

1 Title: Evaluation of TagSeq, a reliable low-cost alternative for RNAseq

2 Authors: Brian K. Lohman^{1*}, Jesse N. Weber^{1,2}, and Daniel I. Bolnick¹

3 ¹Department of Integrative Biology, University of Texas at Austin, One University Station C0990,

4 Austin, TX, 78712, United States

5 ²Current Address: Division of Biological Sciences, The University of Montana-Missoula, Missoula,

6 MT, 59812, United States

7

8 Keywords: RNAseq, 3' Tag-based sequencing

9 *Corresponding Author. E-mail: lohman@utexas.edu

10 Running title: Benchmarking TagSeq; a low-cost alternative for RNAseq

11 **Abstract:**

12 RNAseq is a relatively new tool for ecological genetics that offers researchers insight into changes
13 in gene expression in response to a myriad of natural or experimental conditions. However, standard
14 RNAseq methods (e.g., Illumina TruSeq® or NEBNext®) can be cost prohibitive, especially when
15 study designs require large sample sizes. Consequently, RNAseq is often underused as a method, or
16 is applied to small sample sizes that confer poor statistical power. Low cost RNAseq methods could
17 therefore enable far greater and more powerful applications of transcriptomics in ecological
18 genetics and beyond. Standard mRNAseq is costly partly because one sequences portions of the full
19 length of all transcripts. Such whole-mRNA data is redundant for estimates of relative gene
20 expression. TagSeq is an alternative method that focuses sequencing effort on mRNAs' 3' end,
21 thereby reducing the necessary sequencing depth per sample, and thus cost. Here we present a
22 revised TagSeq protocol, and compare its performance against NEBNext®, the “gold-standard”
23 whole mRNAseq method. We built both TagSeq and NEBNext® libraries from the same biological
24 samples, each spiked with control RNAs. We found that TagSeq measured the control RNA
25 distribution more accurately than NEBNext®, for a fraction of the cost per sample (~10%). The
26 higher accuracy of TagSeq was particularly apparent for transcripts of moderate to low abundance.
27 Technical replicates of TagSeq libraries are highly correlated, and were correlated with NEBNext®
28 results. Overall, we show that our modified TagSeq protocol is an efficient alternative to traditional
29 whole mRNAseq, offering researchers comparable data at greatly reduced cost.

30 **Introduction:**

31 RNAseq has been widely used to describe differences in gene expression among wild populations,
32 as well as changes in captive or wild individuals' expression following exposure to stimuli (mates,
33 predators, parasites, abiotic stress, toxins). This work has helped uncover the genetic basis of
34 complex traits, implicate genes underlying targets of natural selection, and measure the heritable
35 and environmental components of variation in gene expression [1-6]. However, the most widely
36 used RNAseq protocols are cost-prohibitive for many biologists, including but not limited to
37 researchers in ecological genetics.

38 Construction of any whole mRNAseq library for the Illumina platform (including Illumina
39 TruSeq® and NEBNext® kits) involves isolating or enriching for mRNA, which is then fragmented
40 and subject to massively parallel sequencing. The resulting data yields sequences for overlapping
41 portions of the entire lengths of the original messenger RNAs (hence 'whole' mRNAseq). This
42 requires high depth of coverage; although sequencing requirements vary depending on sample type,
43 the ENCODE Consortium suggests ~30 million raw reads per sample as a "best practice" for most
44 RNAseq experiments [7], limiting researchers to a maximum of eight samples per lane of Illumina
45 HiSeq 2500. The high cost of sequencing, combined with high cost of library construction, has
46 forced many studies to use small sample sizes, or pool samples within treatments. This is cause for
47 concern, as meaningful differences in gene expression simply may not be detected with such low-
48 powered sampling designs, and pooled RNAseq may fail to properly account for residual variation
49 in expression.

