

## On the optimal trimming of high-throughput mRNA sequence data

Matthew D MacManes<sup>1,2,3\*</sup>

**1** *University of New Hampshire. Durham, NH 03824*

**2** *Department of Molecular, Cellular & Biomedical Sciences*

**3** *Hubbard Center for Genome Studies*

\* Corresponding author: [macmanes@gmail.com](mailto:macmanes@gmail.com), Twitter: [@PeroMHC](https://twitter.com/PeroMHC)

### Abstract

The widespread and rapid adoption of high-throughput sequencing technologies has changed the face of modern studies of evolutionary genetics. Indeed, newer sequencing technologies, like Illumina sequencing, have afforded researchers the opportunity to gain a deep understanding of genome level processes that underlie evolutionary change. In particular, researchers interested in functional biology and adaptation have used these technologies to sequence mRNA transcriptomes of specific tissues, which in turn are often compared to other tissues, or other individuals with different phenotypes. While these techniques are extremely powerful, careful attention to data quality is required. In particular, because high-throughput sequencing is more error-prone than traditional Sanger sequencing, quality trimming of sequence reads should be an important step in all data processing pipelines. While several software packages for quality trimming exist, no general guidelines for the specifics of trimming have been developed. Here, using empirically derived sequence data, I provide general recommendations regarding the optimal strength of trimming, specifically in mRNA-Seq studies. Although very aggressive quality trimming is common, this study suggests that a more gentle trimming, specifically of those nucleotides whose PHRED score  $<2$  or  $<5$ , is optimal for most studies across a wide variety of metrics.

### 1 Introduction

2 The popularity of genome-enabled biology has increased dramatically, particularly for researchers  
3 studying non-model organisms, over the last few years. For many, the primary goal of these works is to  
4 better understand the genomic underpinnings of adaptive (Linnen et al., 2013; Narum et al., 2013) or  
5 functional (Hsu et al., 2012; Muñoz-Mérida et al., 2013) traits. While extremely promising, the study  
6 of functional genomics in non-model organisms typically requires the generation of a reference  
7 transcriptome to which comparisons are made. Although compared to genome assembly (Bradnam

8 et al., 2013; Earl et al., 2011). transcriptome assembly is less challenging, significant computational  
9 hurdles still exist. Amongst the most difficult of challenges involves the reconstruction of isoforms  
10 (Pyrkosz et al., 2013) and simultaneous assembly of transcripts where read coverage (=expression)  
11 varies by orders of magnitude.

12 These processes are further complicated by the error-prone nature of high-throughput sequencing  
13 reads. With regards to Illumina sequencing, error is distributed non-randomly over the length of the  
14 read, with the rate of error increasing from 5' to 3' end (Liu et al., 2012). These errors are  
15 overwhelmingly substitution errors (Yang et al., 2013), with the global error rate being between 1%  
16 and 3%. Although *de Bruijn* graph assemblers do a remarkable job in distinguishing error from correct  
17 sequence, sequence error does result in assembly error (MacManes and Eisen, 2013). While this type  
18 of error is problematic for all studies, it may be particularly troublesome for SNP-based population  
19 genetic studies. In addition to the biological concerns, sequencing read error may result in problems  
20 of a more technical importance. Because most transcriptome assemblers use a *de Bruijn* graph  
21 representation of sequence connectedness, sequencing error can dramatically increase the size and  
22 complexity of the graph, and thus increase both RAM requirements and runtime.

23 In addition to sequence error correction, which has been shown to improve accuracy of the *de novo*  
24 assembly (MacManes and Eisen, 2013), low quality (=high probability of error) nucleotides are  
25 commonly removed from the sequencing reads prior to assembly, using one of several available tools  
26 (TRIMMOMATIC (Lohse et al., 2012), FASTX TOOLKIT  
27 ([http://hannonlab.cshl.edu/fastx\\_toolkit/index.html](http://hannonlab.cshl.edu/fastx_toolkit/index.html)), BIOPIECES  
28 (<http://www.biopieces.org/>), SOLEXAQA (Cox et al., 2010)). These tools typically use a sliding  
29 window approach, discarding nucleotides falling below a given (user selected) average quality  
30 threshold. The trimmed sequencing read dataset that remains will undoubtedly contain error, though  
31 the absolute number will surely be decreased.

32 Although the process of nucleotide quality trimming is commonplace in HTS analysis pipelines, it's  
33 optimal implementation has not been well defined. Though the rigor with which trimming is  
34 performed may be guided by the design of the experiment, a deeper understanding of the effects of  
35 trimming is desirable. As transcriptome-based studies of functional genomics continue to become more

36 popular, understanding how quality trimming of mRNA-seq reads used in these types of experiments is  
37 urgently needed. Researchers currently working in these field appear to favor aggressive trimming (e.g.  
38 (Looso et al., 2013; Riesgo et al., 2012)), but this may not be optimal. Indeed, one can easily image  
39 aggressive trimming resulting in the removal of a large amount of high quality data (Even nucleotides  
40 removed with the commonly used  $P_{\text{HRED}}=20$  threshold are accurate 99% of the time), just as  
41 lackadaisical trimming (or no trimming) may result in nucleotide errors being incorporated into the  
42 assembled transcriptome.

