# Effects of mapping algorithms on gene selection for RNA-Seq analysis: pulmonary response to acute neonatal hyperoxia

Chin-Yi Chu, Soumyaroop Bhattacharya, Zhongyang Zhou, Min Yee, Ashley Lopez, Valerie A Lunger, Bradley Buczynski, Michael Oreilly, Thomas Mariani

**Background:** A major goal of RNA-Seq data analysis is to reconstruct the full set of gene transcripts expressed in a biological sample in order to quantify their expression levels. The process typically involves multiple steps including mapping short sequence reads to a reference genome, and estimating expression levels based on these mappings. Multiple algorithms and approaches for each processing step exist, and the impact of different methods on estimation of gene expression is not entirely clear.

**Methods:** We evaluated the impact of three common mapping algorithms on differential expression analysis in an RNA-Seq dataset describing the lung response to acute neonatal hyperoxia. RNA-Seq data generated using the Illumina platform were mapped and aligned using CASAVA, TopHat, and SHRiMP against the mouse genome. Significance Analysis of Microarrays and Cuffdiff were used to identify differentially expressed genes between hyperoxia-challenged and age matched control mice.

**Results:** 1403 genes were detected as differentially expressed by least one mapping and gene selection method. A majority of genes (>65%) were identified by all three mapping methods, regardless of the gene selection approach. Expression patterns for 52 genes were examined by quantitative polymerase chain reaction (qPCR). Importantly, we found different validation rates for genes selected by each method; 72% for CASAVA, 69% for TopHat and 63% for SHRiMP. Surprisingly, the validation rate for genes selected by all three mapping methods was no greater than the best single method.

**Conclusion**: The choice of mapping strategy impacts the reliability of gene selection for RNA-Seq data analysis.

2  Effects of Mapping Algorithms on Gene Selection for RNA-Seq Analysis: Pulmonary Response to
3  Acute Neonatal Hyperoxia.
4
5  Chin-Yi Chu[1,2,a], Soumyaroop Bhattacharya[1,2,a], Zhongyang Zhou[1], Min Yee[1], Ashley M Lopez[1],
6  Valerie A Lunger[1], Bradley W Buczynski[1], Michael A O'Reilly[1,3], Thomas J Mariani[1,2]
7
8  [a]These authors contributed equally.
9
10 [1]Division of Neonatology, [2]Pediatric Molecular and Personalized Medicine (PMPM) Program, and
11 [3]Perinatal and Pediatric Origins of Disease (PPOD) Program, Department of Pediatrics, University of
12 Rochester Medical Center, Rochester NY, United States of America
13
14 3 Figures
15 5 Tables
16 4 Supplemental Figures
17 1 Supplemental Table
18
19 Author Emails:
20 Chin-Yi Chu: chinyi_chu@urmc.rochester.edu
21 Soumyaroop Bhattacharya: soumyaroop_bhattacharya@urmc.rochester.edu
22 Zhongyang Zhou: zzhou10@ur.rochester.edu
23 Min Yee: min_yee@urmc.rochester.edu
24 Ashley M Lopez: ashley_lopez@urmc.rochester.edu
25 Valerie A Lunger: valerie_lunger@urmc.rochester.edu
26 Bradley W Buczynski: bradley.buczynski@wilresearch.com
27 Michael A O'Reilly: michael_oreilly@urmc.rochester.edu
28 Thomas J Mariani: tom_mariani@urmc.rochester.edu
29
30 Address for Correspondence:
31
32 Thomas J Mariani, PhD
33 Division of Neonatology and
34 Center for Pediatric Biomedical Research
35 University of Rochester Medical Center,
36 601 Elmwood Ave, Box 850,
37 Rochester, NY 14642, USA
38
39 Phone: 585-276-4616
40 Fax: 585-276-2643
41 Email: Tom_Mariani@URMC.rochester.edu
42
43 Running Title: Comparative mapping of RNA-Seq
44 Keywords: TopHat, SHRiMP, CASAVA, qPCR

46  ABSTRACT

47  **Background:** A major goal of RNA-Seq data analysis is to reconstruct the full set of gene transcripts
48  expressed in a biological sample in order to quantify their expression levels. The process typically
49  involves multiple steps including mapping short sequence reads to a reference genome, and estimating
50  expression levels based on these mappings. Multiple algorithms and approaches for each processing
51  step exist, and the impact of different methods on estimation of gene expression is not entirely clear.

52

53  **Methods:** We evaluated the impact of three common mapping algorithms on differential expression
54  analysis in an RNA-Seq dataset describing the lung response to acute neonatal hyperoxia. RNA-Seq
55  data generated using the Illumina platform were mapped and aligned using CASAVA, TopHat, and
56  SHRiMP against the mouse genome. Significance Analysis of Microarrays and Cuffdiff were used to
57  identify differentially expressed genes between hyperoxia-challenged and age matched control mice.

58

59  **Results:** 1403 genes were detected as differentially expressed by least one mapping and gene selection
60  method. A majority of genes (>65%) were identified by all three mapping methods, regardless of the
61  gene selection approach. Expression patterns for 52 genes were examined by quantitative polymerase
62  chain reaction (qPCR). Importantly, we found different validation rates for genes selected by each
63  method; 72% for CASAVA, 69% for TopHat and 63% for SHRiMP. Surprisingly, the validation rate
64  for genes selected by all three mapping methods was no greater than the best single method.

65

66  **Conclusion**: The choice of mapping strategy impacts the reliability of gene selection for RNA-Seq
67  data analysis.