50 To resolve problems with whole mRNAseq, several low-cost alternatives have been
51 developed. Most notably, Meyer et al. 2011 [8] presented a 3' Tag-based approach to RNAseq,
52 called TagSeq, that requires little input RNA, involves low library construction costs, and requires
53 many fewer raw reads per sample. By focusing on the 3' end of mRNA fragments, TagSeq reduces
54 the sequencing effort required to characterize a population of mRNAs in a biological sample. This
55 cost-saving does come with some constraints: TagSeq cannot distinguish between alternatively

56 spliced transcripts from a single locus, and will not identify polymorphism or allele-specific
57 expression in much of a genes' coding sequence. However, the benefits of precisely measuring
58 locus-level transcriptional differences with high replication may outweigh the lack of splicing or
59 SNP information for many experiments in ecological systems. However, as presented in Meyer et
60 al. 2011 [8], TagSeq uses a number of outdated methods and enzymes, which may skew the
61 distribution of RNA fragments in the library, with respect to both fragment size and GC content [9].
62 In addition, the accuracy of TagSeq has not yet been compared to the industry standard
63 TruSeq[®]/NEBNext[®] which reliably measures moderate and high abundance mRNAs in a sample.

64 Here, we present a modified protocol intended to increase the accuracy and precision of
65 TagSeq, by incorporating recent findings on polymerase performance, fragmentation methods, and
66 bead-based purification technology into the library construction process. We then tested the
67 accuracy of TagSeq against the industry standard NEBNext[®] by sequencing technical replicates of a
68 biological sample, each containing an artificial set of diverse RNAs of known concentration,
69 designed by the External RNA Controls Consortium (hereafter simply "ERCC").

70

71 **Materials and Methods:**

72 *Improvements to TagSeq library construction:*

73 Briefly, our improved TagSeq library construction method involves 11 steps: 1) isolate total RNA,
74 2) remove genomic DNA with DNase (if not included in total RNA isolation), 3) fragment total
75 RNA with Mg⁺ buffer (NEB), 4) synthesize cDNA with a poly-dT oligo, 5) PCR amplify cDNA, 6)
76 purify PCR products with DNA-binding magnetic beads (Agencourt, or made in-house [10]), 7)
77 fluorometrically quantify PCR products (PicoGreen, Life Technologies), 8) normalize among-sample
78 concentrations, 9) add sample-specific barcodes via PCR, 10) pool samples and select a small range
79 of fragment sizes (to maximize output on the Illumina platform) via automated gel extraction (400-
80 500bp, Sage Pippin Prep 2% agarose), 11) quantify concentrations of post-extraction products via
81 Qubit, 12) normalize among pools. This protocol can be completed by a single researcher in three

82 days, and this approach is optimized for 96-well format plates. Improvements over the original
83 protocol are described in Table 1.

84 Total RNA was extracted from six freshly isolated stickleback (*Gasterosteus aculeatus*) head
85 kidneys stored in RNAlater (Ambion). All fish were lab-raised non-gravid females, bred via *in vitro*
86 crosses of wild caught parents. Three fish originated from crosses between parents from Gosling
87 Lake, British Columbia and three fish from crosses between parents from Roberts Lake, British
88 Columbia. Total RNA from all 6 head kidneys was then split, and libraries were constructed with
89 both whole mRNAseq and TagSeq methods. Four whole mRNAseq libraries (NEBNext®
90 directional RNA libraries with poly-A enrichment) were prepared according to the manufacturer's
91 instruction, by the Genomic Sequencing and Analysis Facility at the University of Texas at Austin,
92 with the addition of ERCC before library construction began, according to the manufacturer's
93 instructions. Whole mRNAseq samples were sequenced on a single lane of Illumina HiSeq 2500
94 2x100, producing an average of 40.5 million paired-end reads per sample (81 million reads total per
95 sample). Following the addition of ERCC to one technical replicate per biological sample, TagSeq
96 samples were prepared according to Meyer et al 2011 [8], but with changes detailed in Table 1.
97 Four TagSeq samples had two technical replicates (totally independent library builds from total
98 RNA) and a fifth had three technical replicates. TagSeq libraries (29 total, including 17 outside the
99 scope of this work) were sequenced on 3 lanes of Illumina HiSeq 2500 1x100, producing an
100 average of 10.3 million raw reads per sample.