43 Here, I attempt to provide recommendations regarding the efficient trimming of high-throughput  
44 sequence reads, specifically for mRNASeq reads from the Illumina platform. To do this, I used a  
45 publicly available dataset containing Illumina reads derived from *Mus musculus*. Subsets of these data  
46 (10 million, 20 million, 50 million, 75 million, 100 million reads) were randomly chosen, trimmed to  
47 various levels of stringency, assembled then analyzed for assembly error and content These results aim  
48 to guide researchers through this critical aspect of the analysis of high-throughput sequence data.  
49 While the results of this paper may not be applicable to all studies, that so many researchers are  
50 interested in the genomics of adaptation and phenotypic diversity suggests its widespread utility.

## 51 **Materials and Methods**

52 Because I was interested in understanding the effects of sequence read quality trimming on the  
53 assembly of vertebrate transcriptome assembly, I elected analyze a publicly available (SRR797058)  
54 paired-end Illumina read dataset. This dataset is fully described in a previous publication (Han et al.,  
55 2013), and contains 232 million paired-end 100nt Illumina reads. To investigate how sequencing depth  
56 influences the choice of trimming level, reads data were randomly subsetted into 10 million, 20 million,  
57 50 million, 75 million, 100 million read datasets.

58 Read datasets were trimmed at varying quality thresholds using the software package `TRIMMOMATIC`  
59 (Lohse et al., 2012), which was selected as it appears to be amongst the most popular of read  
60 trimming tools. Specifically, sequences were trimmed at both 5' and 3' ends using  $P_{\text{HRED}} = 0$   
61 (adapter trimming only),  $\leq 2$ ,  $\leq 5$ ,  $\leq 10$ , and  $\leq 20$ . Transcriptome assemblies were generated for  
62 each dataset using the default settings of the program `TRINITY` (Grabherr et al., 2011; Haas et al.,

63 2013). Assemblies were evaluated using a variety of different metrics, many of them comparing  
64 assemblies to the complete collection of *Mus* cDNA's, available at  
65 <http://useast.ensembl.org/info/data/ftp/index.html>.

66 Quality trimming may have substantial effect on assembly quality, and as such, I sought to identify  
67 high quality transcriptome assemblies. Assemblies with few nucleotide errors relative to a known  
68 reference may indicate high quality. The program BLAT (Kent, 2002) was used to identify and count  
69 nucleotide mismatches between reconstructed transcripts and their corresponding reference. To  
70 eliminate spurious short matches between query and template inflating estimates of error, only unique  
71 transcripts that covered more than 90% of their reference sequence were used. Another potential  
72 assessment of assembly quality may be related to the number of paired-end sequencing reads that  
73 concordantly map to the assembly. As the number of reads concordantly mapping increased, so does  
74 assembly quality. To characterize this, I mapped raw (adapter trimmed) sequencing reads to each  
75 assembly using Bowtie2 (Trapnell et al., 2010).

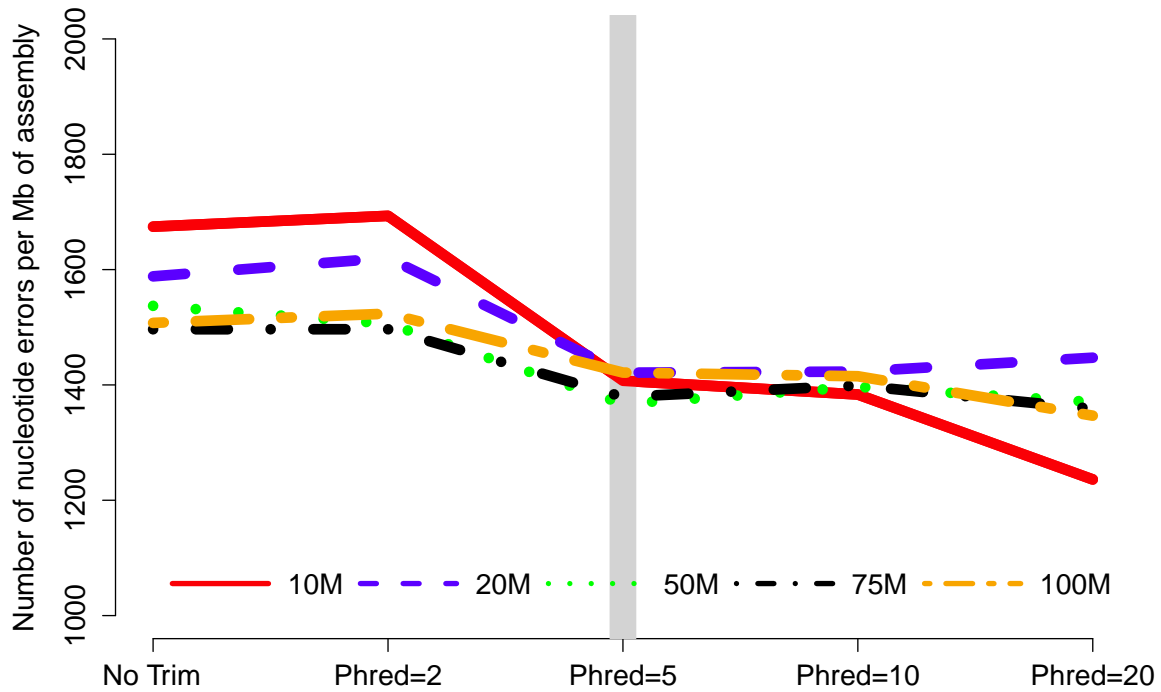
76 Aside from these metrics, measures of assembly content were also assayed. Here, open reading frames  
77 (ORFs) were identified using the program TRANSDCODER  
78 (<http://transdecoder.sourceforge.net/>), and were subsequently translated into amino acid  
79 sequences. The larger the number of complete open reading frames (containing both start and stop  
80 codons) the better the assembly. Lastly, unique transcripts were identified using the blastP program  
81 within the BLAST+ package (Camacho et al., 2009). Blastp hits were retained only if the sequence  
82 similarity was >80% over at least 100 amino acids. As the number of transcripts matching a given  
83 reference increases, so may assembly quality. Code for performing the subsetting, trimming, assembly,  
84 peptide and ORF prediction and blast analyses can be found in the following Github folder  
85 [https://github.com/macmanes/trimming\\_paper/tree/master/scripts](https://github.com/macmanes/trimming_paper/tree/master/scripts).

