

## 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 Durham, NH 03824*

**3** *Hubbard Center for Genome Studies Durham, NH 03824*

\* 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 afforded researchers the opportunity to gain a deep understanding of genome level processes that underlie evolutionary change, and perhaps more importantly, the links between genotype and phenotype. 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 over the last few years. While  
3 researchers involved in the study of model organisms have had the ability to leverage the power of  
4 genomics for nearly a decade, this power is only now available for the study of non-model organisms.  
5 For many, the primary goal of these newer works is to better understand the genomic underpinnings of  
6 adaptive (Linnen et al., 2013; Narum et al., 2013) or functional (Hsu et al., 2012; Muñoz-Mérida  
7 et al., 2013) traits. While extremely promising, the study of functional genomics in non-model

8 organisms typically requires the generation of a reference transcriptome to which comparisons are  
9 made. Although compared to genome assembly transcriptome assembly is less challenging (Bradnam  
10 et al., 2013; Earl et al., 2011), significant computational hurdles still exist. Amongst the most difficult  
11 of challenges in transcriptome assembly involves the reconstruction of isoforms (Pyrkosz et al., 2013),  
12 simultaneous assembly of transcripts where read coverage (=expression) varies by orders of magnitude,  
13 and overcoming biases related to random hexamer (Hansen et al., 2010) and GC content (Dohm  
14 et al., 2008).

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

26 In addition to sequence error correction, which has been shown to improve accuracy of the *de novo*  
27 assembly (MacManes and Eisen, 2013), low quality (=high probability of error) nucleotides are  
28 commonly removed from the sequencing reads prior to assembly, using one of several available tools  
29 (TRIMMOMATIC (Lohse et al., 2012), FASTX TOOLKIT  
30 ([http://hannonlab.cshl.edu/fastx\\_toolkit/index.html](http://hannonlab.cshl.edu/fastx_toolkit/index.html)), BIOPIECES  
31 (<http://www.biopieces.org/>)). These tools typically use either a sliding window approach,  
32 discarding nucleotides falling below a given (user selected) average quality threshold, or trimming of  
33 low-quality nucleotides at one or both ends of the sequencing read. Though the absolute number will  
34 surely be decreased in the trimmed dataset, aggressive quality trimming may remove a substantial  
35 portion of the total read dataset, which in transcriptome studies may disproportionately effect lower

36 expression transcripts.

37 Although the process of nucleotide quality trimming is commonplace, particularly in the  
38 assembly-based HTS analysis pipelines (e.g. SNP development (Helyar et al., 2012; Milano et al.,  
39 2011), functional studies (Ansell et al., 2013; Bhardwaj et al., 2013), and more general studies of  
40 transcriptome characterization (MacManes and Lacey, 2012; Touming et al., 2013)), it's optimal  
41 implementation has not been well defined. Though the rigor with which trimming is performed may be  
42 guided by the design of the experiment, a deeper understanding of the effects of trimming is desirable.  
43 As transcriptome-based studies of functional genomics continue to become more popular,  
44 understanding how quality trimming of mRNA-seq reads used in these types of experiments is urgently  
45 needed. Researchers currently working in these field appear to favor aggressive trimming (e.g. (Looso  
46 et al., 2013; Riesgo et al., 2012)), but this may not be optimal. Indeed, one can easily image  
47 aggressive trimming resulting in the removal of a large amount of high quality data (even nucleotides  
48 removed with the commonly used  $P_{\text{HRED}}=20$  threshold are accurate 99% of the time), just as  
49 lackadaisical trimming (or no trimming) may result in nucleotide errors being incorporated into the  
50 assembled transcriptome.

51 Here, I provide recommendations regarding the efficient trimming of high-throughput sequence reads,  
52 specifically for mRNASeq reads from the Illumina platform. To do this, I used publicly available  
53 datasets containing Illumina reads derived from *Mus musculus*. Subsets of these data (10 million, 20  
54 million, 50 million, 75 million, 100 million reads) were randomly chosen, trimmed to various levels of  
55 stringency, assembled then analyzed for assembly error and content. In addition to this, I develop a set  
56 of metrics that may be generally useful in evaluating the quality of transcriptome assemblies. These  
57 results aim to guide researchers through this critical aspect of the analysis of high-throughput sequence  
58 data. While the results of this paper may not be applicable to all studies, that so many researchers are  
59 interested in the genomics of adaptation and phenotypic diversity suggests its widespread utility.

## 60 **Materials and Methods**

61 Because I was interested in understanding the effects of sequence read quality trimming on the  
62 assembly of vertebrate transcriptome assembly, I elected to analyze a publicly available (SRR797058)

63 paired-end Illumina read dataset. This dataset is fully described in a previous publication (Han et al.,  
64 2013), and contains 232 million paired-end 100nt Illumina reads. To investigate how sequencing depth  
65 influences the choice of trimming level, reads data were randomly subsetted into 10 million, 20 million,  
66 50 million, 75 million, 100 million read datasets. To test the robustness of my findings, I evaluated  
67 both a second dataset (SRR385624, (Macfarlan et al., 2012)) as well as a technical replicate of the  
68 primary dataset, both at the 10M read dataset size.

69 Read datasets were trimmed at varying quality thresholds using the software package TRIMMOMATIC  
70 version 0.30 (Lohse et al., 2012), which was selected as it appears to be amongst the most popular of  
71 read trimming tools. Specifically, sequences were trimmed at both 5' and 3' ends using PHRED =0  
72 (adapter trimming only),  $\leq 2$ ,  $\leq 5$ ,  $\leq 10$ , and  $\leq 20$ . Other parameters (MINLEN=25,  
73 ILLUMINACLIP=barcodes.fa:2:40:15, SLIDINGWINDOW size=4) were held constant. Transcriptome  
74 assemblies were generated for each dataset using the default settings (except group\_pairs.distance flag  
75 set to 999) of the program TRINITY R2013-02-25 (Grabherr et al., 2011; Haas et al., 2013).  
76 Assemblies were evaluated using a variety of different metrics, many of them comparing assemblies to  
77 the complete collection of *Mus* cDNA's, available at  
78 <http://useast.ensembl.org/info/data/ftp/index.html>.