101

102 *Bioinformatics*

103 Raw whole mRNAseq reads were trimmed with Cutadapt v 1.3 [11] to remove any adapter
104 contamination. We then mapped the trimmed reads to version 79 of the stickleback genome (with
105 ERCC sequences appended) using BWA-MEM [12], and counted genes using Bedtools [13],
106 producing 20,678 total genes. TagSeq reads were processed according to the iRNAseq pipeline
107 (https://github.com/z0on/tag-based_RNAseq) [14], producing 19,145 total genes.

108

109 *Statistical analysis of control transcripts*

110 For each sample, we plotted observed counts of artificial ERCC transcripts against expected values,
111 fitting a simple linear model (observed ~ expected). We tested for a difference in mean adjusted R^2
112 value between library construction methods with a paired t-test (paired by biological sample).

113 We calculated the Spearman correlation between observed log transformed counts of ERCC
114 transcripts and expected transcript quantity. We tested for a difference in mean Rho values between
115 library construction methods using a paired t-test. We also considered Rho separately for abundance
116 quartiles.

117

118 *Statistical analysis of stickleback transcripts*

119 We calculated the Spearman correlation among TagSeq technical replicates. We calculated the
120 Spearman correlation between stickleback head kidney samples which had been prepared using
121 both library construction methods.

122

123 *Statistical analysis of inline barcodes*

124 TagSeq, as presented by Meyer et al. and here, uses degenerate inline barcodes on the 5' end of
125 each fragment to identify PCR duplicates. We tested for the random incorporation of these barcodes
126 with a Chi Squared test. We also tested for the effect of increased GC content within each barcode
127 on the number of times that barcode was observed with a Poisson GLM.

128

129 **Results:**

130 We found that, when fitting a linear model between the expected concentrations of ERCC to
131 observed transcript counts, TagSeq had a significantly higher mean adjusted R^2 value ($R^2 = 0.89$)
132 than NEBNext[®] ($R^2 = 0.80$, Figure 1, observed ~ expected, paired t-test, $t = 18.63$, $df = 3$, $p <$
133 0.001). Similarly, the rank correlation between observed and expected ERCC fragments was

134 consistently higher for TagSeq (mean Rho = 0.94) than NEBNext[®] (mean Rho = 0.87, Figure 2,
135 paired T-test, $t = 12.20$, $df = 3$, $p = 0.001$). TagSeq showed higher mean Rho values for all
136 abundance classes except the third quartile. Most notably, whole mRNAseq performed very poorly
137 in the lowest abundance class (relative concentration of 0.014-0.45 attamols/ μ l), and TagSeq
138 substantially outperformed whole mRNAseq in the second abundance class (relative concentration
139 of 0.92-7.3 attamols/ μ l, Figure 3).

140 With respect to stickleback (non-control) sequences, the mean Rho among technical
141 replicates of TagSeq samples was 0.96 ($n=5$, calculate Rho for each biological sample and average).
142 Due to the high cost of NEBNext[®] library generation and sequencing (~\$340 per sample), we did
143 not perform technical replicates using this method. We found a strong significant positive
144 correlation between stickleback gene counts generated with TagSeq and whole mRNAseq (Rho =
145 0.74, $p < 0.001$). This is likely an underestimate of the actual correlation between the two library
146 construction methods because whole mRNAseq performs very poorly when RNAs are in moderate
147 to low abundance (first and second abundance classes, Figure 3). Given that 9,572 loci are in the
148 bottom half of gene counts, even small differences in absolute counts between the methods will
149 strongly influence the rank-based statistic.

