

# Establishing evidenced-based best practice for the *de novo* assembly and evaluation of transcriptomes from non-model organisms

Matthew D. MacManes<sup>1</sup>,

<sup>1</sup> Department of Molecular, Cellular and Biomedical Sciences, University of New Hampshire, Durham NH, USA

\* E-mail: [macmanes@gmail.com](mailto:macmanes@gmail.com)

⊕ Twitter: @macmanes

## 1 Abstract

2 Characterizing transcriptomes in both model and non-model organisms has resulted in a massive increase in  
3 our understanding of biological phenomena. This boon, largely made possible via high-throughput  
4 sequencing, means that studies of functional, evolutionary and population genomics are now being done by  
5 hundreds or even thousands of labs around the world. For many, these studies begin with a *de novo*  
6 transcriptome assembly, which is a technically complicated process involving several discrete steps. Each  
7 step may be accomplished in one of several different ways, using different software packages, each producing  
8 different results. This analytical complexity begs the question – *Which method(s) are optimal?* Using  
9 reference and non-reference based evaluative methods, I propose a set of guidelines that aim to standardize  
10 and facilitate the process of transcriptome assembly. These recommendations include the generation of  
11 between 20 million and 40 million sequencing reads from single individual where possible, error correction of  
12 reads, gentle quality trimming, assembly filtering using **Transrate** and/or gene expression, annotation using  
13 **dammit**, and appropriate reporting. These recommendations have been extensively benchmarked and  
14 applied to publicly available transcriptomes, resulting in improvements in both content and contiguity. To  
15 facilitate the implementation of the proposed standardized methods, I have released a set of version  
16 controlled open-sourced code, **The Oyster River Protocol for Transcriptome Assembly**, available at  
17 <http://oyster-river-protocol.rtfid.org/>.

## 18 Introduction

19 For all biology, modern sequencing technologies has provided for an unprecedented opportunity to gain a  
20 deep understanding of genome level processes that underlie a very wide array of natural phenomenon, from  
21 intracellular metabolic processes to global patterns of population variability. Transcriptome sequencing has  
22 been influential, particularly in functional genomics, and has resulted in discoveries not possible even just a  
23 few years ago. This in large part is due to the scale at which these studies may be conducted. Unlike  
24 studies of adaptation based on one or a small number of candidate genes (e.g. (Fitzpatrick et al., 2005;  
25 Panhuis, 2006)), modern studies may assay the entire suite of expressed transcripts – the transcriptome –  
26 simultaneously. In addition to issues of scale, as a direct result of enhanced dynamic range, newer sequencing  
27 studies have increased ability to simultaneously reconstruct and quantitate lowly- and highly-expressed  
28 transcripts, (Vijay et al., 2013; Wolf, 2013). Lastly, improved methods for the detection of differences in  
29 gene expression (*e.g.*, (Love et al., 2014; Robinson et al., 2010)) across experimental treatments has resulted  
30 in increased resolution for studies aimed at understanding changes in gene expression.

31 As a direct result of their widespread popularity, a diverse toolset for the assembly and analysis of  
32 transcriptome exists. Notable amongst the wide array of tools include several for quality visualization -  
33 **FastQC** (available here) and **SolexaQA** (Cox et al., 2010), read trimming (e.g. **Skewer** (Jiang et al., 2014),  
34 **Trimmomatic** (Bolger et al., 2014) and **Cutadapt** (Martin, 2011)), read normalization (**khmer** (Pell et al.,  
35 2012)), error correction (Le et al., 2013), assembly (**Trinity** (Haas et al., 2013), **SOAPdenovoTrans** (Xie  
36 et al., 2014)), and assembly verification (**Transrate** (Smith-Unna et al., 2015)), **BUSCO** (**Benchmarking**  
37 **Universal Single-Copy Orthologs** - (Simão et al., 2015)), and **RSEM-eval** (Li et al., 2014)). The ease with  
38 which these tools may be used to produce transcriptome assemblies belies the true complexity underlying  
39 the overall process. Indeed, the subtle (and not so subtle) methodological challenges associated with  
40 transcriptome reconstruction may result in highly variable assembly quality. Amongst the most challenging  
41 include isoform reconstruction and simultaneous assembly of low- and high-coverage transcripts (Johnson  
42 et al., 2003; Modrek et al., 2001), which together make accurate transcriptome assembly technically  
43 challenging. Production of an accurate transcriptome assembly requires a large investment in time and  
44 resources. Each step in it's production requires careful consideration. Here, I propose a set of  
45 evidence-based guidelines for assembly and evaluation that will result in the production of the highest  
46 quality transcriptome assembly possible.

