101/021592
- 655 Peng, Y., Leung, H.C.M., Yiu, S.-M., Lv, M.-J., Zhu, X.-G., Chin, F.Y.L., 2013. IDBA-tran: a more
656 robust de novo de Bruijn graph assembler for transcriptomes with uneven expression levels.
657 *Bioinformatics* 29, i326–i334. doi:10.1093/bioinformatics/btt219
- 658 Rahman, A., Pachter, L., 2013. CGAL: computing genome assembly likelihoods. *Genome Biol.* 14,
659 R8. doi:10.1186/gb-2013-14-1-r8
- 660 R Core Team, 2014. R: A language and environment for statistical computing. R Foundation for
661 Statistical Computing, Vienna, Austria.
- 662 Robertson, G., Schein, J., Chiu, R., Corbett, R., Field, M., Jackman, S.D., Mungall, K., Lee, S., Okada,
663 H.M., Qian, J.Q., Griffith, M., Raymond, A., Thiessen, N., Cezard, T., Butterfield, Y.S.,
664 Newsome, R., Chan, S.K., She, R., Varhol, R., Kamoh, B., Prabhu, A.-L., Tam, A., Zhao, Y.,
665 Moore, R.A., Hirst, M., Marra, M.A., Jones, S.J.M., Hoodless, P.A., Birol, I., 2010. De novo
666 assembly and analysis of RNA-seq data. *Nat. Methods* 7, 909–912. doi:10.1038/nmeth.1517
- 667 Schulz, M.H., Zerbino, D.R., Vingron, M., Birney, E., 2012. Oases: robust de novo RNA-seq assembly
668 across the dynamic range of expression levels. *Bioinformatics* 28, 1086–1092.
669 doi:10.1093/bioinformatics/bts094
- 670 Xie, Y., Wu, G., Tang, J., Luo, R., Patterson, J., Liu, S., Huang, W., He, G., Gu, S., Li, S., Zhou, X.,
671 Lam, T.-W., Li, Y., Xu, X., Wong, G.K.-S., Wang, J., 2014. SOAPdenovo-Trans: de novo
672 transcriptome assembly with short RNA-Seq reads. *Bioinformatics* btu077.
673 doi:10.1093/bioinformatics/btu077
- 674 Zaharia, M., Bolosky, W.J., Curtis, K., Fox, A., Patterson, D., Shenker, S., Stoica, I., Karp, R.M.,
675 Sittler, T., 2011. Faster and More Accurate Sequence Alignment with SNAP. *ArXiv11115572*
676 *Cs Q-Bio.*
677





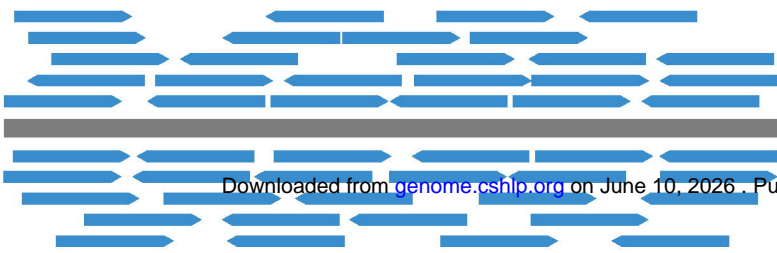
TransRate

1 input data

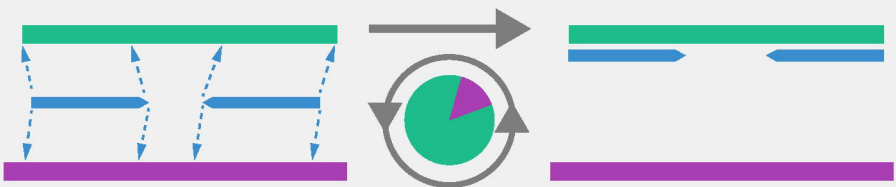
assembled contigs paired-end reads



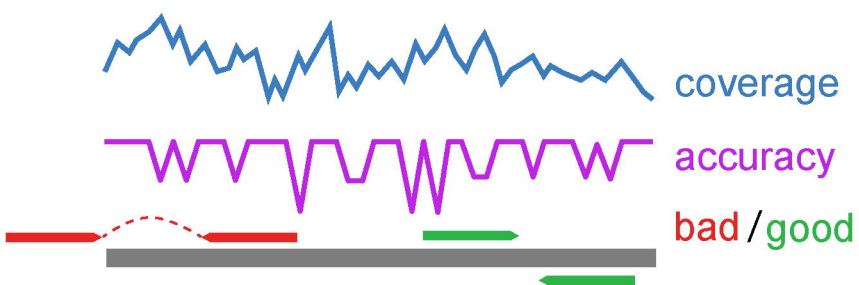
2 align reads to contigs



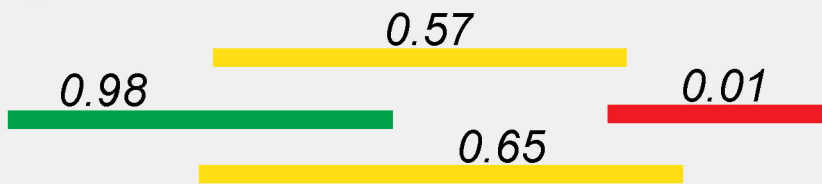
3 assign multimapping reads



4 collect contig score components



5 calculate contig scores



6 calculate assembly score