150 We also wished to compare our new method with that of the original TagSeq protocol, but
151 cannot make a direct comparison with the available samples. Meyer et al. (2011) evaluated their
152 accuracy by comparing fold-differences in differentially expressed genes (between experimental
153 treatments), whereas we measured accuracy using the estimates of relative abundance of ERCC.
154 Keeping in mind these different benchmarking methods, we can draw a rough comparison. The
155 original TagSeq method yielded a correlation of $r = 0.86$ between TagSeq estimates of fold-change
156 expression, and qPCR measures of the same fold change (a “known” benchmark). In contrast, our
157 protocol yields a correlation of Rho = 0.94 between our relative abundance estimates, and the
158 known ERCC relative abundances. We infer that the new protocol performs at least as well, and

159 probably better, than the previous protocol, at generating expression level estimates that resemble
160 known values.

161 The iRNAseq pipeline includes the removal of PCR duplicates, which are a common
162 problem in many library construction methods [9]. Any reads which meet two criteria are called
163 PCR duplicates and removed: 1) identical in-line barcodes; four degenerate bases at the start of each
164 read, and 2) the first 30 bases of sequence after the in-line barcode are identical. The removal of
165 PCR duplicates substantially reduces the number of TagSeq reads in each library (mean reduction
166 of 70.3%, $n = 12$). However, this avoids potential bias introduced by PCR, namely over
167 representation of smaller fragments. We found that inline barcodes were incorporated non-randomly
168 (Chi Square = 10,500,000, $df = 63$, $p \ll 0.001$). We found that increased GC content in the inline
169 barcode significantly reduced the number of times a barcode was observed. For every G or C added
170 to the inline barcode, the expected value of the number of observed barcodes is reduced by ~2.9%
171 ($\text{count} \sim \text{gcContent}$, family = poisson, $\beta_{\text{gcContent}} = -0.133$, $p < 0.001$).

172

173 **Discussion:**

174 We present a number of methodological improvements to the TagSeq method of Meyer et al. 2011,
175 and taken the important next step of comparing the new protocol to the NEBNext® kit, the industry-
176 standard for whole mRNAseq. Overall, our results illustrate that the updated TagSeq method offers
177 researchers the ability to dramatically increase sample sizes in gene expression analyses, which will
178 facilitate testing for more subtle transcriptional differences than traditional whole mRNAseq
179 methods.

180 While TagSeq has been used predominantly in corals [14, 15], it should be applicable to
181 nearly all metazoans. However, we caution researchers to perform several basic checks during
182 TagSeq library construction, most especially ensuring the size distribution of RNA fragments is as
183 narrow as possible during total RNA fragmentation. We recommend evaluating the results of
184 various total RNA fragmentation times via BioAnalyzer. Fragments should be larger than 100 bp

185 and smaller than 500 bp (see supplementary materials). Here we were interested in evaluating the
186 robustness of the TagSeq method for threespine stickleback, and therefore sequenced stickleback
187 transcripts more deeply than required for an accurate estimate of gene expression across the
188 majority of expressed loci (we generated an average of 10.3 million raw reads per sample). We
189 recommend that researchers aim for ~5-6 M raw reads per sample if the goal is to measure the top
190 75% of all expressed mRNAs in a sample, as this has produce sufficient gene counts for robust
191 statistical power in an invertebrate, a plant, and stickleback (M. Matz and T. Jeunger, personal
192 communications).

193 In this project, we intentionally under-loaded our TagSeq libraries on the HiSeq lane by
194 15% (0.0017 pmols loaded), anticipating that low base diversity on the 5' end of the fragments (the
195 inline barcode used to detect and remove PCR duplicates) would lead to poor clustering. However,
196 quality metrics from the HiSeq run indicate that this is not a problem. We observed ~500-600
197 clusters per mm² on each tile, and the majority of these clusters passed filtering (low base diversity
198 or overclustering would generate large numbers of clusters with few passing filtering). We therefore
199 recommend that users load the expected quantity or even 10-20% extra material on each lane of
200 HiSeq (see supplementary material). Overloading TagSeq libraries may help to increase raw read
201 yield, relative to NEBNext[®] (optimally clustered at ~1000 clusters per mm² when 0.002 pmols
202 loaded). We also emphasize that small fragments need to be removed from TagSeq libraries, as they
203 will more easily cluster on the HiSeq, reducing read output. These may be identified by
204 BioAnalyzer and removed with additional bead clean-ups.