47 Currently, a very large number of labs and research programs depend, often critically, on the production  
48 of accurate transcriptome resources. That no current best practices exists – particularly for those working in

49 non-model systems – has resulted in an untenable situation where each laboratory makes up it's own  
50 computational pipeline. These pipelines, often devoid of rigorous quality evaluation, may have important  
51 downstream consequences. This manuscript, by proposing a specific evidence-based process, significantly  
52 enhances the technical quality and reproducibility of transcriptome studies, which is critical for this  
53 emerging field of research.

## 54 **Methods**

55 To demonstrate the merits of my recommendations, a large number of assemblies were produced using a  
56 variety of methods. For all assemblies performed, Illumina sequencing adapters were removed from both  
57 ends of the sequencing reads, as were nucleotides with quality Phred  $\leq 2$ , using the program **Trimmomatic**  
58 version 0.32 (Bolger et al., 2014). The reads were assembled using Trinity release 2.1.1 (Haas et al., 2013)  
59 using default settings. **Trinity** was used as the default assembler as it has been previously reported to be  
60 best in class (Li et al., 2014; Smith-Unna et al., 2015). Assemblies were characterized using **Transrate**  
61 version 1.0.1 (Smith-Unna et al., 2015). Using this software, I generated three kinds of metrics: contig  
62 metrics; mapping metrics which used as input the same reads that were fed into the assembler for each  
63 assembly; and comparative metrics which used as input the *Mus musculus* version 75 transcriptome. In  
64 addition to the metrics provided by **Transrate**, I evaluated completeness of each assembly by use of BUSCO,  
65 a software package that searches for highly conserved, near-universal, single copy orthologs. All assemblies  
66 generated are available here, and will be moved to Dryad on acceptance.

67 To understand the influence of read depth on assembly quality, I produced subsets of size  
68 1,2,5,10,20,40,60,80,100 million paired end reads of two publicly available paired-end datasets - A *Mus*  
69 dataset -SRR797058 described in (Han et al., 2013) and a human dataset - SRR1659968. The subsampling  
70 procedure was accomplished via the software package **seqtk** (<https://github.com/lh3/seqtk>). For the  
71 evaluation of the effects of sequence polymorphism on assembly quality, I use reads from BioProject  
72 PRJNA157895 described in (MacManes and Lacey, 2012), a *Ctenomys* dataset which consists of 10 read  
73 files from the hypothalami of 10 different individuals. This dataset was assembled two ways. First, the reads  
74 from all 10 individuals were jointly assembled in one large assembly [CODE]. This assembly was compared  
75 to the assembly of a single individual [CODE]. Assemblies were generated and evaluated as per above.

76 To evaluate the effects of error correction, I used the subsampled read datasets, which were subsequently  
77 error corrected using the following software packages: **SEECER** version 0.1.3 (Le et al., 2013), **Lighter**  
78 version 1.0.7 (Song et al., 2014), **SGA** version 0.10.13 (Simpson and Durbin, 2012), **bfc** version r177 (Li,

79 2015), `RCorrector` (Song and Florea, 2015), and `BLESS` version 0.24 (Heo et al., 2014). In correction  
80 algorithms (`SGA`, `BLESS`, `bfc`) that allowed for the use of larger *kmer* lengths, I elected to error correct with  
81 a small ( $k = 31$ ) and a long ( $k = 55$ ) *kmer*, while for the other software (`RCorrector`, `SEECER` and `Lighter`)  
82 that does not allow for longer *kmer* values, I set  $k = 31$ . `bfc` requires interleaved reads, which was  
83 accomplished using `khmer` version 2.0 (Alameldin et al., 2015; Brown et al., 2012; McDonald and Brown,  
84 2013). Code for performing these steps is available [here].

85 The effects of `khmer` digital normalization (Pell et al., 2012) were characterized by generating three 20  
86 million and three 100 million read subsets of the larger *Mus* dataset. Digital normalization was performed  
87 using a median *kmer* abundance threshold of 30. The resulting datasets were assembled using `Trinity`, and  
88 evaluated using `BUSCO` and `Transrate`. Code for performing these steps is available in the `diginorm` target  
89 of the [Makefile].