68   INTRODUCTION
69   Genome-wide expression profiling is used to assess the expression of thousands of genes
70   simultaneously, with the ultimate goal of creating a comprehensive description of expression at the
71   mRNA level. Historically DNA microarray technology has been used to measure genome-wide
72   expression. With the advent of next-generation sequencing, high-throughput sequence-based
73   approaches for expression analysis are becoming an increasingly popular alternative to microarrays.
74   RNA-Seq, otherwise known as Whole Transcriptome Shotgun Sequencing (WTSS), refers to the use of
75   high-throughput technologies to sequence cDNA in order to get information about a sample's RNA
76   content. Applications of RNA-Seq may include analyzing the genome for coding and non-coding RNA
77   [1]. RNA-Seq has been used in studies of transcription in yeast [2], Arabidopsis [3], mouse [4, 5], and
78   human [6, 7]. RNA-Seq is a promising replacement for microarrays as initial studies have shown that
79   RNA-Seq expression estimates are highly reproducible[6] and are often more accurate, based on
80   assessments by either quantitative PCR (qPCR) or spike-in experiments [2, 5]. The primary advantages
81   of RNA-Seq are its large dynamic range (spanning five orders of magnitude), low background noise,
82   reduced input sample (RNA) requirement and ability to detect novel transcripts, when studied at
83   appropriate coverage depth [8]. However, similar to the early days of microarray analysis, there are still
84   a number of experimental and computational issues remaining to be resolved for RNA-Seq.
85
86   The process of RNA-Seq involves four major steps: (i) generation of cDNA libraries from the RNA
87   samples, (ii) generating millions of short sequence "reads" from these cDNA libraries (that are
88   proportional to the mRNA diversity of the initial RNA sample), (iii) mapping the sequence reads to a
89   reference genome, and (iv) summarizing the mapped reads into raw expression level measures, for
90   subsequent normalization and analysis. Methods and computational tools are being rapidly developed
91   to meet the challenges of these steps. Presently, more than 60 mappers are available to accomplish this
92   task [9, 10]. Our choices of mappers represent a spliced (TopHat) and a contemporary unspliced
93   (SHRiMP) freeware aligner, as well as the vendor recommended default aligner. While both CASAVA
94   and TopHat are spliced read aligner which is a critical feature when aligning RNA-seq reads to a
95   reference genome, we chose to compare them with SHRiMP, which is an unspliced read aligner[11].
96   Consensus Assessment of Sequence And Variation, (CASAVA) is the part of Illumina's sequencing
97   analysis software that performs alignment of a sequencing run to a reference genome and subsequent
98   variant analysis and read counting [12]. Its underlying alignment algorithm is ELAND (Efficient
99   Large-Scale Alignment of Nucleotide Databases). ELAND is very fast and can perform multi-seed and
100  gapped alignment. TopHat is a fast and popular spliced aligner for RNA-Seq reads [13]. It aligns RNA-
101  Seq reads to mammalian-sized genomes by using its underlying mapping engine, either Bowtie or
102  Bowtie2. TopHat consists of three mapping steps: transcriptome mapping, genome mapping, and
103  spliced mapping. This three-step pipeline assures the best and unique alignment for each read. It can
104  identify novel splice sites with direct mapping to known transcripts, producing sensitive and accurate
105  alignments, even for highly repetitive genomes or in the presence of pseudogenes. Short Read Mapping
106  Program (SHRiMP) is an unspliced aligner and it aligns reads against a target genome using an
107  algorithm distinct from that TopHat uses [14]. It was initially developed with the multitudinous short
108  reads of next generation sequencers and Applied Biosystem's colorspace genomic representation. The
109  algorithm can align reads with extensive polymorphism and sequencing errors. It can also perform the
110  multi-seed alignment to speed up the alignment.
111
112  Here we studied the biological impact of these three mapping algorithms on identification of
113  differentially expressed genes in a typical experimental model dataset. We included molecular
114  validation for genes identified in individually mapped datasets as a means of determining the estimated

115  reliability of each method. Similar studies in the past have used empirical and simulated datasets [15-
116  18], however, these studies did not have further validation results.
117
118

119 METHODS
120 *RNA-Seq Data*
121 This data set represents high throughput sequencing of seven cDNA libraries generated from whole
122 lung tissue RNA recovered from mice treated under hyperoxic conditions in the newborn period (n=3)
123 or age-matched controls (n=4). Mice were exposed to 100% oxygen, 40-70% humidity between birth
124 and postnatal day 10 as shown previously [19]. Mice were housed in sterile microisolator cages in a
125 specified pathogen-free environment and exposed to super-physiological levels of oxygen according to
126 a protocol (Protocol No. 2007-121R) approved by the University Committee on Animal Resources at
127 the University of Rochester. The University Committee on Animal Resources (UCAR) at the
128 University of Rochester reviewed and approved these studies. The data were generated on the Illumina
129 Genome Analyzer II platform at an average sequence depth of 20 million reads (65 bases in length) per
130 sample as reported previously [20].
131
132 *Short-Read Alignment*
133 Sequences were mapped to the mouse genome (mm10), containing 21,718 unique genes, using
134 multiple alignment algorithms. CASAVA (v 1.7) pipeline software was implemented with the
135 manufacturer default settings, using the ELAND aligner program in Multiplexer setting, for separating
136 bar-coded reads into different bins. Tophat (v 2.0.9) was implemented using Bowtie (v 2.1.0) and
137 SAMtools (v 0.1.18) with default parameters (b2-sensitive, report-secondary-alignments, library-type
138 fr-unstranded). SHRiMP (v 2.2.3) was implemented using SAMtools (v 0.1.18) with default
139 parameters (gmapper – ls –qv-offset 33 single_en_reads-). For TopHat and SHRiMP mapping, HTSeq
140 (v 0.5.3p3) was used to generate the count matrices with the following parameters: 'htseq-count -m
141 intersection-strict -s no'. For both SHRiMP and TopHat alignments, mouse gene annotation file (GTF)
142 for mouse genome build 10 obtained from UCSC was used for alignment. Codes for mapping and
143 generating the raw counts are available in shown in Supplemental Table 1.
144
145 *Normalization*
146 Alignment counts obtained from CASAVA, SHRiMP and TopHat were further normalized using reads
147 per million bases (RPM) or trimmed-mean (TM), independently. RPM normalization involves dividing
148 the raw count for each individual gene in a particular sample by the sum total of counts for all the
149 genes in that sample (Equation 1). It is a modification of the RPKM (reads per kilobase transcript per
150 million reads) approach, excluding normalization for transcript length [5].
151

$$Count_{RPM} = \frac{Count_{Raw}}{\sum_{sampleA} Count} \qquad \text{.....(1)}$$

156 TM normalization involves calculation of the mean of 5th-95th percentile of raw counts for each
157 sample, which is used to determine a sample-specific normalization factor. The raw counts for all
158 genes for a sample are then multiplied by this factor (Equation 2).