205 Several of our methodological changes aimed to mitigate the number of PCR duplicates,
206 which are artifacts of all PCR-related methods. First, we predicted that adding two additional
207 degenerate bases to the inline barcodes (the first four sequenced bases of every adapter, which were
208 coded as degenerate bases in the old TagSeq method) would not only increase our ability to detect
209 independent transcripts from PCR duplicates, and also increase base diversity on the 5' end of
210 fragments, thereby increasing the number of clusters passing Illumina's quality filters. However,

211 this alteration did not ameliorate the problem of PCR duplicates or increase the number of raw reads
212 generated in each lane (data not shown). In the future we recommend that protocol users consider
213 replacing the degenerate bases in inline barcodes with 3-nitropyrrole, as this should better
214 randomize which bases are incorporated during initial round of PCR [16]. Second, we limited our
215 number of PCR cycles to 12. Empirically testing the effects of PCR cycle number on TagSeq
216 accuracy was outside the scope of the present study. However, it is widely accepted that the best
217 way to limit bias is to reduce the number of PCR cycles during cDNA amplification as much as
218 possible [9].

219 In summary, we show that the improved TagSeq method has both benefits and drawbacks
220 compared to traditional whole mRNA sequencing. While our TagSeq libraries did not generate
221 optimal numbers of clusters on the HiSeq platforms, but we identify several potential solutions to
222 the problem. Regardless of the slightly lower number of raw reads, our improved TagSeq method
223 remains far and away much more cost effective than whole mRNAseq. At maximal efficiency (32
224 individuals per sequencing lane), our method was able to produce highly accurate, transcriptome-
225 wide gene counts for only ~\$33 per sample, including sequencing costs (one lane of HiSeq 2500 V3
226 chemistry with ~5.6 M raw reads per sample). This low cost and high reliability offer molecular
227 ecologists the opportunity to vastly increase sample sizes and increase replication to uncover new
228 patterns in gene expression.

229

230 **Acknowledgements:** All live animal research was approved by the UT Austin IACUC [protocol
231 AUP-2013-00012]; and collections were approved by the British Columbia Ministry of
232 Environment [Scientific Fish Collection permit NA09-52421]. We wish to thank Mikhail Matz,
233 Marie Strader, and the Juenger lab for fruitful discussion on improvements to the TagSeq method
234 both during library construction and analysis. Figures 2 and 3 were generated using plotting
235 functions written by Luke Reding (<https://github.com/lukereding/readingPlot>). This work was
236 supported by the Howard Hughes Medical Institute (DIB).

237

238 **Data Accessibility:** Meta data, code for raw read processing, gene counts, code for statistical
239 analysis, and plotting of data, BioAnalyzer .XAD files, and HiSeq quality metrics, and detailed
240 protocol are located in DRYAD entry: <http://dx.doi.org/10.5061/dryad.vq275>

241

242 Raw sequence reads will be stored on Corral, a permanent data repository with multiple,
243 independent backups, located and owned by the University of Texas at Austin Texas Advanced
244 Computing Center. Users can download data at any time via secure copy.