90 Post-assembly processing was evaluated using several assembly datasets of various sizes, generated above.  
91 Each assembly was evaluated using `Transrate`. `Transrate` produces a score based on contig and mapping  
92 metrics, as well as a more optimal assembly where poorly supported contigs (putative assembly artifacts)  
93 are removed. Both the original and `Transrate` optimal assembly are evaluated using `BUSCO`, to help better  
94 understand if filtration results in the loss of non-artifactual transcripts. In addition to `Transrate` filtration,  
95 an additional, or alternative filtration step is performed using estimates of gene expression  
96 (TPM=transcripts per million). TPM is estimated by two different software packages that implement two  
97 distinct methods - `Salmon` (Patro et al., 2015) and `Kallisto` (Bray et al., 2015). Transcripts whose  
98 expression is estimated to be greater than a given threshold, typically TPM=1 or TPM=0.5 are retained.  
99 As above, the filtered assemblies are evaluated using `BUSCO`, to help better understand if filtration results in  
100 the loss of non-artifactual transcripts. Code for performing these steps is available in the `QC` target of the  
101 makefile available [here].

## 102 Recommendations

### 103 0.1 Input Data

104 **Summary Statement: Sequence 1 or more tissues from 1 individual to a depth of between 20**  
105 **million and 40 million 100bp or longer paired-end reads.**

106 When planning to construct a transcriptome, the first question to ponder is the type and quantity of

107 data required. While this will be somewhat determined by the specific goals of the study and availability of  
108 tissues, there are some general guiding principals. As of 2014, Illumina continues to offer the most flexibility  
109 in terms of throughput, analytical tractability, and cost (Glenn, 2011) and as a result, the recommendations  
110 here are primarily related to assembly using Illumina data. It is worth noting however, that long-read (e.g.  
111 PacBio) transcriptome sequencing is just beginning to emerge as an alternative (Au et al., 2013),  
112 particularly for researchers interested in understanding isoform complexity. Though currently lacking the  
113 throughput for accurate quantitation of gene expression, long read technologies, much like they have done  
114 for *de novo* genome assembly, seem likely to replace short-read-based *de novo* transcriptome assembly at  
115 some point in the future.

116 For the typical transcriptome study, one should plan to generate a reference based on 1 or more tissue  
117 types, with each tissue adding unique tissue-specific transcripts and isoforms. Though increasing the  
118 amount of sequence data collected does increase the accuracy and completeness of the assembly (Figure 1,  
119 3) albeit marginally, a balance between cost and quality exists. For the datasets examined here (vertebrate  
120 tissues), sequencing more than between 20M and 40M paired-end reads is associated with the discovery of  
121 very few additional transcripts, and only minor improvement in other assembly metrics. Read length should  
122 be at least 100bp, with longer reads likely aiding in isoform reconstruction and contiguity (Garber et al.,  
123 2011). In the case where multiple tissues are sequenced, it is likely best to combine reads from each tissue  
124 together to produce a joint assembly.

125 Because sequence polymorphism increases the complexity of the *de bruijn* graph (Iqbal et al., 2012;  
126 Studholme, 2010), and therefore may negatively effect the assembly itself, the reference transcriptome  
127 should be generated from reads corresponding to as homogeneous a sample as possible. For outbred,  
128 non-model organisms, this usually means generating reads from a single individual. When more then one  
129 individual is required to meet other requirements (*e.g.*, for differential expression replicates or experimental  
130 treatment conditions), keeping the number of individuals to a minimum is paramount. For instance, when  
131 performing an experiment where a distinct set of genes may be expressed in different treatments (or sexes),  
132 the recommendation is to sequence one individual from each treatment class.

133 To illustrate this effect, I examined the effects of assembling reads from 10 individuals jointly, versus  
134 assembling a representative individual. This individual was selected based on having the highest number of  
135 reads. Using 30 threads on a standard Ubuntu workstation, the individual assembly of 38 million paired end  
136 read took approximately 23 hours and 20Gb of RAM, while the joint assembly took five days and 150Gb of  
137 RAM. In addition to this, to eliminate the potential confounding factor of increased coverage, I assembled  
138 another dataset that consisted of a random subsample of 3.8M paired-end reads from each of the 10 samples.