79 Quality trimming may have substantial effect on assembly quality, and as such, I sought to identify  
80 high quality transcriptome assemblies. Assemblies with few nucleotide errors relative to a known  
81 reference may indicate high quality. The program BLAT v34 (Kent, 2002) was used to identify and  
82 count nucleotide mismatches between reconstructed transcripts and their corresponding reference. To  
83 eliminate spurious short matches between query and template inflating estimates of error, only unique  
84 transcripts that covered more than 90% of their reference sequence were used. Next, because kmers  
85 represent the fundamental unit of assembly, kmers ( $k=25$ ) were counted for each dataset using the  
86 program Jellyfish v1.1.11 (Marçais and Kingsford, 2011). Another potential assessment of assembly  
87 quality may be related to the number of paired-end sequencing reads that concordantly map to the  
88 assembly. As the number of reads concordantly mapping increased, so does assembly quality. To  
89 characterize this, I mapped the full dataset (not subsampled) of adapter trimmed sequencing reads to  
90 each assembly using Bowtie2 v2.1.0 (Trapnell et al., 2010) using default settings, except for maximum

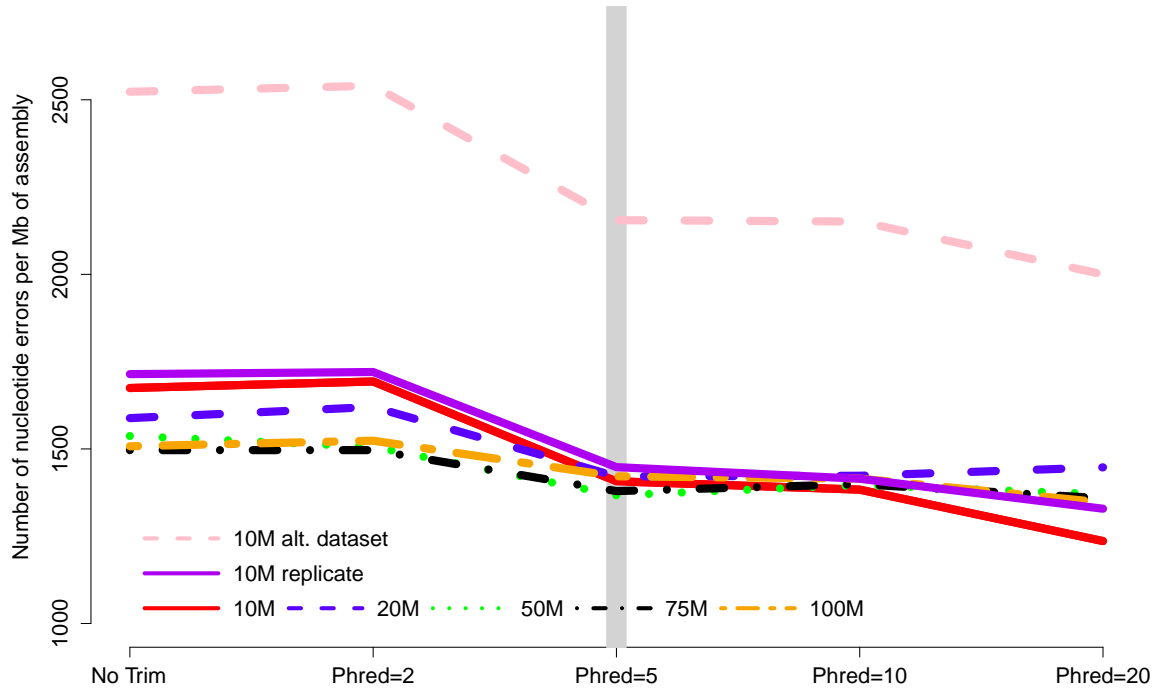
91 insert size (-X 999) and number of multiple mappings (-k 30).

92 Aside from these metrics, measures of assembly content were also assayed. Here, open reading frames  
93 (ORFs) were identified using the default settings of the program TRANSDECODER R20131110  
94 (<http://transdecoder.sourceforge.net/>), and were subsequently translated into amino acid  
95 sequences, both using default settings. The larger the number of complete open reading frames  
96 (containing both start and stop codons) the better the assembly. Next, unique transcripts were  
97 identified using the blastP program within the BLAST+ package version 2.2.28 (Camacho et al.,  
98 2009). Blastp hits were retained only if the sequence similarity was >80% over at least 100 amino  
99 acids, and evaluate  $<10^{-10}$ . As the number of transcripts matching a given reference increases, so may  
100 assembly quality. Lastly, because the effects of trimming may vary with expression, I estimated  
101 expression (e.g. FPKM) for each assembled contig using default settings of the the program EXPRESS  
102 v1.5.0 (Roberts and Pachter, 2013) and the BAM file produced by Bowtie2 as described above. Code  
103 for performing the subsetting, trimming, assembly, peptide and ORF prediction and blast analyses can  
104 be found in the following Github folder  
105 [https://github.com/macmanes/trimming\\_paper/tree/recreate\\_ms\\_analyses/scripts](https://github.com/macmanes/trimming_paper/tree/recreate_ms_analyses/scripts).

## 106 Results

107 Quality trimming of sequence reads had a relatively large on the total number of errors contained in  
108 the final assembly (**Figure 1**), which was reduced by between 9 and 26% when comparing the  
109 assemblies of untrimmed versus PHRED=20 trimmed sequence reads. Most of the improvement in  
110 accuracy is gained when trimming at the level of PHRED=5 or greater, with modest improvements  
111 potentially garnered with more aggressive trimming at certain coverage levels (**Table 1**).