## 86 Results

87 Quality trimming of sequence reads had a relatively large on the total number of errors contained in  
88 the final assembly (Figure 1), which was reduced by between 9 and 26% when comparing the  
89 assemblies of untrimmed versus PHRED=20 trimmed sequence reads. Most of the improvement in

90 accuracy is gained when trimming at the level of  $P_{\text{HRED}}=5$  or greater, with modest improvements  
 91 potentially garnered with more aggressive trimming at certain coverage levels (Table 1).

92 **Figure 1**

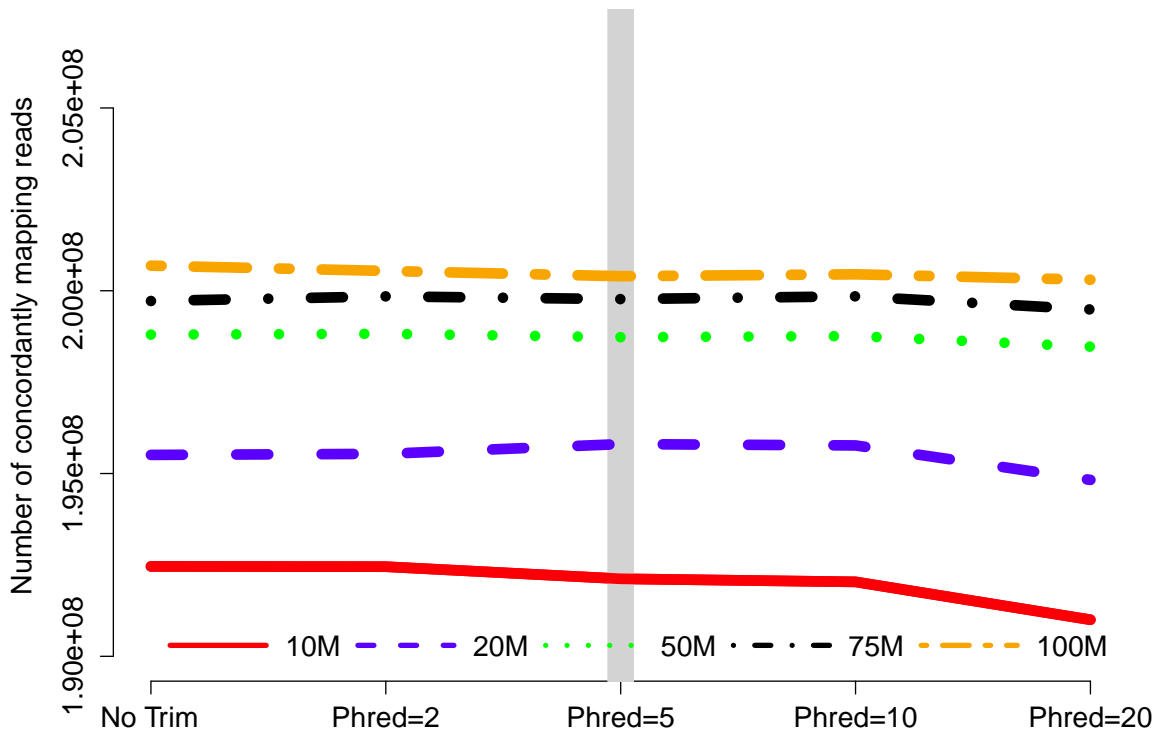


93 Figure 1. The number of nucleotide errors contained in the final transcriptome assembly,  
 94 normalized to assembly size, is related to the strength of quality trimming (Trimming of nucleotides  
 95 whose error scores are:  $P_{\text{HRED}} > 20$ , 10, 5, 2, or no trimming, though most benefits are observed  
 96 at a modest level of trimming. This patterns is largely unchanged with varying depth of sequencing  
 97 coverage (10 million to 100 million sequencing reads). Trimming at  $P_{\text{HRED}} = 5$  may be optimal,  
 98 given the potential untoward effects of more stringent quality trimming.