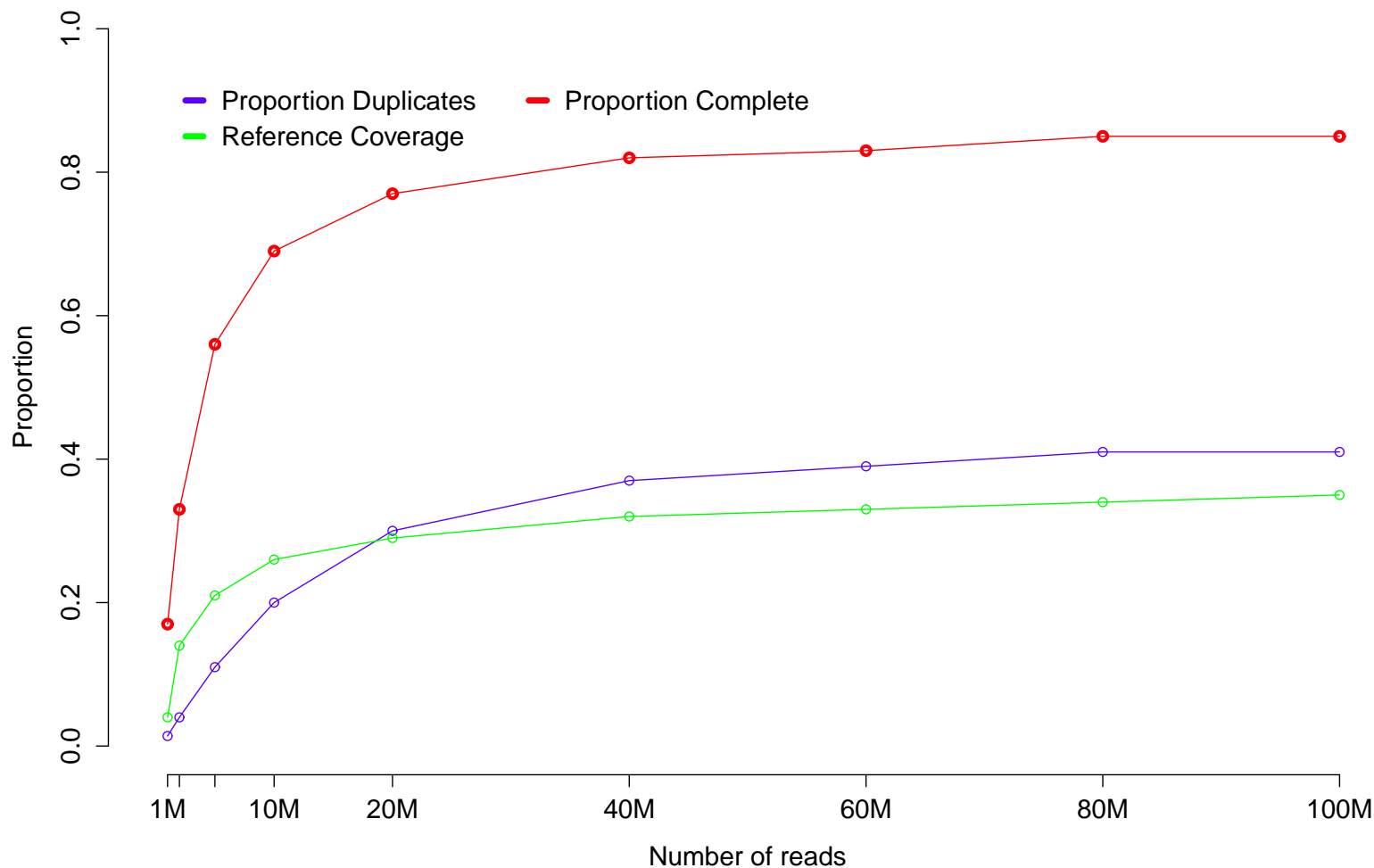
139 This subsampled assembly used similar resources as did the single-individual assembly. Per Table 1, the  
140 joint assembly used more than eight times more reads, and is more than four times larger than the assembly  
141 of a single individual. Despite the additional read data, the **Transrate** score is markedly decreased,  
142 although the BUSCO statistics are slightly better. The large joint assembly suffers from major structural  
143 problems that are unfixable via the proposed filtering procedures. Specifically, read-mapping data suggests  
144 that 28.7% of the contigs in the joint assembly and 18% in the subsampled assembly could be merged,  
145 versus 15% in the single assembly. This structural problem is likely the result of sequence polymorphism  
146 and may cause significant issues for many common downstream processes.