## 112 Figure 1

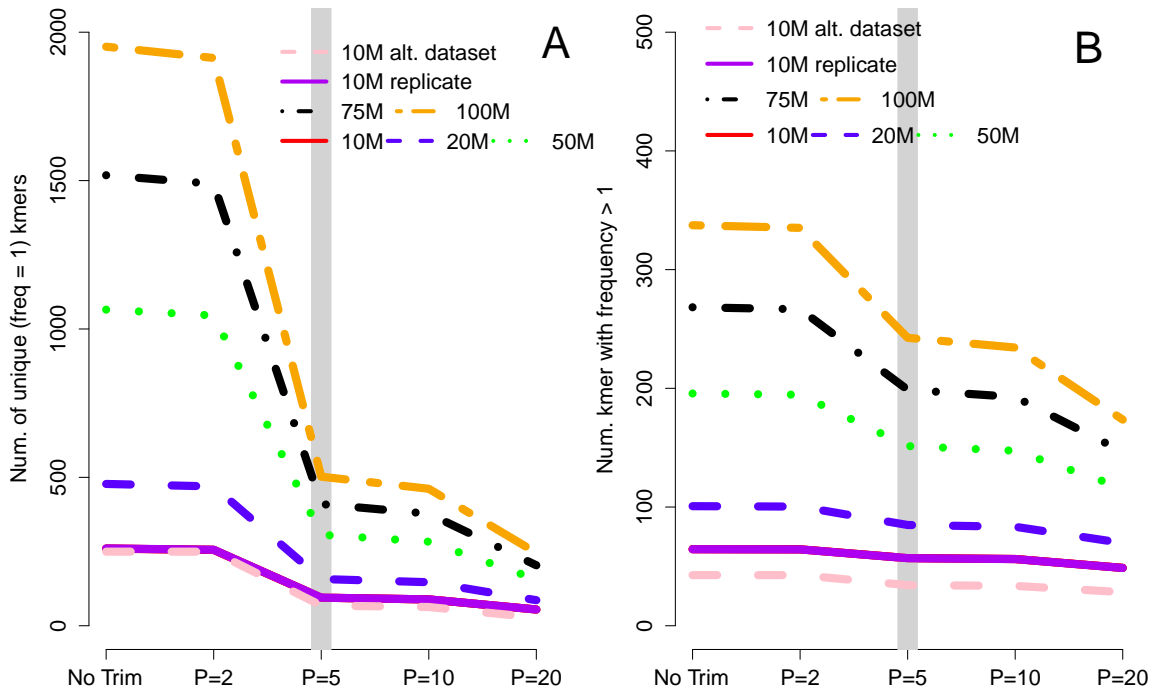


113 Figure 1. The number of nucleotide errors contained in the final transcriptome assembly,  
 114 normalized to assembly size, is related to the strength of quality trimming (Trimming of nucleotides  
 115 whose error scores are: PHRED >20, 10, 5, 2, or no trimming, though most benefits are observed  
 116 at a modest level of trimming. This patterns is largely unchanged with varying depth of sequencing  
 117 coverage (10 million to 100 million sequencing reads). Trimming at PHRED = 5 may be optimal,  
 118 given the potential untoward effects of more stringent quality trimming. 10M, 20M, 50M, 75M,  
 119 100M refer to the subsamples size. 10M replicate is the technical replicate, 10M alt. dataset is the  
 120 secondary dataset. Note that to enhance clarity, the Y-axis does not start at zero.

121 In *de Bruijn* graph-based assemblers, the kmer is the fundamental unit of assembly. Even in  
 122 transcriptome datasets, unique kmers are likely to be formed as a results of sequencing error, and  
 123 therefore may be removed during the trimming process. Figure 1A shows the pattern of unique kmer  
 124 loss across the various trimming levels and read datasets. What is apparent, is that trimming at  
 125 PHRED=5 removes a large fraction of unique kmers, with either less- or more-aggressive trimming  
 126 resulting in smaller effects. In contrast to the removal of unique kmers, those kmers whose frequency  
 127 is >1 are more likely to be real, and therefore should be retained. Figure 1B shows that while

128 PHRED=5 removes unique kmers, it may also reduce the number of non-unique kmers, which many  
 129 hard the assembly process.

130 **Figure 2**

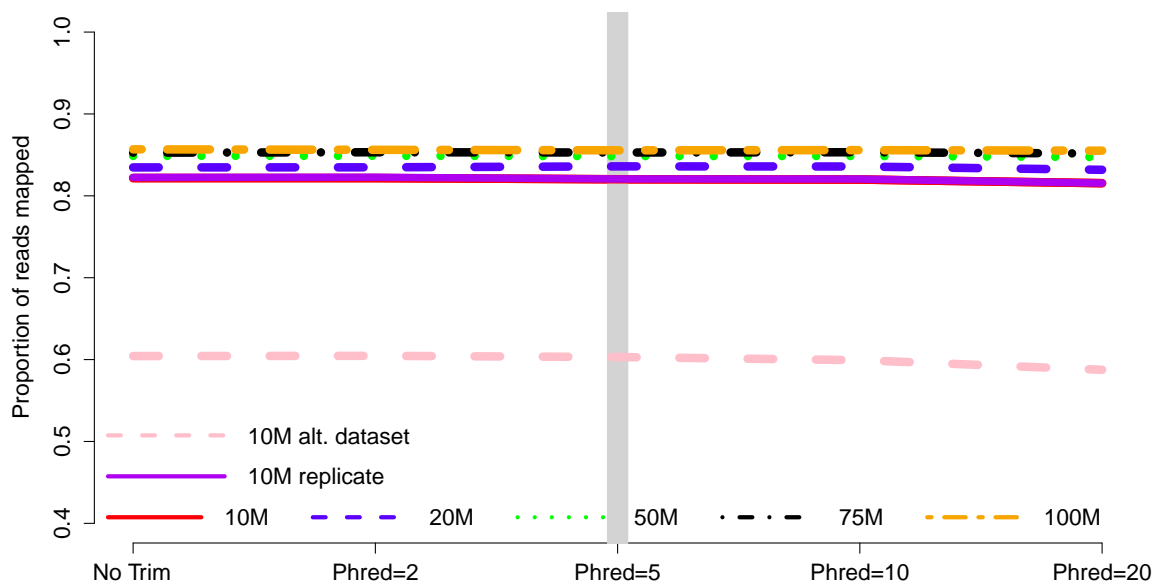


131 Figure 2A. The number of unique kmers removed with various trimming levels across all datasets.  
 132 Trimming at Phred=5 results in a substantial loss of likely erroneous kmers, while the effect of  
 133 more and less aggressive trimming is more diminished. 2B depicts the relationship between  
 134 trimming and non-unique kmers, whose pattern is similar to that of unique kmers.

135 In addition to looking at nucleotide error and kmer distributions, assembly quality may be measured by  
 136 the the proportion of sequencing reads that map concordantly to a given transcriptome assembly  
 137 (Hunt et al., 2013). As such, the analysis of assembly quality includes study of the mapping rates.  
 138 Here, I found small but important effects of trimming. Specifically, assembling with aggressively  
 139 quality trimmed reads decreased the proportion of reads that map concordantly to a given contig  
 140 (Figure 3). Though the patterns are not visually striking, mapping to an assembly of aggressively  
 141 trimmed reads results in several hundred thousand fewer reads mapped compared to mapping against

142 the assembly of less aggressively trimmed reads.