$$Count_{TrimMEAN} = Count_{Raw}(\frac{1}{\mu_{Trim(5-95\%)}}) \qquad \text{......(2)}$$

162 *Gene Selection*
163 Genes that were not detected in all samples of at least one experimental group (hyperoxia-treated or
164 control) were excluded from analysis. Multiple gene selection approaches were applied on each
165 normalized and filtered dataset, from each of the three mapping algorithms. The filtered and
166 normalized data, were however not log-transformed. Significance Analysis for Microarrays (SAM) is a

167 frequentist approach used for identification of differentially expressed genes that uses a modified t-
168 statistic with permutations [21]. SAM was applied using the minimum median False Discovery Rate
169 (FDR) possible, in MultiExperiment Viewer (MeV) v4.8.1 (http://www.tm4.org/mev.html). For the
170 purposes of this study, SAM threshold of median FDR of 0 was applied. Cuffdiff uses the Cufflinks
171 transcript quantification to calculate gene expression levels in different conditions, and tests them for
172 significant differences [22]. It uses FPKM (fragments per kilobase of exon per million fragments
173 mapped) normalization. Cuffdiff was applied using mouse gene annotation file (GTF) for mouse
174 genome build 10 obtained from UCSC with a significance threshold of FDR adjusted q <0.05.
175
176 *Molecular Validation*
177 cDNA were synthesised from RNA samples isolated from individual lungs using the iScript Reverse
178 Transcription Kit (BioRad, Hercules, CA). qPCR was performed on a Viia7 (Applied Biosystems,
179 Santa Clara, CA) using SYBR green chemistry as previously described [23]. Gene-specific assays
180 primer sequences were retrieved from the MGH Primer Bank (http://pga.mgh.harvard.edu/primerbank).
181 Gene expression levels (dCt) were calculated relative to the measured Ct value of PPIA (peptidyl
182 prolyl isomerase A or cyclophilin A) as an internal, endogenous control and were analyzed for relative
183 expression changes by the ddCt method as previously described [23]. QPCR data were assessed using
184 both the students T-test and the Mann-Whitney U test at a p<0.05.

185 RESULTS
186 The main objective here is to compare the effect of different, state-of-the-art alignment and mapping
187 procedures upon the accuracy of gene selection. The data set we used for this analysis consists of 3
188 experimental treatment samples (separate pools of RNA isolated from lungs of neonatal mice
189 challenged with hyperoxia) and 4 matched, normoxic controls [20]. As indicated in Figure 1, in order
190 to complete the analysis, Fastq files were independently mapped using 3 different alignment methods;
191 CASAVA, TopHat and SHRiMP. The mapped data were normalized and filtered, using tools
192 appropriate for the mapping methods, and analyzed for differential expression using Significance
193 Analysis for Microarrays (SAM) or Cuffdiff. In particular, we used reads per million (RPM) or
194 trimmed mean (TM) normalization for SAM, and "fragments per kilobase of exon per million
195 fragments mapped" (FPKM) normalization for Cuffdiff. Differential expression was estimated using
196 SAM, for all three mapping methods, or Cuffdiff for TopHat and SHRiMP mapping. qPCR was used to
197 assess the reliability of differential expression estimation.
198
199 *Mapping Statistics*
200 For each mapping algorithm we calculated the total number of reads, percentage of genes detected,
201 alignment rate, number of aligned reads, and the number of reads assigned to genes. Table 1 shows a
202 summary of these mapping statistics. SHRiMP gave highest total number of reads followed by those of
203 CASAVA and TopHat. On an average, SHRiMP gave highest percent genes, detected (80%) which
204 was slightly above those of CASAVA (77%) and TopHat (76%). In addition, the total number of reads
205 assigned to mapped genes was highest in SHRiMP (12.75 x $10^6$), followed by TopHat (12.56 x $10^6$)
206 and CASAVA (11.18 x $10^6$).
207
208 The correlation analysis confirmed strong general concordance on the gene expression measurements
209 across mappers. Pearson correlation coefficients between the raw or either of normalized counts
210 generated by the TopHat and CASAVA was found to be well above 0.8 indicating the data were of
211 comparable quality (Figure 1). While the pearson correlation coefficients between the raw or either of
212 normalized counts generated by the SHRiMP and CASAVA were found to be 0.46 or more, the
213 correlation coefficients between the raw or either of normalized counts generated by the SHRiMP and
214 TopHat were 0.58 or more. Figure 2 shows representative plots (from one control and one hyperoxia
215 sample) for same-sample correlations among the normalized counts obtained using three mappers.
216 Correlation plots for raw and both normalized counts can be seen in Supplemental Figure 1.
217
218 *CASAVA (v 1.7)*
219 On average, expression of 77% of the genes in the genome was detected in the samples using
220 CASAVA. After removing signal from genes not present in all samples of at least one experimental
221 group, the expression of 16,079 genes were assessed for differential expression (Table 2). As described
222 elsewhere [20] SAM identified 1020 genes in the TM-normalized data and 813 genes in the RPM-
223 normalized data. A total of 798 genes were common, and 300 of these had a fold-change greater than
224 or equal to 2. There were a greater number of genes showing significant decreases in expression in
225 response to treatment (56%) than were increased.
226
227 *SHRiMP (v 2.2.3)*
228 On average, expression of 80% of the genes in the genome was detected in the lung tissue samples
229 using SHRiMP. After removing signal from genes not present in all samples of at least one
230 experimental group, the expression of 17,814 genes were assessed for differential expression (Table 3).
231 SAM identified 879 genes in the TM-normalized data and 937 genes in the RPM-normalized data. A