99 In addition to looking at nucleotide errors, assembly quality may be measured by the the proportion of  
 100 sequencing reads that map concordantly to a given transcriptome assembly (Hunt et al., 2013). As  
 101 such, the analysis of assembly quality includes study of the mapping rates. Here, we found small but  
 102 significant effects of trimming. Specifically, assembling with aggressively quality trimmed reads  
 103 decreased the proportion of reads that map concordantly to a given contig (Figure 2). The pattern is

104 particularly salient with trimming at the  $PHRED = 20$  level. Here, several hundred thousand fewer  
 105 reads mapped compared to mapping against the assembly of untrimmed reads.

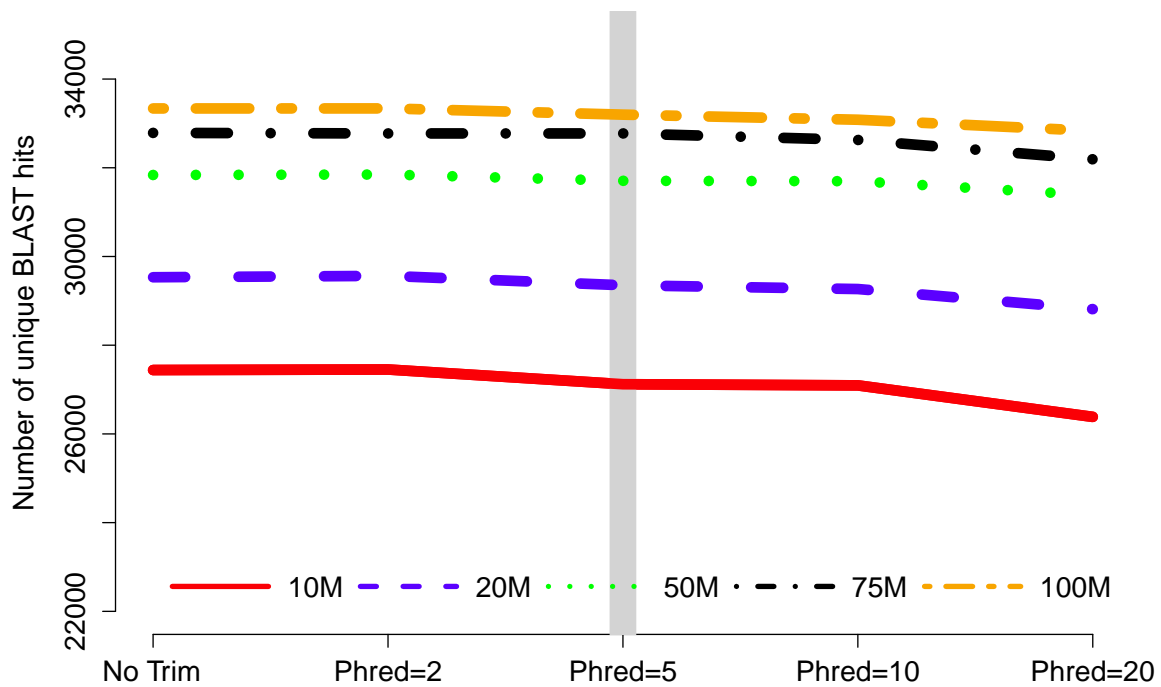
106 **Figure 2**



107 Figure 2. The number of concordantly mapping reads was reduced by trimming. The pattern is  
 108 particularly salient with trimming at  $PHRED=20$  which was always associated with the successful  
 109 mapping of hundreds of thousands of fewer reads.

110 Analysis of assembly content painted a similar picture, with trimming having a relatively small, though  
 111 tangible effect. The number of BLAST+ matches decreased with stringent trimming (Figure 3), with  
 112 trimming at  $PHRED=20$  associated with particularly poor performance.