143 **Figure 3**

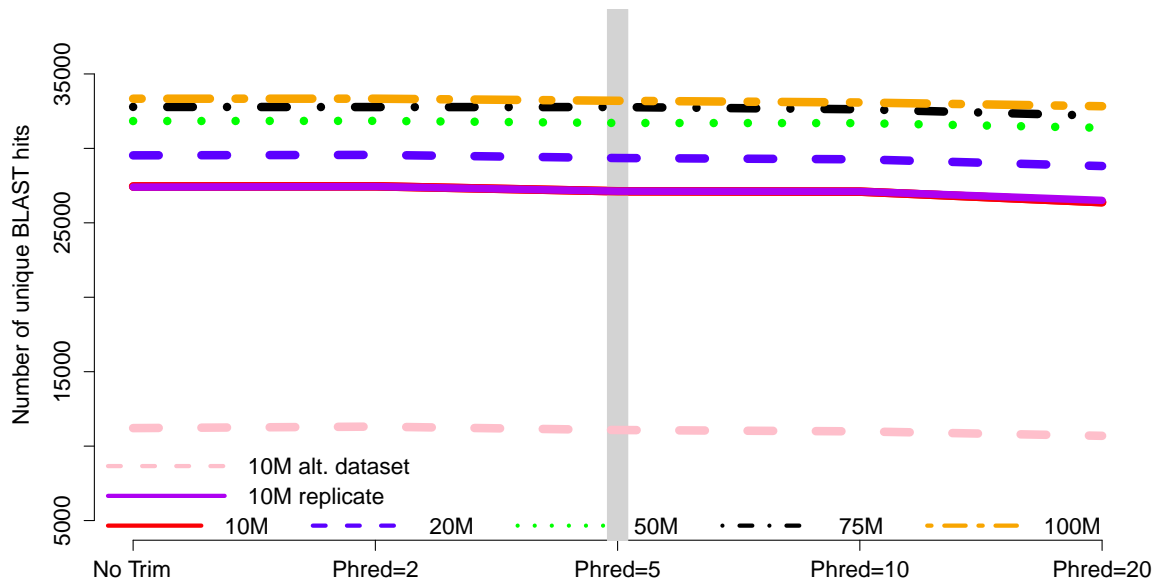


144 Figure 3. The proportion of concordantly mapping reads was reduced by trimming. The pattern is  
 145 particularly salient with trimming at  $PHRED=20$  which was always associated with the successful  
 146 mapping of hundreds of thousands of fewer reads. 10M, 20M, 50M, 75M, 100M refer to the  
 147 subsamples size. 10M replicate is the technical replicate, 10M alt. dataset is the secondary dataset.  
 148 Note that to enhance clarity, the Y-axis does not start at zero.

149 Analysis of assembly content painted a similar picture, with trimming having a relatively small, though  
 150 tangible effect. The number of BLAST+ matches decreased with stringent trimming (Figure 4), with  
 151 trimming at  $PHRED=20$  associated with particularly poor performance.

152 **Figure 4**

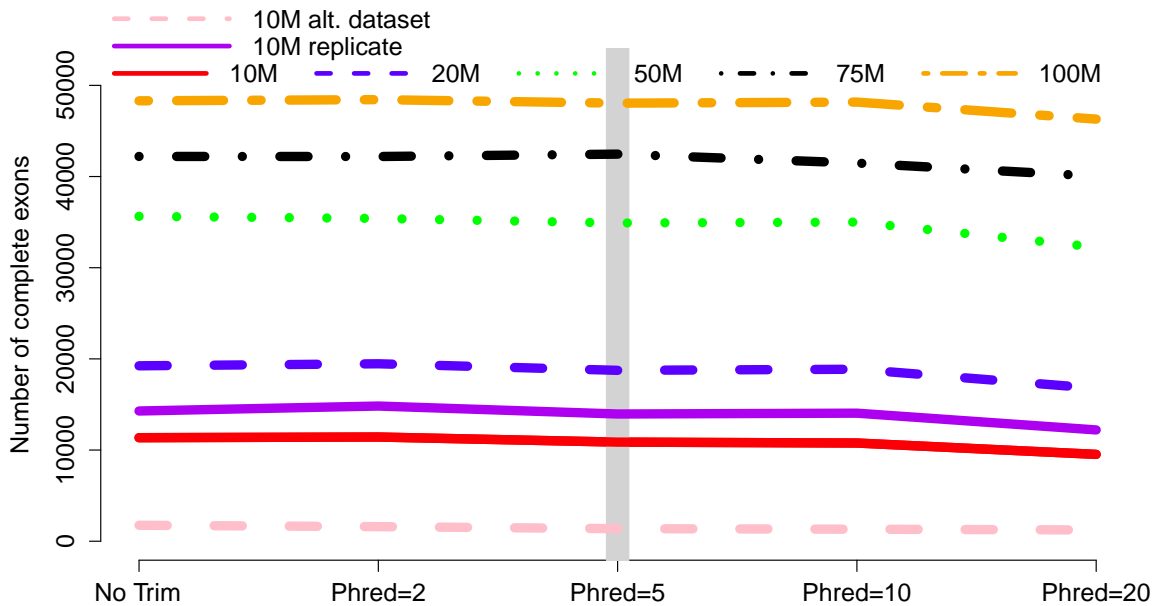




153 Figure 4. The number of unique BLAST matches contained in the final transcriptome assembly is  
 154 related to the strength of quality trimming for any of the studied sequencing depths. A gentle  
 155 trimming strategy typically yielded the most number of unique matches, while trimming at  
 156 PHRED=20 was always associated with much poorer assembly content. 10M, 20M, 50M, 75M,  
 157 100M refer to the subsamples size. 10M replicate is the technical replicate, 10M alt. dataset is the  
 158 secondary dataset. Note that to enhance clarity, the Y-axis does not start at zero.

159 When counting complete open reading frames, low and moderate coverage datasets (10M, 20M, 50M)  
 160 were all worsened by aggressive trimming (Figure 5). Trimming at PHRED=20 was the most poorly  
 161 performing level at all read depths.

162 **Figure 5**



163 Figure 5. The number of complete exons contained in the final transcriptome assembly is not  
 164 strongly related to the strength of quality trimming for any of the studies sequencing depths,  
 165 though trimming at  $\text{PHRED}=20$  was always associated with fewer identified exons. 10M, 20M,  
 166 50M, 75M, 100M refer to the subsamples size. 10M replicate is the technical replicate, 10M alt.  
 167 dataset is the secondary dataset.