232 total of 857 genes were common, and 386 of these had a fold-change greater than or equal to 2. Similar
233 to CASAVA mapped data, there were more genes showing significant decreases in expression in
234 response to treatment (60%) than were increased. Cuffdiff identified 3886 genes as differentially
235 expressed, and 1087 of those had a magnitude of change greater than or equal to 2. In addition to
236 selecting a greater number of genes as significantly affected, Cuffdiff identified a greater number of
237 genes showing significant increases in expression in response to treatment (57%) than were decreased.
238
239 *TopHat (v 2.0.9)*
240 On average, expression of 76% of the genes in the genome was detected in the lung tissue samples
241 using TopHat. After removing signal from genes not present in all samples of at least one experimental
242 group, the expression of 16,892 genes were assessed for differential expression (Table 4). SAM
243 identified 880 genes in the TM-normalized data and 951 genes in the RPM-normalized data. A total of
244 860 genes were common, and 396 of these had a fold-change greater than or equal to 2. Cuffdiff
245 identified 2831 genes as differentially expressed, and 1044 of those had a fold-change greater than or
246 equal to 2. Again (similar to SHRiMP mapped data) Cuffdiff identified a greater number of genes
247 showing significant increases in expression in response to treatment (57%) than were decreased.
248
249 *Consistency of Gene Selection*
250 We tested the consistency of individual gene selection tools. SAM identified 699 genes that were
251 selected as differentially expressed by all three mapping methods, of which 251 had a fold change
252 greater than or equal to 2  (Supplemental Figure 2). CASAVA tended to be more conservative in the
253 number of genes identified, with a substantial majority of these genes also selected by SHRiMP and
254 TopHat.
255
256 Cuffdiff identified 2719 genes that were selected as differentially expressed using both SHRiMP and
257 TopHat, of which 919 had a fold change greater than or equal to 2 (Supplemental Figure 3). Again,
258 there was a high degree of consistency between SHRiMP and TopHat, with nearly 90% of genes
259 identified using data mapped by both methods.
260
261 There were a total of 240 genes identified by all these analyses (Supplemental Figure 4). SAM,
262 implemented as described, was much more highly conservative in gene selection. Most genes identified
263 by SAM were also identified by Cuffdiff, while a majority of genes identified by Cuffdiff (74%) were
264 not identified by SAM.
265
266 In order to compare gene selection estimation with other mapping algorithms (TopHat and SHRiMP),
267 we subsequently, focused on Cuffdiff as well. A total of 267 genes were identified as differentially
268 expressed by SAM on CASAVA, and by Cuffdiff on both SHRiMP and TopHat , which had a fold
269 change greater than or equal to 2 (Figure 3).
270
271 *Molecular Validation*
272 Predicted changes were evaluated by qPCR for 52 of the genes (Table 5). Out of these 32 genes were
273 identified using CASAVA mapped data separately as reported previously [20]. We chose 10 additional
274 genes each that were uniquely selected by SHRiMP or TopHat for qPCR validation. These genes were
275 chosen based on prior knowledge of their relevance to oxidative stress response or lung biology in
276 general. Genes, that had a significant difference between the two groups (hyperoxia and controls) by
277 either t-test or Mann-Whitney U test at p-value less than 0.05, were designated as successfully
278 validated. We report detailed qPCR results for 32 genes identified using CASAVA mapped data

279  separately [20]. We validated differential expression for 23 of 32 (72%) of these genes selected by
280  CASAVA. Of these 32 genes, 29 were selected as differentially expressed by TopHat, and 20 of those
281  genes (69%) showed significant differences in expression by qPCR. Of these 32 genes, 28 were
282  selected as differentially expressed by SHRiMP, and 19 of those genes (63%) showed significant
283  differences in expression by qPCR. Of these 32 genes, 27 were selected as differentially expressed by
284  all three mapping methods, and 18 of those genes (67%) showed significant differences in expression
285  by qPCR.
286
287  Two of the 32 genes were uniquely selected by CASAVA, and both of those genes (100%) showed
288  significant differences in expression by qPCR. For TopHat, 7 of the 10 (70%) additional unique genes,
289  and a total of 27 of 39 (69%) genes tested showed significant differences in expression by qPCR. For
290  SHRiMP, 5 of the 10 (50%) additional unique genes, and a total of 24 of 38 (63%) genes tested showed
291  significant differences in expression by qPCR.

292   DISCUSSION
293   RNA-Seq is becoming the method of choice for genome-wide transcriptomics analysis, and has been
294   used to identify post-transcriptional changes and other modifications in human diseases such as cancer
295   [24]. Many methods for data processing and analysis are available, but the best approaches to estimate
296   differential expression for specific types of data sets are not at all clear. Here, we describe an empirical
297   assessment of the impact three different mapping algorithms upon the reliability of differential
298   expression estimation, in a typical genome-wide expression data set from an animal model of disease.
299   Our initial analysis of this data set indicated a substantial gene expression response associated with the
300   experimental challenge, such that it was appropriate for the current studies. The rationale behind
301   choosing these three included multiple factors, such as manufacturer recommendation, usage among
302   research community, cost, efficiency and ease of usage, among others.
303

304   Researchers can consider multiple mapping algorithms in their RNA-Seq analysis, which brings up
305   issues of interoperability. To achieve interoperability, input and output formats need to be
306   standardized. Currently the level of interoperability is high since most of the mapping algorithms
307   accept FASTQ format input files and generate SAM/BAM files as output. Prior publications
308   comparing various combinations of mapping algorithms and tests for differential gene expression have
309   found high level of consistency among the results [15, 16]. However, these prior studies did not include
310   subsequent attempts of molecular validation, which is a critical step in any differential expression
311   analysis.
312

313   It is important to point out that, of the three different mapping algorithms compared in this study, two
314   of them (SHRiMP and TopHat) were developed at academic institutions, and are freely available,
315   while the third (CASAVA) is commercial software that is provided by the manufacturer of the
316   sequencing instrument. We observed that SHRiMP mapping resulted in higher rates of genes detected
317   as expressed in the samples (80% for SHRiMP, 77% for CASAVA and 76% for TopHat), when
318   compared to either CASAVA or TopHat. This is likely due to the fact that SHRiMP allows reads to be
319   mapped to multiple loci, unlike TopHat and CASAVA, which require reads are mapped to a single
320   locus. Interestingly, a higher number of genes detected as expressed led to higher estimation of
321   differential expression for SHRiMP, but a somewhat lower level of accuracy as defined by qPCR.
322