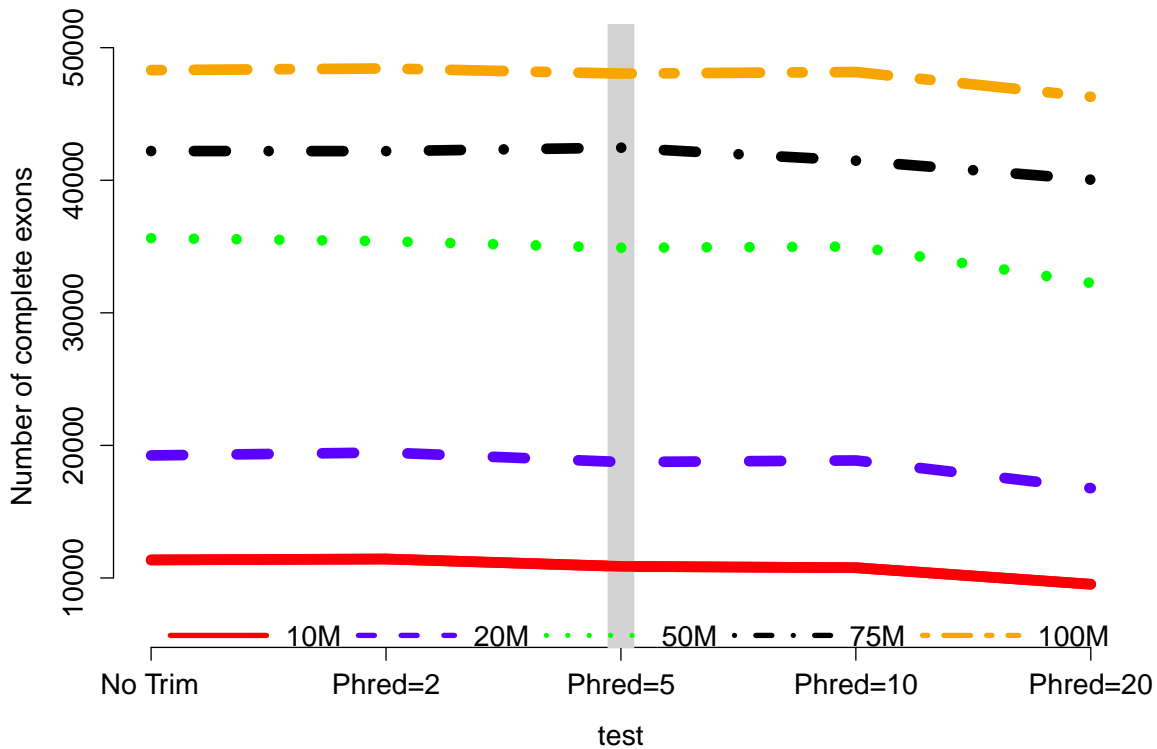
113 **Figure 3**



114 Figure 3. The number of unique BLAST matches contained in the final transcriptome assembly is  
 115 related to the strength of quality trimming for any of the studied sequencing depths. A gentle  
 116 trimming strategy typically yielded the most number of unique matches, while trimming at  
 117 PHRED=20 was always associated with much poorer assembly content

118 When counting complete open reading frames, low and moderate coverage datasets (10M, 20M, 50M)  
 119 were all worsened by aggressive trimming (Figure 4). Trimming at PHRED=20 was the most poorly  
 120 performing level at all read depths.

121 **Figure 4**



122 Figure 4. The number of complete exons contained in the final transcriptome assembly is not  
 123 strongly related to the strength of quality trimming for any of the studies sequencing depths,  
 124 though trimming at PHRED=20 was always associated with fewer identified exons.

125 Of note, all assembly files will be deposited in Dryad upon acceptance for publication. Until then, they  
 126 can be accessed via <https://www.dropbox.com/sh/oiem0v5jgr5c5ir/TYQdGcpYwP>

## 127 Discussion

128 Although the process of nucleotide quality trimming is commonplace in HTS analysis pipelines, it's  
 129 optimal implementation has not been well defined. Though the rigor with which trimming is performed  
 130 seems to vary, there seems to be a bias towards stringent trimming (Ansell et al., 2013; Barrett and  
 131 Davis, 2012; Straub et al., 2013; Tao et al., 2013). This study provides strong evidence that stringent  
 132 quality trimming of nucleotides whose quality scores are  $\leq 20$  results in a poorer transcriptome  
 133 assembly across the majority metrics. Instead, researchers interested in assembling transcriptomes *de*  
 134 *novo* should elect for a much more gentle quality trimming, or no trimming at all. Table 1 summarizes



135 my finding across all experiments, where the numbers represent the trimming level that resulted in the  
 136 most favorable result. What is apparent, is that for typically-sized datasets, trimming at PHRED=2 or  
 137 PHRED=5 optimizes assembly quality. The exception to this rule appears to be in studies where the  
 138 identification of SNP markers from high (or very low) coverage datasets is the primary goal.

139 **Table 1**

DATASET SIZE	ERROR	MAP	ORF	BLAST
10M	20	0	2	2
20M	5	5	2	2
50M	5	5	5	2
75M	20	10	5	0
100M	20	0	2	2

141 Table 1. The PHRED trimming levels that resulted in optimal assemblies across the 4 metrics  
 142 tested in the different size datasets. Error= the number of nucleotide errors in the assembly.  
 143 Map= the number of concordantly mapped reads. ORF= the number of ORFs identified.  
 144 BLAST= the number of unique BLAST hits.

145 The results of this study were surprising. In fact, much of my own work assembling transcriptomes  
 146 included a vigorous trimming step. That trimming had generally small effects, and even negative  
 147 effects when trimming at PHRED=20 was unexpected. To understand if trimming changes the  
 148 distribution of quality scores along the read, we generated plots with the program SolexaQA (Cox  
 149 et al., 2010). Indeed, the program modifies the distribution of PHRED scores in the predicted fashion  
 150 yet downstream effects are minimal. This should be interpreted as speaking to the performance of the  
 151 the bubble popping algorithms included in TRINITY and other *de Bruijn* assemblers.