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

## 170 Discussion

171 Although the process of nucleotide quality trimming is commonplace in HTS analysis pipelines,  
 172 particularly those involving assembly, its optimal implementation has not been well defined. Though  
 173 the rigor with which trimming is performed seems to vary, there seems to be a bias towards stringent  
 174 trimming (Ansell et al., 2013; Barrett and Davis, 2012; Straub et al., 2013; Tao et al., 2013). This  
 175 study provides strong evidence that stringent quality trimming of nucleotides whose quality scores are  
 176  $\leq 20$  results in a poorer transcriptome assembly across the majority metrics. Instead, researchers  
 177 interested in assembling transcriptomes *de novo* should elect for a much more gentle quality trimming,

178 or no trimming at all. **Table 1** summarizes my finding across all experiments, where the numbers  
 179 represent the trimming level that resulted in the most favorable result. What is apparent, is that for  
 180 typically-sized datasets, trimming at PHRED=2 or PHRED=5 optimizes assembly quality. The  
 181 exception to this rule appears to be in studies where the identification of SNP markers from high (or  
 182 very low) coverage datasets is the primary goal.

183 **Table 1**

	DATASET SIZE	ERROR	MAP	ORF	BLAST
	10M	20	0	2	2
	10M REP.	20	2	2	2
184	10M ALT	20	2	0	0
	20M	5	5	2	2
	50M	5	10	5	2
	75M	20	10	5	0
	100M	20	0	2	2

185 Table 1. The PHRED trimming levels that resulted in optimal assemblies across the 4 metrics  
 186 tested in the different size datasets. Error= the number of nucleotide errors in the assembly.  
 187 Map= the number of concordantly mapped reads. ORF= the number of ORFs identified.  
 188 BLAST= the number of unique BLAST hits. 10M rep. is the technical replicate, 10M alt. is the  
 189 secondary dataset.

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

197 The majority of the results presented here stem from the analysis of a single Illumina dataset and

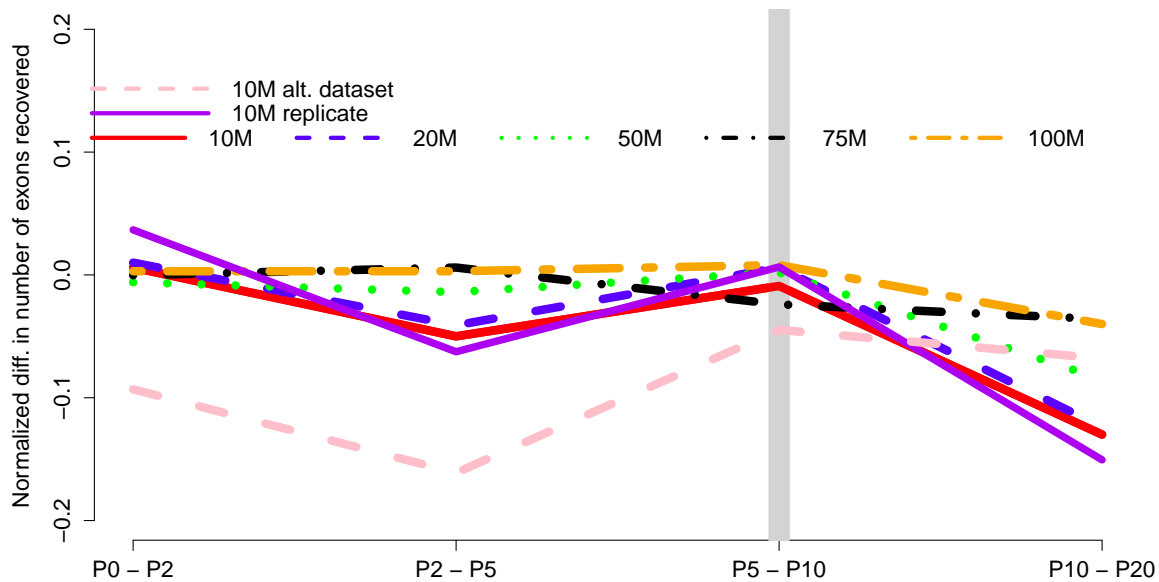
198 specific properties of that dataset may have biased the results. Though the dataset was selected for its  
199 'typical' Illumina error profile, other datasets may produce different results. To evaluate this possibility,  
200 a second dataset was evaluated at the 10M subsampling level. Interestingly, although the assemblies  
201 based on this dataset contained more error (e.g. [Figure 1](#)), aggressive trimming did not improve quality  
202 for any of the assessed metrics, though like other datasets, the absolute number of errors were reduced.

203 In addition to the specific dataset, the subsampling procedure may have resulted in undetected biases.  
204 To address these concerns, a technical replicate of the original dataset was produced at the 10M  
205 subsampling level. This level was selected as a smaller sample of the total dataset is more likely to  
206 contain an unrepresentative sample than larger samples. The results, depicted in all figures as the solid  
207 purple line, are concordant. Therefore, sampling bias is unlikely do drive the patterns reported on here.

208 WHAT IS MISSING IN TRIMMED DATASETS? — The question of differences in recovery of specific  
209 contigs is a difficult question to answer. Indeed, these relationships are complex, and could involve a  
210 stochastic process, or be related to differences in expression (low expression transcripts lost in trimmed  
211 datasets) or length (longer contigs lost in trimmed datasets). To investigate this, I attempted to  
212 understand how contigs recovered in the 10 million reads untrimmed dataset, but not in the  
213 PHRED=20 trimmed dataset were different. Using the information on FPKM and length generated by  
214 the program EXPRESS, it was clear that the transcripts unique to the untrimmed dataset were more  
215 lowly expressed (mean FPKM=3.2) when compared to the entire untrimmed dataset (mean  
216 FPKM=11.1;  $W = 18591566$ ,  $p\text{-value} = 7.184e\text{-}13$ , non-parametric Wilcoxon test).

217 Indeed, I believe that the untoward effects of trimming are linked to a reduction in coverage. For the  
218 datasets tested here, trimming at PHRED=20 resulted in the loss of nearly 25% of the dataset,  
219 regardless of the size of the initial dataset. This relationship does suggest, however, that the  
220 magnitude of the negative effects of trimming should be reduced in larger datasets, and in fact may be  
221 completely erased with ultra-deep sequencing. Indeed, when looking the the differences in the  
222 magnitude of negative effects in the datasets presented here, it is apparent that trimming at  
223 PHRED=20 is 'less bad' in the 100M read dataset than it is in the 10M read dataset ([Figure 6](#)).