245

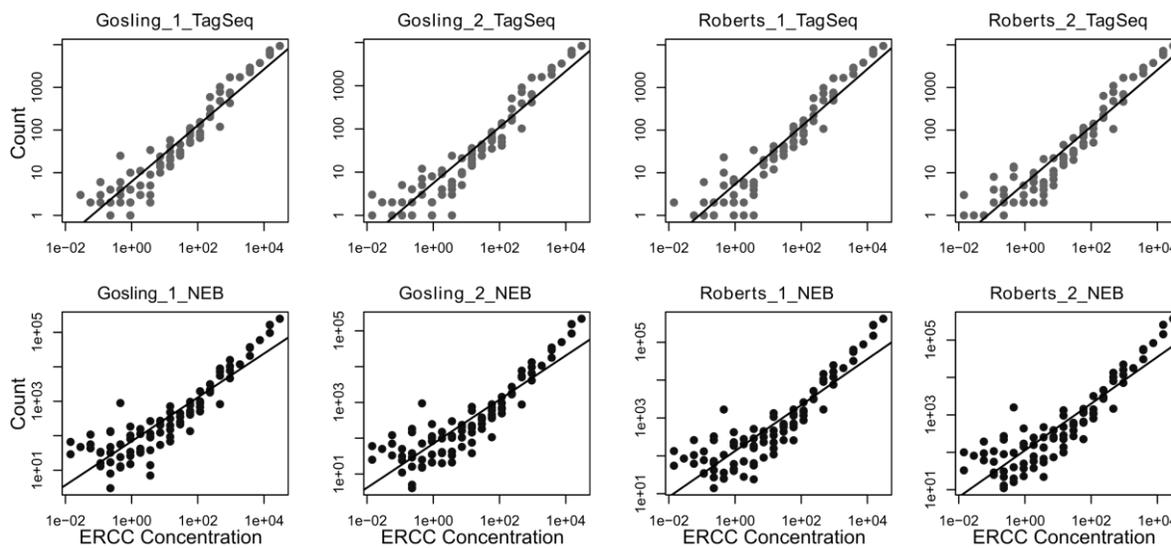
246 **Author Contributions:** BKL, JNW, and DIB jointly designed the research. BKL carried out all
247 library construction improvements. BKL, and JNW analyzed data. BKL wrote the manuscript, with
248 comments from JNW and DIB. All authors approved the final version.

249 **Table 1.** Changes to Meyer et al. 2011. We identified a number of areas where the Meyer protocol
250 could be improved and implemented changes to address these concerns.
251

Problem	Solution
Quantification of DNA/RNA by spectroscopy is inaccurate.	Fluorescent based quantification with Quant-iT assays.
Genomic DNA contamination leads to nonspecific amplification.	Increase DNase treatment to 1.5x concentration at 37°C for 1 hour.
Fragmentation of total RNA with Tris buffer produces a wide distribution of fragment sizes.	Precisely fragment total RNA with a specialized Mg ⁺ buffer.
Yield of first strand synthesis is too variable.	Normalize RNA input to 1µg.
Variable GC content among fragments can cause dropout of transcripts [9].	Use AccuPrime Taq polymerase and associated thermal profile for PCR steps [9].
Excessive PCR amplification increases the number of PCR duplicates.	Reduce number of PCR cycles to 12 or less.
Purification using solid-phase methods (e.g. spin columns) is not high throughput compatible, inefficient and costly.	Clean with Agencourt AMPure beads, which can be made in-house [10]
Post-PCR cDNA amplification yield is highly variable.	Normalize input to 40ng total.
Size selection by standard gel extraction is highly variable.	Precise size selection with Pippin Prep automated gel extraction.
Mixing individual libraries based on qPCR is slow and expensive.	Normalize lanes of Pippin Prep with Qubit Fluorimeter.