152 The results presented here stem from the analysis of a single Illumina dataset and specific properties of  
 153 that dataset may have biased the results. This dataset was selected from several evaluated SRA  
 154 datasets for it's 'typical' error profile. The preliminary analysis of a 10 million read subset of another  
 155 typical dataset were concordant with those presented here. Taken together, this suggests that the  
 156 results presented here do not appear to be dependent on the particulars of this dataset, but instead are

157 typical of Illumina mRNAseq datasets.

158 WHAT IS MISSING IN TRIMMED DATASETS? — The question of differences in recovery of specific  
159 contigs is a difficult question to answer. Indeed, these relationships are complex, and could involve a  
160 stochastic process, or be related to differences in expression (low expression transcripts lost in trimmed  
161 datasets) or length (longer contigs lost in trimmed datasets). To investigate this, I attempted to  
162 understand how contigs recovered in the 10 million reads untrimmed dataset but not in the  
163 PHRED=20 trimmed dataset were different. Using the information on FPKM and length generated by  
164 the program EXPRESS, it was clear that the transcripts unique to the untrimmed dataset were more  
165 lowly expressed (mean FPKM=3.2) when compared to the entire untrimmed dataset (mean  
166 FPKM=11.1;  $t = -2.2255$ ,  $df = 70773$ ,  $p\text{-value} = 0.02605$ ). Of note, a similar result was found when  
167 using the non-parametric Wilcoxon test ( $W = 18591566$ ,  $p\text{-value} = 7.184e-13$ ).

168 Turning my attention to length, when comparing uniquely recovered transcripts to the entire  
169 untrimmed dataset of 10 million reads, it appears to be the shorter contigs (mean length 857nt versus  
170 954nt;  $t = -2.1285$ ,  $df = 650.05$ ,  $p\text{-value} = 0.03367$ ,  $W = 26790212$ ,  $p\text{-value} < 2.2e-16$ ) that are  
171 differentially recovered in the untrimmed dataset relative to the PHRED=20 trimmed dataset.

172 EFFECTS OF COVERAGE — Though the experiment was not designed to evaluate the effects of  
173 sequencing depth on assembly, the data speak well to this issue. Contrary to other studies, suggesting  
174 that 30 million paired end reads were sufficient to cover eukaryote transcriptomes (Francis et al.,  
175 2013), the results of the current study suggest that assembly content was more complete as  
176 sequencing depth increased; a pattern that holds at all trimming levels. Though the suggested 30  
177 million read depth was not included in this study, all metrics, including the number of assembly errors  
178 was dramatically reduced, and the number of exons, and BLAST hits were increased as read depth  
179 increased. While generating more sequence data is expensive, given the assembled transcriptome  
180 reference often forms the core of future studies, this investment may be warranted.

181 In summary, the process of nucleotide quality trimming is commonplace in HTS analysis pipelines, but  
182 it's optimal implementation has not been well defined. A very aggressive strategy, where sequence  
183 reads are trimmed when PHRED scores fall below 20 is common. My analyses suggest that for studies

184 whose primary goal is transcript discovery, that a more gentle trimming strategy (e.g. PHRED=2 or  
 185 PHRED=5) that removes only the lowest quality bases is optimal. In particular, it appears as if the  
 186 shorter and more lowly expressed transcripts are particularly vulnerable to loss in studies involving  
 187 more harsh trimming. The one potential exception to this general recommendation may be in studies  
 188 of population genomics, where deep sequencing is leveraged to identify SNPs. Here, a more stringent  
 189 trimming strategy may be warranted.