224 **Figure 6**



225 Figure 6. The normalized difference in the number of complete exons contained in the  
 226 transcriptome assembly between trimming levels for the five different read datasets. Positive  
 227 numbers indicate increased exons recovered, while negative values indicate decreased recovery. At  
 228 the most aggressive trimming levels, the negative effects are greatest for the smaller datasets, while  
 229 more mitigated for larger datasets. This supports the hypothesis that the relationship between  
 230 assembly quality and trimming may be driven by differences in coverage. P0 - P2 indicates  
 231 differences between Phred=0 and Phred 2 trimming, P2 - P5 are differences between Phred=2 and  
 232 Phred=5 trimming, etc.

233 Turning my attention to length, when comparing uniquely recovered transcripts to the entire  
 234 untrimmed dataset of 10 million reads, it appears to be the shorter contigs (mean length 857nt versus  
 235 954nt;  $W = 26790212$ ,  $p\text{-value} < 2.2e-16$ ) that are differentially recovered in the untrimmed dataset  
 236 relative to the PHRED=20 trimmed dataset.

237 EFFECTS OF COVERAGE ON TRANSCRIPTOME ASSEMBLY — Though the experiment was not  
 238 designed to evaluate the effects of sequencing depth on assembly, the data speak well to this issue.  
 239 Contrary to other studies, suggesting that 30 million paired end reads were sufficient to cover  
 240 eukaryote transcriptomes (Francis et al., 2013), the results of the current study suggest that assembly

241 content was more complete as sequencing depth increased; a pattern that holds at all trimming levels.  
242 Though the suggested 30 million read depth was not included in this study, all metrics, including the  
243 number of assembly errors, as well as the number of exons, and BLAST hits were improved as read  
244 depth increased. While generating more sequence data is expensive, given the assembled  
245 transcriptome reference often forms the core of future studies, this investment may be warranted.

246 SHOULD QUALITY TRIMMING BE REPLACED BY UNIQUE KMER FILTERING? — For transcriptome  
247 studies that revolve around assembly, quality control of sequence data has been thought to be an  
248 crucial step. Though the removal of erroneous nucleotides is the goal, how best to accomplish this is  
249 less clear. As described above, quality trimming has been a common method, but in its commonplace  
250 usage, may be detrimental to assembly. What if, instead of relying on quality scores, we instead rely  
251 on the distribution of kmers to guide our quality control endeavors? In transcriptomes of typical  
252 complexity, sequenced to even moderate coverage, it is reasonable to expect that all but the most  
253 exceptionally rare mRNA molecules are sequenced at a depth  $>1$ . Following this, all kmer whose  
254 frequency is  $<1$  are putative errors, and should be removed before assembly. This idea and its  
255 implementation are fodder for future research.

256 In summary, the process of nucleotide quality trimming is commonplace in many HTS analysis  
257 pipelines, but its optimal implementation has not been well defined. A very aggressive strategy, where  
258 sequence reads are trimmed when PHRED scores fall below 20 is common. My analyses suggest that  
259 for studies whose primary goal is transcript discovery, that a more gentle trimming strategy (*e.g.*  
260 PHRED=2 or PHRED=5) that removes only the lowest quality bases is optimal. In particular, it  
261 appears as if the shorter and more lowly expressed transcripts are particularly vulnerable to loss in  
262 studies involving more harsh trimming. The one potential exception to this general recommendation  
263 may be in studies of population genomics, where deep sequencing is leveraged to identify SNPs. Here,  
264 a more stringent trimming strategy may be warranted.

## 265 Acknowledgments

266 This paper was greatly improved by suggestions of C. Titus Brown and a second anonymous reviewer.  
 267 In addition, the paper was first released as a bioRxiv preprint, and I received several comments based  
 268 on that work both on that website as well as via Twitter. Let it be said here, that early use of a  
 269 preprint archive, open access publication, and Twitter based discussion is a powerful way to rapidly  
 270 disseminate (and get feedback on) work. I highly encourage its use!

## 271 References

- 272 Ansell, B.R.E., Schnyder, M., Deplazes, P., Korhonen, P.K., Young, N.D., Hall, R.S., Mangiola, S.,  
 273 Boag, P.R., Hofmann, a., Sternberg, P.W., Jex, A.R., Gasser, R.B., 2013. Insights into the  
 274 immuno-molecular biology of *Angiostrongylus vasorum* through transcriptomics-Prospects for new  
 275 interventions. *Biotechnology Advances* 31, 1486–1500.
- 276 Barrett, C.F., Davis, J.I., 2012. The plastid genome of the mycoheterotrophic *Corallorhiza striata*  
 277 (Orchidaceae) is in the relatively early stages of degradation. *American Journal of Botany* 99,  
 278 1513–1523.
- 279 Bhardwaj, J., Bhardwaj, J., Chauhan, R., Chauhan, R., Swarnkar, M.K., Swarnkar, M.K., Chahota,  
 280 R.K., Chahota, R.K., Singh, A.K., Singh, A.K., Shankar, R., Shankar, R., Yadav, S.K., Yadav, S.K.,  
 281 2013. Comprehensive transcriptomic study on horse gram (*Macrotyloma uniflorum*): *De novo*  
 282 assembly, functional characterization and comparative analysis in relation to drought stress. *BMC*  
 283 *Genomics* 14, 647.
- 284 Bradnam, K.R., Fass, J.N., Alexandrov, A., Baranay, P., Bechner, M., Birol, I., Boisvert, S., Chapman,  
 285 J.A., Chapuis, G., Chikhi, R., Chitsaz, H., Chou, W.C., Corbeil, J., Del Fabbro, C., Docking, T.R.,  
 286 Durbin, R., Earl, D., Emrich, S., Fedotov, P., Fonseca, N.A., Ganapathy, G., Gibbs, R.A., Gnerre,  
 287 S., Godzaridis, E., Goldstein, S., Haimel, M., Hall, G., Haussler, D., Hiatt, J.B., Ho, I.Y., Howard,  
 288 J., Hunt, M., Jackman, S.D., Jaffe, D.B., Jarvis, E., Jiang, H., Kazakov, S., Kersey, P.J., Kitzman,  
 289 J.O., Knight, J.R., Koren, S., Lam, T.W., Lavenier, D., Laviolette, F., Li, Y., Li, Z., Liu, B., Liu, Y.,  
 290 Luo, R., Maccallum, I., MacManes, M.D., Maillet, N., Melnikov, S., Naquin, D., Ning, Z., Otto,  
 291 T.D., Paten, B., Paulo, O.S., Phillippy, A.M., Pina-Martins, F., Place, M., Przybylski, D., Qin, X.,  
 292 Qu, C., Ribeiro, F.J., Richards, S., Rokhsar, D.S., Ruby, J.G., Scalabrin, S., Schatz, M.C., Schwartz,  
 293 D.C., Sergushichev, A., Sharpe, T., Shaw, T.I., Shendure, J., Shi, Y., Simpson, J.T., Song, H.,  
 294 Tsarev, F., Vezzi, F., Vicedomini, R., Vieira, B.M., Wang, J., Worley, K.C., Yin, S., Yiu, S.M.,  
 295 Yuan, J., Zhang, G., Zhang, H., Zhou, S., Korf, I.F., 2013. Assemblathon 2: evaluating *de novo*  
 296 methods of genome assembly in three vertebrate species. *GigaScience* 2, 10.
- 297 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., Madden, T.L., 2009.  
 298 BLAST+: architecture and applications. *BMC Bioinformatics* 10, 421.
- 299 Cox, M.P., Peterson, D.A., Biggs, P.J., 2010. SolexaQA: At-a-glance quality assessment of Illumina  
 300 second-generation sequencing data. *BMC Bioinformatics* 11, 485.
- 301 Dohm, J.C., Lottaz, C., Borodina, T., Himmelbauer, H., 2008. Substantial biases in ultra-short read  
 302 data sets from high-throughput DNA sequencing. *Nucleic Acids Research* 36, e105–e105.