323   We report here only a subset of the possible permutations of analysis that could be completed with the
324   mapping, normalization and gene selection methods we have included. In addition to comparing the
325   mappers, we also looked at the effects of methods of count generation by running correlation analysis
326   among the raw and normalized counts of the same samples, and found that even among methods using
327   different count generation approaches (CASAVA and TopHat), there was a high level of correlation
328   among the counts (Supplemental Figure 1). This indicated to us that the counting methods, independent
329   of mapping algorithm, may not have a big impact. Our analysis also revealed that there was a high
330   level of consistency among the two normalization methods (RPM and TM) on each mapped version of
331   data, when it comes to gene selection, irrespective of the test for differential expression used. We,
332   however did notice higher number of genes being identified by Cuffdiff when using FPKM normalized
333   counts. Even though SAM appeared to be effective in the current data set, applying this analytical
334   approach to other RNA-Seq datasets identified a number of limitations. For this, and other reasons
335   (e.g., free access to software, difficulties in generating bam files from CASAVA), we decided to use
336   SAM analysis of CASAVA data as our benchmark, and focus on the efficiency of Cuffdiff selection
337   using SHRiMP and TopHat mapping.
338

339  Encouragingly, we find that all mapping approaches performed similarly. Overall validation rates for
340  estimation of differential expression were comparable. CASAVA was most conservative with gene
341  selection, consistently estimating fewer genes as differentially expressed as compared to SHRiMP and
342  TopHat. However, this was associated with a higher frequency of validation; 72% for all genes, 100%
343  for unique genes. Conversely, SHRiMP tended to be slightly more liberal than TopHat, but had a
344  reduced validation rate than the other mapping approaches (63% for all genes, 50% for unique genes).
345  It is likely that this is at least partially due to its tolerance for including non-unique mapping for
346  individual reads. Interestingly, TopHat selected genes as differentially expressed at nearly the same
347  rate as SHRiMP, but its validation rate was much closer to that of CASAVA (69% for all genes, 70%
348  for unique genes). It is significant to note that individual genes selected by all mapping methods did
349  not demonstrate a higher validation rate than those identified by a single mapping method alone.
350
351  CONCLUSIONS
352  In summary, these data describe the differences that can be expected in the performance of these three
353  common mapping strategies, when applied to a typical genome-wide expression data set comparing
354  biological paradigms, such as an animal- or cell model response. Our data highlight the importance of
355  considering the analytical goals when choosing a data analysis approach. For instance, focused
356  analyses may want to consider a more conservative approach with a slightly higher validation rate,
357  while discovery approaches may be more tolerant to, and benefit from, more liberal approaches with a
358  slightly lower validation rate.

359 ACKNOWLEDMENTS
360 We would like to thank University of Rochester Center for Integrated Research Computing (CIRC) for
361 providing computational resources. We would also like to thank University of Rochester Medical
362 Center Genomics Research Center (GRC) for their support and Stephen L. Welle, Ph.D. for his
363 guidance and suggestions in preparation of the manuscript.
364
365 COMPETING INTERESTS
366 The authors declare that they have no competing interests.
367
368 AUTHOR CONTRIBUTIONS
369 CC carried out the sequence alignment, differential expression analysis, and drafted the manuscript. SB
370 carried out the differential expression analysis and drafted the manuscript. ZZ carried out qPCR
371 validation. AML, VAL participated in cDNA synthesis and qPCR. MY, BWB participated in
372 generating the animal data. MAO conceived of the study, and participated in its design. TJM conceived
373 of the study, participated in its design and coordination and helped to draft the manuscript. All authors
374 read and approved the final manuscript.
375

376 REFERENCES

1. Ingolia NT, Brar GA, Rouskin S, McGeachy AM, Weissman JS: **The ribosome profiling strategy for monitoring translation in vivo by deep sequencing of ribosome-protected mRNA fragments**. *Nature protocols* 2012, **7**(8):1534-1550.

2. Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, Snyder M: **The transcriptional landscape of the yeast genome defined by RNA sequencing**. *Science* 2008, **320**(5881):1344-1349.

3. Lister R, O'Malley RC, Tonti-Filippini J, Gregory BD, Berry CC, Millar AH, Ecker JR: **Highly integrated single-base resolution maps of the epigenome in Arabidopsis**. *Cell* 2008, **133**(3):523-536.

4. Cloonan N, Grimmond SM: **Transcriptome content and dynamics at single-nucleotide resolution**. *Genome biology* 2008, **9**(9):234.

5. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B: **Mapping and quantifying mammalian transcriptomes by RNA-Seq**. *Nature methods* 2008, **5**(7):621-628.

6. Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y: **RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays**. *Genome research* 2008, **18**(9):1509-1517.

7. Morin R, Bainbridge M, Fejes A, Hirst M, Krzywinski M, Pugh T, McDonald H, Varhol R, Jones S, Marra M: **Profiling the HeLa S3 transcriptome using randomly primed cDNA and massively parallel short-read sequencing**. *BioTechniques* 2008, **45**(1):81-94.

8. Wang Z, Gerstein M, Snyder M: **RNA-Seq: a revolutionary tool for transcriptomics**. *Nature reviews Genetics* 2009, **10**(1):57-63.

9. Fonseca NA, Rung J, Brazma A, Marioni JC: **Tools for mapping high-throughput sequencing data**. *Bioinformatics* 2012, **28**(24):3169-3177.

10. Lindner R, Friedel CC: **A comprehensive evaluation of alignment algorithms in the context of RNA-seq**. *PloS one* 2012, **7**(12):e52403.

11. Garber M, Grabherr MG, Guttman M, Trapnell C: **Computational methods for transcriptome annotation and quantification using RNA-seq**. *Nature methods* 2011, **8**(6):469-477.

12. Gambera D, Carta S, Crainz E, Fortina M, Maniscalco P, Ferrata P: **Metallosis due to impingement between the socket and the femoral head in a total hip prosthesis. A case report**. *Acta bio-medica : Atenei Parmensis* 2002, **73**(5-6):85-91.

13. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL: **TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions**. *Genome biology* 2013, **14**(4):R36.

14. David M, Dzamba M, Lister D, Ilie L, Brudno M: **SHRiMP2: sensitive yet practical SHort Read Mapping**. *Bioinformatics* 2011, **27**(7):1011-1012.

15. Nookaew I, Papini M, Pornputtapong N, Scalcinati G, Fagerberg L, Uhlen M, Nielsen J: **A comprehensive comparison of RNA-Seq-based transcriptome analysis from reads to differential gene expression and cross-comparison with microarrays: a case study in Saccharomyces cerevisiae**. *Nucleic acids research* 2012, **40**(20):10084-10097.