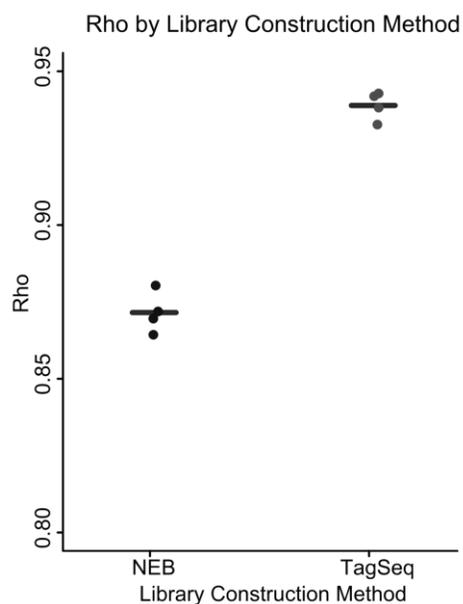
252

253 **Figure 1.** Regression of observed vs. expected ERCC transcripts shows TagSeq has higher adjusted
254 R^2 values for samples prepared with both methods (paired T-test, $t = 18.63$, $df = 3$, $p < 0.001$).
255



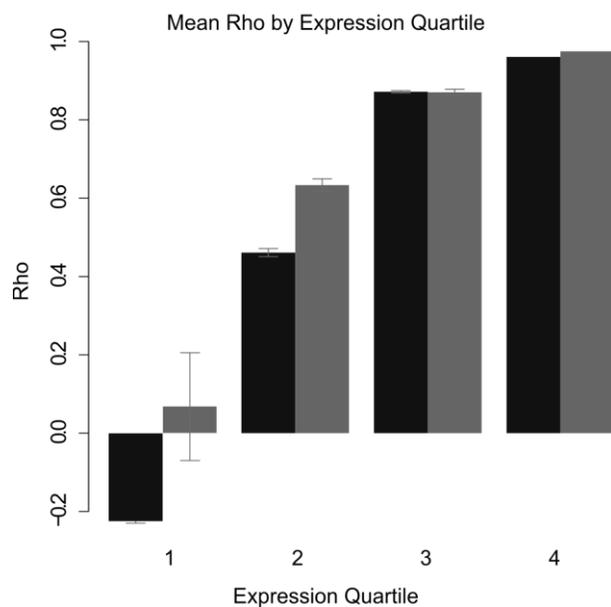
256

257 **Figure 2.** TagSeq more accurately recovers a known distribution of control mRNA fragments
258 (ERCC) than whole mRNAseq (mean Rho for TagSeq is higher than mean Rho for whole
259 mRNAseq, paired T-test, $t = 12.20$, $df = 3$, $p = 0.001$).
260



261

262 **Figure 3.** Breakdown of control mRNAs by abundance class shows that TagSeq recovers mRNAs
263 better than TruSeq, especially at lower abundances. Light grey bars are TagSeq, dark grey bars are
264 whole mRNAseq. Fences indicate standard error.
265