## 190 Acknowledgments

## 191 References

- 192 Ansell, B.R.E., Schnyder, M., Deplazes, P., Korhonen, P.K., Young, N.D., Hall, R.S., Mangiola, S.,  
 193 Boag, P.R., Hofmann, a., Sternberg, P.W., Jex, A.R., Gasser, R.B., 2013. Insights into the  
 194 immuno-molecular biology of *Angiostrongylus vasorum* through transcriptomics-Prospects for new  
 195 interventions. *Biotechnology Advances* .
- 196 Barrett, C.F., Davis, J.I., 2012. The plastid genome of the mycoheterotrophic *Corallorhiza striata*  
 197 (Orchidaceae) is in the relatively early stages of degradation. *American Journal of Botany* 99,  
 198 1513–1523.
- 199 Bradnam, K.R., Fass, J.N., Alexandrov, A., Baranay, P., Bechner, M., Birol, I., Boisvert, S., Chapman,  
 200 J.A., Chapuis, G., Chikhi, R., Chitsaz, H., Chou, W.C., Corbeil, J., Del Fabbro, C., Docking, T.R.,  
 201 Durbin, R., Earl, D., Emrich, S., Fedotov, P., Fonseca, N.A., Ganapathy, G., Gibbs, R.A., Gnerre,  
 202 S., Godzaridis, E., Goldstein, S., Haimel, M., Hall, G., Haussler, D., Hiatt, J.B., Ho, I.Y., Howard,  
 203 J., Hunt, M., Jackman, S.D., Jaffe, D.B., Jarvis, E., Jiang, H., Kazakov, S., Kersey, P.J., Kitzman,  
 204 J.O., Knight, J.R., Koren, S., Lam, T.W., Lavenier, D., Laviolette, F., Li, Y., Li, Z., Liu, B., Liu, Y.,  
 205 Luo, R., Maccallum, I., MacManes, M.D., Maillet, N., Melnikov, S., Naquin, D., Ning, Z., Otto,  
 206 T.D., Paten, B., Paulo, O.S., Phillippy, A.M., Pina-Martins, F., Place, M., Przybylski, D., Qin, X.,  
 207 Qu, C., Ribeiro, F.J., Richards, S., Rokhsar, D.S., Ruby, J.G., Scalabrin, S., Schatz, M.C., Schwartz,  
 208 D.C., Sergushichev, A., Sharpe, T., Shaw, T.I., Shendure, J., Shi, Y., Simpson, J.T., Song, H.,  
 209 Tsarev, F., Vezzi, F., Vicedomini, R., Vieira, B.M., Wang, J., Worley, K.C., Yin, S., Yiu, S.M.,  
 210 Yuan, J., Zhang, G., Zhang, H., Zhou, S., Korf, I.F., 2013. Assemblathon 2: evaluating *de novo*  
 211 methods of genome assembly in three vertebrate species. *GigaScience* 2, 10.
- 212 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., Madden, T.L., 2009.  
 213 BLAST+: architecture and applications. *BMC Bioinformatics* 10, 421.
- 214 Cox, M.P., Peterson, D.A., Biggs, P.J., 2010. SolexaQA: At-a-glance quality assessment of Illumina  
 215 second-generation sequencing data. *BMC Bioinformatics* 11, 485.
- 216 Earl, D., Bradnam, K., St John, J., Darling, A., Lin, D., Fass, J., Yu, H.O.K., Buffalo, V., Zerbino,  
 217 D.R., Diekhans, M., Nguyen, N., Ariyaratne, P.N., Sung, W.K., Ning, Z., Haimel, M., Simpson,  
 218 J.T., Fonseca, N.A., Birol, I., Docking, T.R., Ho, I.Y., Rokhsar, D.S., Chikhi, R., Lavenier, D.,

- 219 Chapuis, G., Naquin, D., Maillet, N., Schatz, M.C., Kelley, D.R., Phillippy, A.M., Koren, S., Yang,  
 220 S.P., Wu, W., Chou, W.C., Srivastava, A., Shaw, T.I., Ruby, J.G., Skewes-Cox, P., Betegon, M.,  
 221 Dimon, M.T., Solovyev, V., Seledtsov, I., Kosarev, P., Vorobyev, D., Ramirez-Gonzalez, R., Leggett,  
 222 R., MacLean, D., Xia, F., Luo, R., Li, Z., Xie, Y., Liu, B., Gnerre, S., Maccallum, I., Przybylski, D.,  
 223 Ribeiro, F.J., Yin, S., Sharpe, T., Hall, G., Kersey, P.J., Durbin, R., Jackman, S.D., Chapman, J.A.,  
 224 Huang, X., Derisi, J.L., Caccamo, M., Li, Y., Jaffe, D.B., Green, R.E., Haussler, D., Korf, I., Paten,  
 225 B., 2011. Assemblathon 1: a competitive assessment of *de novo* short read assembly methods.  
 226 *Genome Research* 21, 2224–2241.
- 227 Francis, W.R., Christianson, L.M., Kiko, R., Powers, M.L., Shaner, N.C., D Haddock, S.H., 2013. A  
 228 comparison across non-model animals suggests an optimal sequencing depth for *de novo*  
 229 transcriptome assembly. *BMC Genomics* 14, 167.
- 230 Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L.,  
 231 Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, a., Rhind, N., di Palma,  
 232 F., Birren, B.W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., Regev, A., 2011. Full-length  
 233 transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology* 29,  
 234 644–652.
- 235 Haas, B.J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P.D., Bowden, J., Couger, M.B.,  
 236 Eccles, D., Li, B., Lieber, M., MacManes, M.D., Ott, M., Orvis, J., Pochet, N., Strozzi, F., Weeks,  
 237 N., Westerman, R., William, T., Dewey, C.N., Henschel, R., Leduc, R.D., Friedman, N., Regev, A.,  
 238 2013. *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for  
 239 reference generation and analysis. *Nature protocols* 8, 1494–1512.
- 240 Han, H., Irimia, M., Ross, P.J., Sung, H.K., Alipanahi, B., David, L., Golipour, A., Gabut, M.,  
 241 Michael, I.P., Nachman, E.N., Wang, E., Trcka, D., Thompson, T., O'Hanlon, D., Slobodeniuc, V.,  
 242 Barbosa-Morais, N.L., Burge, C.B., Moffat, J., Frey, B.J., Nagy, a., Ellis, J., Wrana, J.L., Blencowe,  
 243 B.J., 2013. MBNL proteins repress ES-cell-specific alternative splicing and reprogramming. *Nature*  
 244 498, 241–245.
- 245 Hsu, J.C., Chien, T.Y., Hu, C.C., Chen, M.J.M., Wu, W.J., Feng, H.T., Haymer, D.S., Chen, C.Y.,  
 246 2012. Discovery of genes related to insecticide resistance in *Bactrocera dorsalis* by functional  
 247 genomic analysis of a *de novo* assembled transcriptome. *PLOS one* 7, e40950.
- 248 Hunt, M., Kikuchi, T., Sanders, M., Newbold, C., Berriman, M., Otto, T.D., 2013. REAPR: a  
 249 universal tool for genome assembly evaluation. *Genome Biology* 14, R47.
- 250 Kent, W.J., 2002. BLAT—the BLAST-like alignment tool. *Genome Research* 12, 656–664.
- 251 Linnen, C.R., Poh, Y.P., Peterson, B.K., Barrett, R.D.H., Larson, J.G., Jensen, J.D., Hoekstra, H.E.,  
 252 2013. Adaptive evolution of multiple traits through multiple mutations at a single gene. *Science*  
 253 (New York, NY) 339, 1312–1316.
- 254 Liu, B., Yuan, J., Yiu, S.M., Li, Z., Xie, Y., Chen, Y., Shi, Y., Zhang, H., Li, Y., Lam, T.W., Luo, R.,  
 255 2012. COPE: an accurate k-mer-based pair-end reads connection tool to facilitate genome assembly.  
 256 *Bioinformatics* (Oxford, England) 28, 2870–2874.
- 257 Lohse, M., Bolger, A.M., Nagel, A., Fernie, A.R., Lunn, J.E., Stitt, M., Usadel, B., 2012. RobiNA: a  
 258 user-friendly, integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids*  
 259 *Research* 40, W622–7.