16. Soneson C, Delorenzi M: **A comparison of methods for differential expression analysis of RNA-seq data**. *BMC bioinformatics* 2013, **14**:91.

17. Benjamin AM, Nichols M, Burke TW, Ginsburg GS, Lucas JE: **Comparing reference-based RNA-Seq mapping methods for non-human primate data**. *BMC genomics* 2014, **15**(1):570.

421  18.  Zhang ZH, Jhaveri DJ, Marshall VM, Bauer DC, Edson J, Narayanan RK, Robinson GJ,
422       Lundberg AE, Bartlett PF, Wray NR *et al*: **A Comparative Study of Techniques for**
423       **Differential Expression Analysis on RNA-Seq Data**. *PloS one* 2014, **9**(8):e103207.
424  19.  Yee M, Vitiello PF, Roper JM, Staversky RJ, Wright TW, McGrath-Morrow SA, Maniscalco
425       WM, Finkelstein JN, O'Reilly MA: **Type II epithelial cells are critical target for hyperoxia-**
426       **mediated impairment of postnatal lung development**. *American journal of physiology Lung*
427       *cellular and molecular physiology* 2006, **291**(5):L1101-1111.
428  20.  Bhattacharya S, Zhou Z, Yee M, Chu CY, Lopez AM, Lunger VA, Solleti SK, Resseguie E,
429       Buczynski B, Mariani TJ *et al*: **The genome-wide transcriptional response to neonatal**
430       **hyperoxia identifies Ahr as a key regulator**. *American journal of physiology Lung cellular*
431       *and molecular physiology* 2014, **307**(7):L516-523.
432  21.  Tusher VG, Tibshirani R, Chu G: **Significance analysis of microarrays applied to the**
433       **ionizing radiation response**. *Proceedings of the National Academy of Sciences of the United*
434       *States of America* 2001, **98**(9):5116-5121.
435  22.  Trapnell C, Hendrickson DG, Sauvageau M, Goff L, Rinn JL, Pachter L: **Differential analysis**
436       **of gene regulation at transcript resolution with RNA-seq**. *Nature biotechnology* 2013,
437       **31**(1):46-53.
438  23.  Bhattacharya S, Go D, Krenitsky DL, Huyck HL, Solleti SK, Lunger VA, Metlay L, Srisuma S,
439       Wert SE, Mariani TJ *et al*: **Genome-wide transcriptional profiling reveals connective tissue**
440       **mast cell accumulation in bronchopulmonary dysplasia**. *American journal of respiratory*
441       *and critical care medicine* 2012, **186**(4):349-358.
442  24.  Maher CA, Kumar-Sinha C, Cao X, Kalyana-Sundaram S, Han B, Jing X, Sam L, Barrette T,
443       Palanisamy N, Chinnaiyan AM: **Transcriptome sequencing to detect gene fusions in cancer**.
444       *Nature* 2009, **458**(7234):97-101.
445
446

447    Figure 1: Analysis Workflow. RNA-Seq raw reads were mapped using three separate methods;
448    CASAVA, SHRiMP and TopHat. Mapped data were normalized using RPM or TM for all three
449    mapping methods, or FPKM for SHRiMP and TopHat only. Normalized data were filtered and tested
450    for differential expression using either SAM for all three mapping methods, or CuffDiff for SHRiMP
451    and TopHat. A subset of the differentially expressed genes was assessed by qPCR.

452

453    Figure 2: Correlation between mapped reads of the sample across three mappers. RPM normalized
454    counts from individual mappers for the same sample were plotted to identify the correlation between
455    the expression levels. Shown here are dot plots of two samples (Control: RA 11+12 and Hyperoxia: O2
456    11+12) between TopHat and CASAVA counts (A:Control, D:Hyperoxia), SHRiMP and CASAVA
457    (B:Control, E:Hyperoxia) and TopHat and SHRiMP (C:Control, F:Hyperoxia).

458

459    Figure 3: Summary of gene selection comparison. A comparison of genes with a fold-change greater
460    than or equal to 2 identified as differentially expressed by SAM using CASAVA mapped data, and by
461    Cuffdiff on SHRiMP and TopHat mapped data.

462 Table 1: Summary of mapping statistics.
463

| | CASAVA | SHRiMP | TopHat |
|---|---|---|---|
| No of genes in genome | 21718 | 23278 | 23278 |
| Median number of genes mapped | 16731 | 18674 | 17789 |
| Median number of raw reads (x $10^6$) | 16.53 | 16.53 | 16.53 |
| Median number of reads assigned to genes (x $10^6$) | 11.18 | 12.75 | 12.56 |
| Median alignment rate (%) | NA | 85.26 | 74.94 |
| Median gene detection rate (%) | 77.03 | 80.21 | 76.42 |

464
465 Table 2: Differential expression estimated by CASAVA.
466

| CASAVA | SAM | | |
|---|---|---|---|
| | RPM | TM | Overlap |
| Significant Genes | 1017 | 1019 | 798 |
| FC > 2 | 472 | 446 | 300 |
| No. of up-regulated | 210 | 180 | 132 |
| No. of down-regulated | 262 | 266 | 168 |
| *% of up-regulated* | *0.44* | *0.40* | *0.44* |
| *% of down-regulated* | *0.56* | *0.60* | *0.56* |

467
468 Table 3: Differential expression estimation by SHRiMP.
469

| SHRiMP | SAM | | | Cuffdiff |
|---|---|---|---|---|
| | RPM | TM | Overlap | FPKM |
| Significant Gene | 937 | 879 | 857 | 3886 |
| FC > 2 | 401 | 386 | 379 | 1087 |
| No. of up-regulated | 170 | 157 | 153 | 623 |
| No. of down-regulated | 231 | 229 | 226 | 464 |
| *% of up-regulated* | *0.42* | *0.41* | *0.40* | *0.57* |
| *% of down-regulated* | *0.58* | *0.59* | *0.60* | *0.43* |