266

References

- 267
268 1. Videvall, E., C.K. Cornwallis, V. Palinauskas, G. Valkiūnas, and O. Hellgren, *The Avian*
269 *Transcriptome Response to Malaria Infection*. *Molecular biology and evolution*, 2015.
270 **32**(5): p. 1255-1267.
- 271 2. Foth, B.J., I.J. Tsai, A.J. Reid, A.J. Bancroft, S. Nichol, A. Tracey, N. Holroyd, J.A. Cotton,
272 E.J. Stanley, and M. Zarowiecki, *Whipworm genome and dual-species transcriptome*
273 *analyses provide molecular insights into an intimate host-parasite interaction*. *Nature*
274 *genetics*, 2014.
- 275 3. Barribeau, S.M., B.M. Sadd, L. du Plessis, and P. Schmid-Hempel, *Gene expression*
276 *differences underlying genotype-by-genotype specificity in a host-parasite system*.
277 *Proceeding of the National Academy of Sciences*, 2014. **111**: p. 3496-3501.
- 278 4. Lenz, T.L., C. Eizaguirre, B. Rotter, M. Kalbe, and M. Milinski, *Exploring local*
279 *immunological adaptation of two stickleback ecotypes by experimental infection and*
280 *transcriptome-wide digital gene expression analysis*. *Molecular Ecology*, 2013. **22**(3): p.
281 774-86.
- 282 5. Lovell, J.T., J.L. Mullen, D.B. Lowry, K. Awole, J.H. Richards, S. Sen, P.E. Verslues, T.E.
283 Juenger, and J.K. McKay, *Exploiting Differential Gene Expression and Epistasis to*
284 *Discover Candidate Genes for Drought-Associated QTLs in Arabidopsis thaliana*. *The*
285 *Plant Cell*, 2015. **27**(4): p. 969-983.
- 286 6. Pickrell, J.K., J.C. Marioni, A.A. Pai, J.F. Degner, B.E. Engelhardt, E. Nkadori, J.-B.
287 Veyrieras, M. Stephens, Y. Gilad, and J.K. Pritchard, *Understanding mechanisms*
288 *underlying human gene expression variation with RNA sequencing*. *Nature*, 2010.
289 **464**(7289): p. 768-772.
- 290 7. The ENCODE Consortium. *Standards, Guidelines, and Best Practices for RNA-seq*. 2011;
291 Available from:
292 [https://genome.ucsc.edu/ENCODE/protocols/dataStandards/ENCODE_RNAseq_Standards](https://genome.ucsc.edu/ENCODE/protocols/dataStandards/ENCODE_RNAseq_Standards_V1.0.pdf)
293 [_V1.0.pdf](https://genome.ucsc.edu/ENCODE/protocols/dataStandards/ENCODE_RNAseq_Standards_V1.0.pdf).
- 294 8. Meyer, E., G.V. Aglyamova, and M.V. Matz, *Profiling gene expression responses of coral*
295 *larvae (Acropora millepora) to elevated temperature and settlement inducers using a novel*
296 *RNA-Seq procedure*. *Molecular Ecology*, 2011: p. no-no.
- 297 9. Aird, D., M.G. Ross, W.S. Chen, M. Danielsson, T. Fennell, C. Russ, D.B. Jaffe, C.
298 Nusbaum, and A. Gnirke, *Analyzing and minimizing PCR amplification bias in Illumina*
299 *sequencing libraries*. *Genome Biology*, 2011. **12**: p. R18.
- 300 10. Rohland, N. and D. Reich, *Cost-effective, high-throughput DNA sequencing libraries for*
301 *multiplexed target capture*. *Genome research*, 2012. **22**(5): p. 939-946.
- 302 11. Martin, M., *Cutadapt removes adapter sequences from high-throughput sequencing reads*.
303 *EMBnet. journal*, 2011. **17**(1): p. pp. 10-12.
- 304 12. Li, H. and R. Durbin, *Fast and accurate long-read alignment with Burrows–Wheeler*
305 *transform*. *Bioinformatics*, 2010. **26**(5): p. 589-595.
- 306 13. Quinlan, A.R. and I.M. Hall, *BEDTools: a flexible suite of utilities for comparing genomic*
307 *features*. *Bioinformatics*, 2010. **26**(6): p. 841-842.
- 308 14. Dixon, G.B., S.W. Davies, G.A. Aglyamova, E. Meyer, L.K. Bay, and M.V. Matz,
309 *Genomic determinants of coral heat tolerance across latitudes*. *Science*, 2015. **348**(6242):
310 p. 1460-1462.
- 311 15. Des Marais, D.L., W.D. Skillern, and T.E. Juenger, *Deeply diverged alleles in the*
312 *Arabidopsis AREB1 transcription factor drive genome-wide differences in transcriptional*
313 *response to the environment*. *Molecular biology and evolution*, 2015. **32**(4): p. 956-969.
- 314 16. Schweyen, H., A. Rozenberg, and F. Leese, *Detection and Removal of PCR Duplicates in*
315 *Population Genomic ddRAD Studies by Addition of a Degenerate Base Region (DBR) in*
316 *Sequencing Adapters*. *The Biological Bulletin*, 2014. **227**(2): p. 146-160.
- 317