147 **Table 1**

<b>Name</b>	<b>Num. Reads</b>	<b>Num. Contigs</b>	<b>Assembly Size</b>	<b>Score</b>	<b>BUSCO</b>
<b>Single Ind.</b>	38M	205812	131.6Mb	0.3064	C:81%,D:41%,M:9%
<b>Subsampled</b>	38M	304162	183.8Mb	0.2619	C:84%,D:47%,M:8.4%
<b>10 Ind.</b>	269M	913295	440.2Mb	0.22011	C:88%,D:51%,M:5%

149 Table 1. A comparison of the raw assemblies resulting from a single individuals versus the joint  
150 assembly of 10 individuals as well as a joint - but subsampled dataset consisting of a random 3.8M reads  
151 from each of the 10 samples. The individual assembly of 38 million reads resulted in a high quality  
152 assembly as evidenced by the **Transrate** score of 0.3064 (per (Smith-Unna et al., 2015), this is a score  
153 better than 50% of published transcriptomes), and a high BUSCO score. The assembly of 10 individuals  
154 scores lower using **Transrate**, though a small number of new transcripts are discovered. The  
155 joint-subsampled assembly recovered an equivalent number of transcripts, but resulted in a lower  
156 **Transrate** score, indicating a lower quality assembly.

157 **Figure 1**



158 Figure 1. Assembly of multiple subsetting datasets suggests that sequencing beyond 20-40 million paired  
159 end reads does not result in further sequence discovery. Proportion complete indicates the proportion  
160 of BUSCOs that were found to be full length. Proportion duplicates are those BUSCOs that were found  
161 multiple times in the assembly dataset. Reference coverage is a **Transrate** generated metric indicating  
162 the proportion of the reference *Mus* transcriptome found in the *de novo* assembly. Higher numbers for  
163 reference coverage and proportion complete indicate a more complete assembly.

## 164 0.2 Quality Control of Sequence Read Data

165 **Summary Statement:** Visualize your read data. Error correct reads using **bfc** for low to  
166 moderately sized datasets and **RCorrector** for higher coverage datasets. Remove adapters, and