- 260 Looso, M., Preussner, J., Sousounis, K., Bruckskotten, M., Michel, C.S., Lignelli, E., Reinhardt, R.,  
 261 Höffner, S., Krüger, M., Tsonis, P.A., Borchartd, T., Braun, T., 2013. A *de novo* assembly of the  
 262 newt transcriptome combined with proteomic validation identifies new protein families expressed  
 263 during tissue regeneration. *Genome Biology* 14, R16.
- 264 MacManes, M.D., Eisen, M.B., 2013. Improving transcriptome assembly through error correction of  
 265 high-throughput sequence reads. *PeerJ* 1, e113.
- 266 Muñoz-Mérida, A., González-Plaza, J.J., Cañada, a., Blanco, A.M., García-López, M.d.C., Rodríguez,  
 267 J.M., Pedrola, L., Sicardo, M.D., Hernández, M.L., De la Rosa, R., Belaj, A., Gil-Borja, M., Luque,  
 268 F., Martínez-Rivas, J.M., Pisano, D.G., Trelles, O., Valpuesta, V., Beuzón, C.R., 2013. *De novo*  
 269 assembly and functional annotation of the olive (*Olea europaea*) transcriptome. *DNA Research* 20,  
 270 93–108.
- 271 Narum, S.R., Campbell, N.R., Meyer, K.A., Miller, M.R., Hardy, R.W., 2013. Thermal adaptation and  
 272 acclimation of ectotherms from differing aquatic climates. *Molecular Ecology* 22, 3090–3097.
- 273 Pyrkosz, A.B., Cheng, H., Brown, C.T., 2013. RNA-Seq Mapping Errors When Using Incomplete  
 274 Reference Transcriptomes of Vertebrates. *arXiv.org* [arXiv:1303.2411v1](https://arxiv.org/abs/1303.2411v1).
- 275 Riesgo, A., Perez-Porro, A.R., Carmona, S., Leys, S.P., Giribet, G., 2012. Optimization of preservation  
 276 and storage time of sponge tissues to obtain quality mRNA for next-generation sequencing.  
 277 *Molecular ecology resources* 12, 312–322.
- 278 Straub, S.C.K., Cronn, R.C., Edwards, C., Fishbein, M., Liston, A., 2013. Horizontal transfer of DNA  
 279 from the mitochondrial to the plastid genome and its subsequent evolution in milkweeds  
 280 (*Apocynaceae*). *Genome Biology and Evolution* 5, 1872–1885.
- 281 Tao, T., Zhao, L., Lv, Y., Chen, J., Hu, Y., Zhang, T., Zhou, B., 2013. Transcriptome Sequencing  
 282 and Differential Gene Expression Analysis of Delayed Gland Morphogenesis in *Gossypium australe*  
 283 during Seed Germination. *PLOS one* .
- 284 Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg, S.L.,  
 285 Wold, B.J., Pachter, L., 2010. Transcript assembly and quantification by RNA-Seq reveals  
 286 unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnology* 28,  
 287 511–515.
- 288 Yang, X., Chockalingam, S.P., Aluru, S., 2013. A survey of error-correction methods for  
 289 next-generation sequencing. *Briefings In Bioinformatics* 14, 56–66.