470
471

472 Table 4: Differential expression estimation by TopHat.
473

| TopHat | SAM | | | Cuffdiff |
|---|---|---|---|---|
| | RPM | TM | Overlap | FPKM |
| Significant Gene | 951 | 880 | 860 | 2831 |
| FC > 2 | 418 | 412 | 396 | 1044 |
| No. of up-regulated | 179 | 172 | 165 | 593 |
| No. of down-regulated | 239 | 240 | 231 | 451 |
| *% of up-regulated* | *0.43* | *0.42* | *0.42* | *0.57* |
| *% of down-regulated* | *0.57* | *0.58* | *0.58* | *0.43* |

474
475 Table 5: qPCR validation rate. Genes were significantly different by either T-Test or MWU at p<0.05
476

| Mapping Programs | No. of Gene Chosen | No. of Gene Validated | Validation Rate (%) | Validation Rate of Unique Gene (%) |
|---|---|---|---|---|
| CASAVA | 32 | 23 | 71.87% | 100% |
| TopHat | 39 | 27 | 69.23% | 70% |
| SHRiMP | 38 | 24 | 63.16% | 50% |
| Overlap | 27 | 18 | 66.67% | |

477

**Figure 1**

Analysis Workflow.

RNA-Seq raw reads were mapped using three separate methods; CASAVA, SHRiMP and TopHat. Mapped data were normalized using RPM or TM for all three mapping methods, or FPKM for SHRiMP and TopHat only. Normalized data were filtered and tested for differential expression using either SAM for all three mapping methods, or CuffDiff for SHRiMP and TopHat. A subset of the differentially expressed genes was assessed by qPCR.