- 303 Earl, D., Bradnam, K., St John, J., Darling, A., Lin, D., Fass, J., Yu, H.O.K., Buffalo, V., Zerbino,  
304 D.R., Diekhans, M., Nguyen, N., Ariyaratne, P.N., Sung, W.K., Ning, Z., Haimel, M., Simpson,  
305 J.T., Fonseca, N.A., Birol, I., Docking, T.R., Ho, I.Y., Rokhsar, D.S., Chikhi, R., Lavenier, D.,  
306 Chapuis, G., Naquin, D., Maillet, N., Schatz, M.C., Kelley, D.R., Phillippy, A.M., Koren, S., Yang,  
307 S.P., Wu, W., Chou, W.C., Srivastava, A., Shaw, T.I., Ruby, J.G., Skewes-Cox, P., Betegon, M.,  
308 Dimon, M.T., Solovyev, V., Seledtsov, I., Kosarev, P., Vorobyev, D., Ramirez-Gonzalez, R., Leggett,  
309 R., MacLean, D., Xia, F., Luo, R., Li, Z., Xie, Y., Liu, B., Gnerre, S., Maccallum, I., Przybylski, D.,  
310 Ribeiro, F.J., Yin, S., Sharpe, T., Hall, G., Kersey, P.J., Durbin, R., Jackman, S.D., Chapman, J.A.,  
311 Huang, X., Derisi, J.L., Caccamo, M., Li, Y., Jaffe, D.B., Green, R.E., Haussler, D., Korf, I., Paten,  
312 B., 2011. Assemblathon 1: a competitive assessment of *de novo* short read assembly methods.  
313 Genome Research 21, 2224–2241.
- 314 Francis, W.R., Christianson, L.M., Kiko, R., Powers, M.L., Shaner, N.C., D Haddock, S.H., 2013. A  
315 comparison across non-model animals suggests an optimal sequencing depth for *de novo*  
316 transcriptome assembly. BMC Genomics 14, 167.
- 317 Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L.,  
318 Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, a., Rhind, N., di Palma,  
319 F., Birren, B.W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., Regev, A., 2011. Full-length  
320 transcriptome assembly from RNA-Seq data without a reference genome. Nature Biotechnology 29,  
321 644–652.
- 322 Haas, B.J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P.D., Bowden, J., Couger, M.B.,  
323 Eccles, D., Li, B., Lieber, M., MacManes, M.D., Ott, M., Orvis, J., Pochet, N., Strozzi, F., Weeks,  
324 N., Westerman, R., William, T., Dewey, C.N., Henschel, R., Leduc, R.D., Friedman, N., Regev, A.,  
325 2013. *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for  
326 reference generation and analysis. Nature protocols 8, 1494–1512.
- 327 Han, H., Irimia, M., Ross, P.J., Sung, H.K., Alipanahi, B., David, L., Golipour, A., Gabut, M.,  
328 Michael, I.P., Nachman, E.N., Wang, E., Trcka, D., Thompson, T., O'Hanlon, D., Slobodeniuc, V.,  
329 Barbosa-Morais, N.L., Burge, C.B., Moffat, J., Frey, B.J., Nagy, a., Ellis, J., Wrana, J.L., Blencowe,  
330 B.J., 2013. MBNL proteins repress ES-cell-specific alternative splicing and reprogramming. Nature  
331 498, 241–245.
- 332 Hansen, K.D., Brenner, S.E., Dudoit, S., 2010. Biases in Illumina transcriptome sequencing caused by  
333 random hexamer priming. Nucleic Acids Research 38, e131–e131.
- 334 Helyar, S.J., Helyar, S.J., Limborg, M.T., Limborg, M.T., Bekkevold, D., Babbucci, M., Babbucci, M.,  
335 van Houdt, J., van Houdt, J., Maes, G.E., Bargelloni, L., Bargelloni, L., Nielsen, R.O., Nielsen, R.O.,  
336 Taylor, M.I., Taylor, M.I., Ogden, R., Ogden, R., Cariani, A., Cariani, A., Carvalho, G.R., Carvalho,  
337 G.R., Consortium, F., Consortium, F., Panitz, F., 2012. SNP Discovery Using Next Generation  
338 Transcriptomic Sequencing in Atlantic Herring (*Clupea harengus*). PLOS ONE 7, e42089.
- 339 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.,  
340 2012. Discovery of genes related to insecticide resistance in *Bactrocera dorsalis* by functional  
341 genomic analysis of a *de novo* assembled transcriptome. PLOS ONE 7, e40950.
- 342 Hunt, M., Kikuchi, T., Sanders, M., Newbold, C., Berriman, M., Otto, T.D., 2013. REAPR: a  
343 universal tool for genome assembly evaluation. Genome Biology 14, R47.
- 344 Kent, W.J., 2002. BLAT—the BLAST-like alignment tool. Genome Research 12, 656–664.