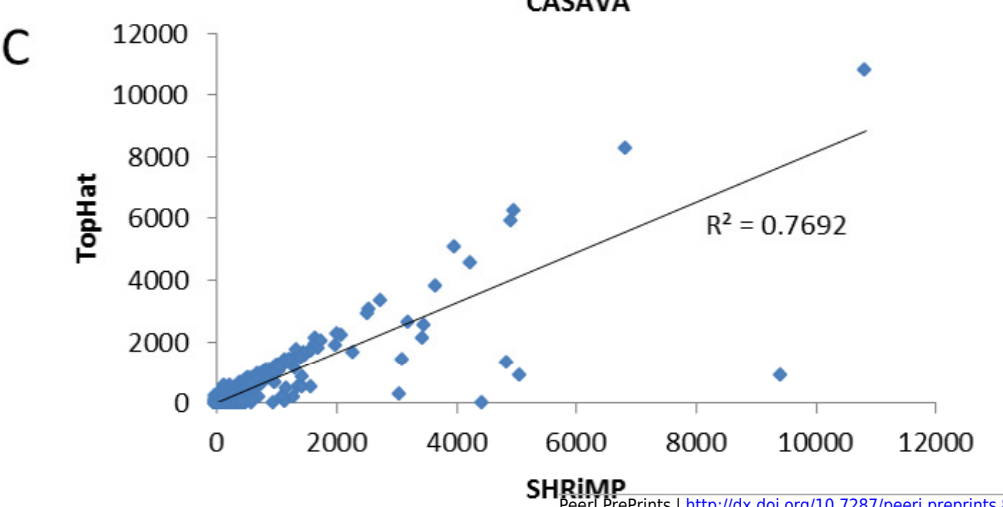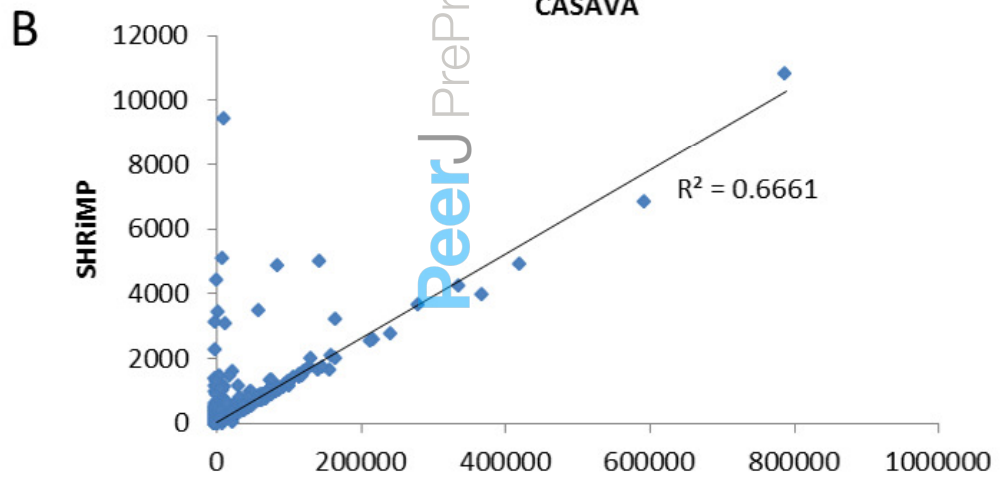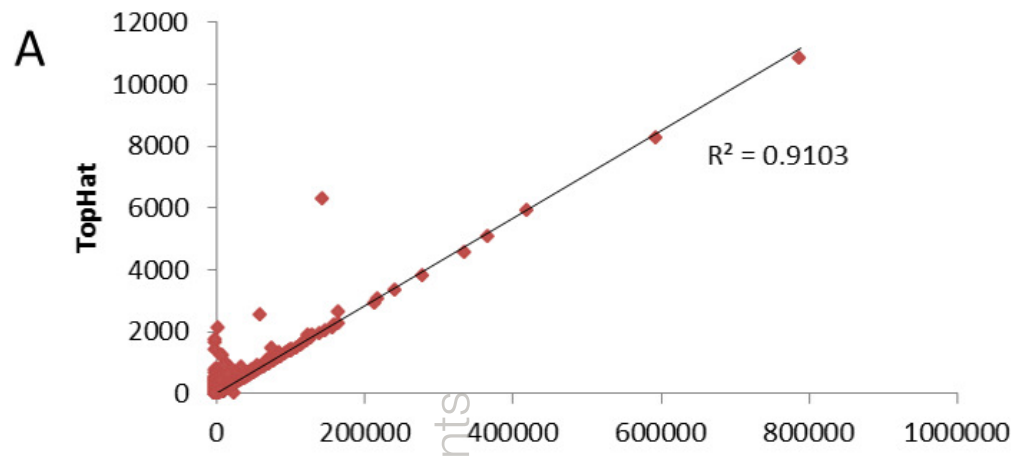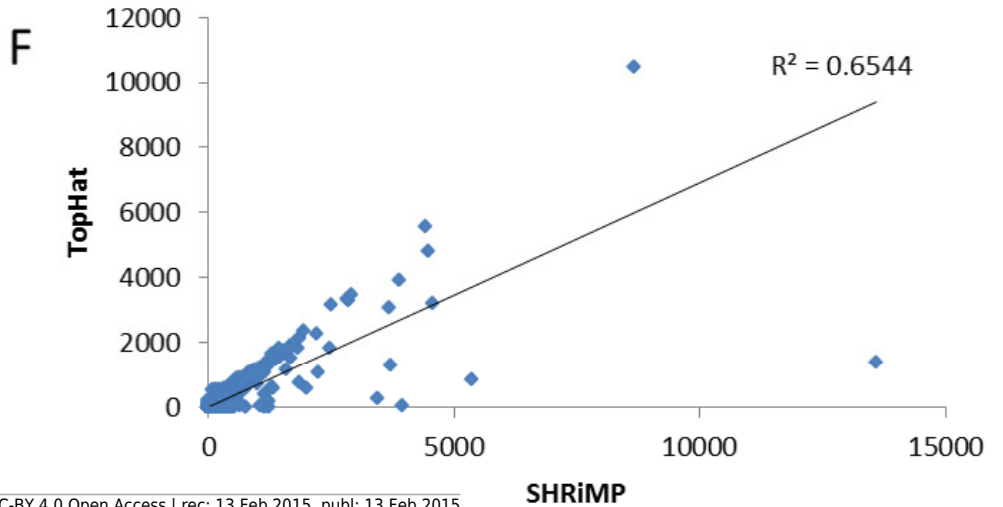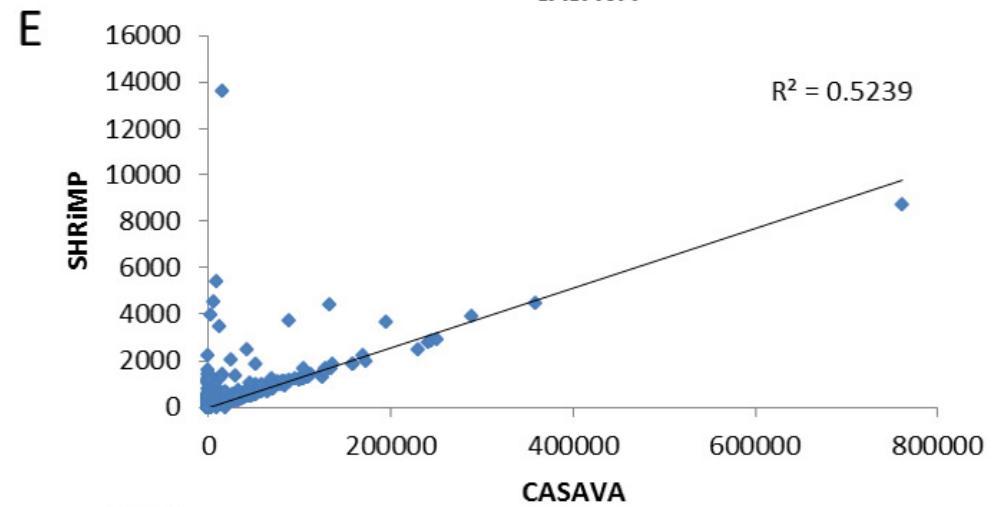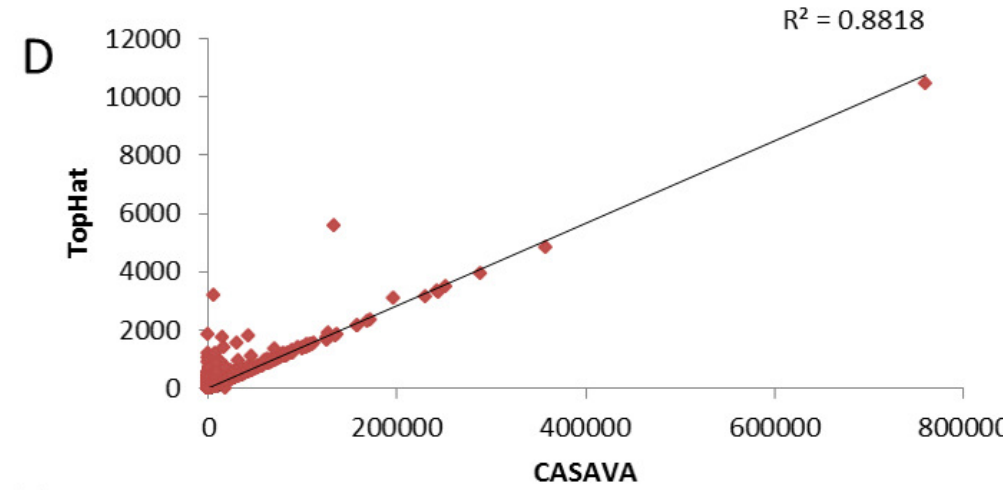
**Figure 2**

Correlation between mapped reads of the sample across three mappers.

RPM normalized counts from individual mappers for the same sample were plotted to identify the correlation between the expression levels. Shown here are dot plots of two samples (Control: RA 11+12 and Hyperoxia: O2 11+12) between TopHat and CASAVA counts (A:Control, D:Hyperoxia), SHRiMP and CASAVA (B:Control, E:Hyperoxia) and TopHat and SHRiMP (C:Control, F:Hyperoxia). @0�

# Control

## Hyperoxia

A

$R^2 = 0.9103$

B

$R^2 = 0.6661$

C

$R^2 = 0.7692$

D

$R^2 = 0.8818$

E

$R^2 = 0.5239$

F

$R^2 = 0.6544$

## Figure 3(on next page)

Summary of gene selection comparison.

A comparison of genes with a fold-change greater than or equal to 2 identified as differentially expressed by SAM using CASAVA mapped data, and by Cuffdiff on SHRiMP and TopHat mapped data.