- 345 Linnen, C.R., Poh, Y.P., Peterson, B.K., Barrett, R.D.H., Larson, J.G., Jensen, J.D., Hoekstra, H.E.,  
 346 2013. Adaptive evolution of multiple traits through multiple mutations at a single gene. *Science*  
 347 (New York, NY) 339, 1312–1316.
- 348 Liu, B., Yuan, J., Yiu, S.M., Li, Z., Xie, Y., Chen, Y., Shi, Y., Zhang, H., Li, Y., Lam, T.W., Luo, R.,  
 349 2012. COPE: an accurate k-mer-based pair-end reads connection tool to facilitate genome assembly.  
 350 *Bioinformatics* (Oxford, England) 28, 2870–2874.
- 351 Lohse, M., Bolger, A.M., Nagel, A., Fernie, A.R., Lunn, J.E., Stitt, M., Usadel, B., 2012. RobiNA: a  
 352 user-friendly, integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids*  
 353 *Research* 40, W622–7.
- 354 Looso, M., Preussner, J., Sousounis, K., Bruckskotten, M., Michel, C.S., Lignelli, E., Reinhardt, R.,  
 355 Höffner, S., Krüger, M., Tsonis, P.A., Borchardt, T., Braun, T., 2013. A *de novo* assembly of the  
 356 newt transcriptome combined with proteomic validation identifies new protein families expressed  
 357 during tissue regeneration. *Genome Biology* 14, R16.
- 358 Macfarlan, T.S., Gifford, W.D., Driscoll, S., Lettieri, K., Rowe, H.M., Bonanomi, D., Firth, A., Singer,  
 359 O., Trono, D., Pfaff, S.L., 2012. Embryonic stem cell potency fluctuates with endogenous retrovirus  
 360 activity. *Nature* 487, 57–63.
- 361 MacManes, M.D., Eisen, M.B., 2013. Improving transcriptome assembly through error correction of  
 362 high-throughput sequence reads. *PeerJ* 1, e113.
- 363 MacManes, M.D., Lacey, E.A., 2012. The Social Brain: Transcriptome Assembly and Characterization  
 364 of the Hippocampus from a Social Subterranean Rodent, the Colonial Tuco-Tuco (*Ctenomys*  
 365 *sociabilis*). *PLOS ONE* 7, e45524.
- 366 Marçais, G., Kingsford, C., 2011. A fast, lock-free approach for efficient parallel counting of  
 367 occurrences of k-mers. *Bioinformatics* (Oxford, England) 27, 764–770.
- 368 Milano, I., Babbucci, M., Panitz, F., Ogden, R., Nielsen, R.O., Taylor, M.I., Helyar, S.J., Carvalho,  
 369 G.R., Espiñeira, M., Atanassova, M., Tinti, F., Maes, G.E., Patarnello, T., FishPopTrace  
 370 Consortium, Bargelloni, L., 2011. Novel tools for conservation genomics: comparing two  
 371 high-throughput approaches for SNP discovery in the transcriptome of the European hake. *PLOS*  
 372 *ONE* 6, e28008.
- 373 Muñoz-Mérida, A., González-Plaza, J.J., Cañada, a., Blanco, A.M., García-López, M.d.C., Rodríguez,  
 374 J.M., Pedrola, L., Sicardo, M.D., Hernández, M.L., De la Rosa, R., Belaj, A., Gil-Borja, M., Luque,  
 375 F., Martínez-Rivas, J.M., Pisano, D.G., Trelles, O., Valpuesta, V., Beuzón, C.R., 2013. *De novo*  
 376 assembly and functional annotation of the olive (*Olea europaea*) transcriptome. *DNA Research* 20,  
 377 93–108.
- 378 Narum, S.R., Campbell, N.R., Meyer, K.A., Miller, M.R., Hardy, R.W., 2013. Thermal adaptation and  
 379 acclimation of ectotherms from differing aquatic climates. *Molecular Ecology* 22, 3090–3097.
- 380 Pyrkosz, A.B., Cheng, H., Brown, C.T., 2013. RNA-Seq Mapping Errors When Using Incomplete  
 381 Reference Transcriptomes of Vertebrates. *arXiv.org* [arXiv:1303.2411v1](https://arxiv.org/abs/1303.2411v1).
- 382 Riesgo, A., Perez-Porro, A.R., Carmona, S., Leys, S.P., Giribet, G., 2012. Optimization of preservation  
 383 and storage time of sponge tissues to obtain quality mRNA for next-generation sequencing.  
 384 *Molecular ecology resources* 12, 312–322.

- 385 Roberts, A., Pachter, L., 2013. Streaming fragment assignment for real-time analysis of sequencing  
386 experiments. *Nature Methods* 10, 71–73.
- 387 Straub, S.C.K., Cronn, R.C., Edwards, C., Fishbein, M., Liston, A., 2013. Horizontal transfer of DNA  
388 from the mitochondrial to the plastid genome and its subsequent evolution in milkweeds  
389 (Apocynaceae). *Genome Biology and Evolution* 5, 1872–1885.
- 390 Tao, T., Zhao, L., Lv, Y., Chen, J., Hu, Y., Zhang, T., Zhou, B., 2013. Transcriptome Sequencing  
391 and Differential Gene Expression Analysis of Delayed Gland Morphogenesis in *Gossypium australe*  
392 during Seed Germination. *PLOS ONE* 8, e75323.
- 393 Touming, L., Touming, L., Siyuan, Z., Siyuan, Z., Qingming, T., Qingming, T., Ping, C., Ping, C.,  
394 Yongting, Y., Yongting, Y., Shouwei, T., Shouwei, T., 2013. *De novo* assembly and characterization  
395 of transcriptome using Illumina paired-end sequencing and identification of CesA gene in ramie  
396 (*Boehmeria nivea* L. Gaud). *BMC Genomics* 14, 125.
- 397 Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg, S.L.,  
398 Wold, B.J., Pachter, L., 2010. Transcript assembly and quantification by RNA-Seq reveals  
399 unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnology* 28,  
400 511–515.
- 401 Yang, X., Chockalingam, S.P., Aluru, S., 2013. A survey of error-correction methods for  
402 next-generation sequencing. *Briefings In Bioinformatics* 14, 56–66.