167 **employ gentle quality filtering using PHRED  $\leq 2$  as a threshold.**

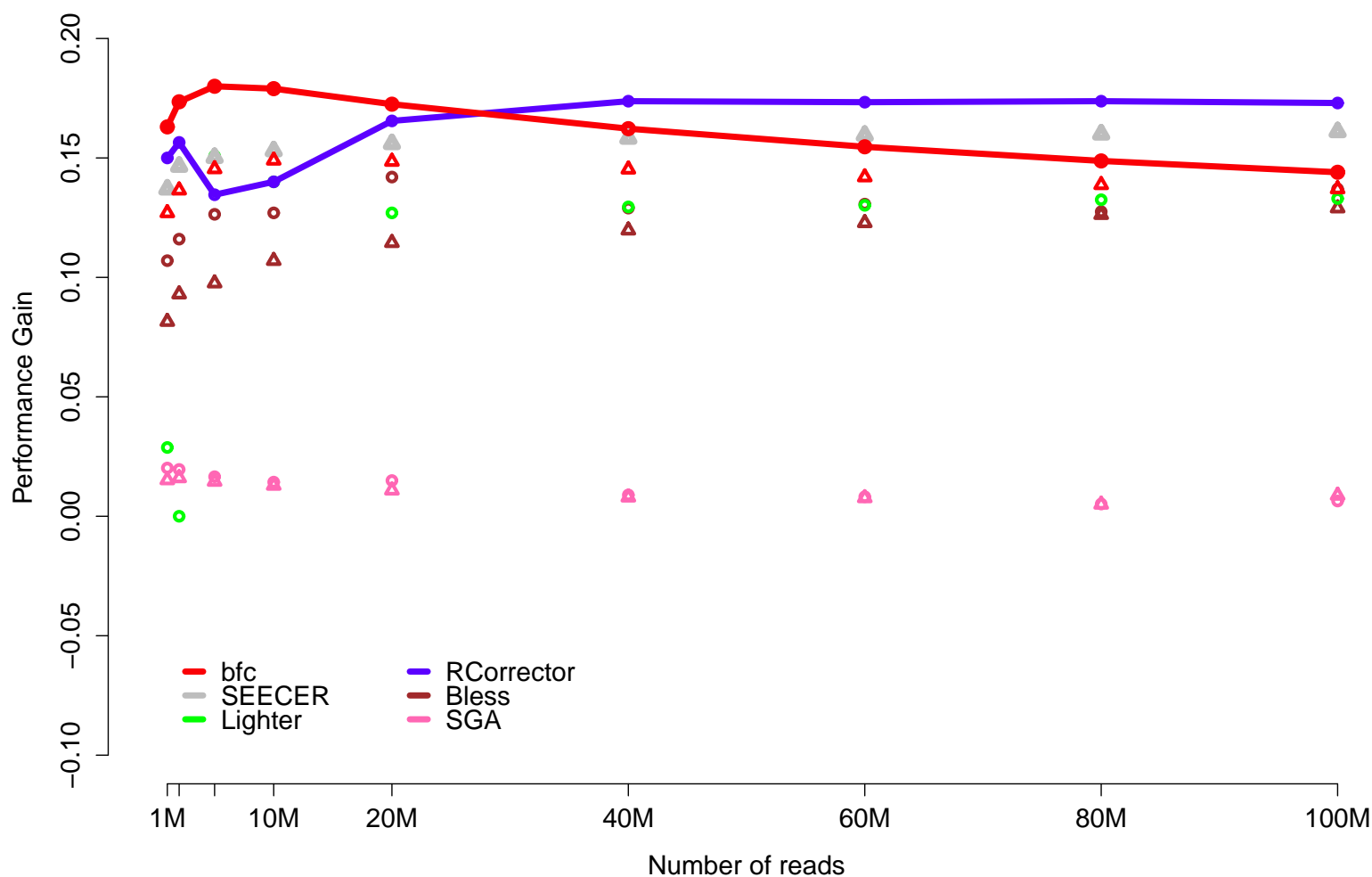
168 Before assembly, it is critical that appropriate quality control steps are implemented. It is often helpful  
169 to generate some metrics of read quality on the raw data. Several software packages are available – I am  
170 fond of **SoLextaQA** (Cox et al., 2010) and **FastQC**. Immediately upon download of the read dataset from the  
171 sequence provider, metrics of read quality, generated by either of these two software packages, should be  
172 generated. Of note – a copy of the raw reads should be compressed and archived, preferably on a physically  
173 separated device for long term archival storage. For this, I have successfully used Amazon S3 cloud storage,  
174 though many options exist.

175 Immediately after visualizing the raw data, error correction of the sequencing reads should be done  
176 (MacManes and Eisen, 2013). A very large number of read correction software packages exist, and several of  
177 them are benchmarked here using the *Mus* (Figure 2, and Tables S1-S11) and *Homo* datasets (Tables  
178 S12-S21). In all evaluated datasets, the error correction **bfc** was the best when correcting less than  
179 approximately 20M paired-end reads. When correcting more, the software **RCorrector** provided the  
180 optimal correction. The effects of error correction on assembly were evaluated using **BUSCO** and **Transrate**.  
181 While error correction did not result in significant improvements in BUSCO metrics, the transrate scores were  
182 substantially improved (Figure 3). These scores were largely improved by the fact that assemblies using  
183 error corrected reads had fewer low-covered bases and contigs, and a slightly higher mapping rate.

184 The error corrected reads are then subjected to vigorous adapter sequence removal, typically using  
185 **Trimmomatic** (Bolger et al., 2014) or **Skewer** (Jiang et al., 2014). With adapter sequence removal may be a  
186 quality trimming step. Here, substantial caution is required, as aggressive trimming has detrimental effects  
187 on assembly quality. Specifically, I recommend trimming at Phred=2 (MacManes, 2014), a threshold  
188 associated with removal of only the lowest quality bases. After adapter removal and quality trimming, the  
189 previously error corrected reads are now ready for *de novo* transcriptome assembly.

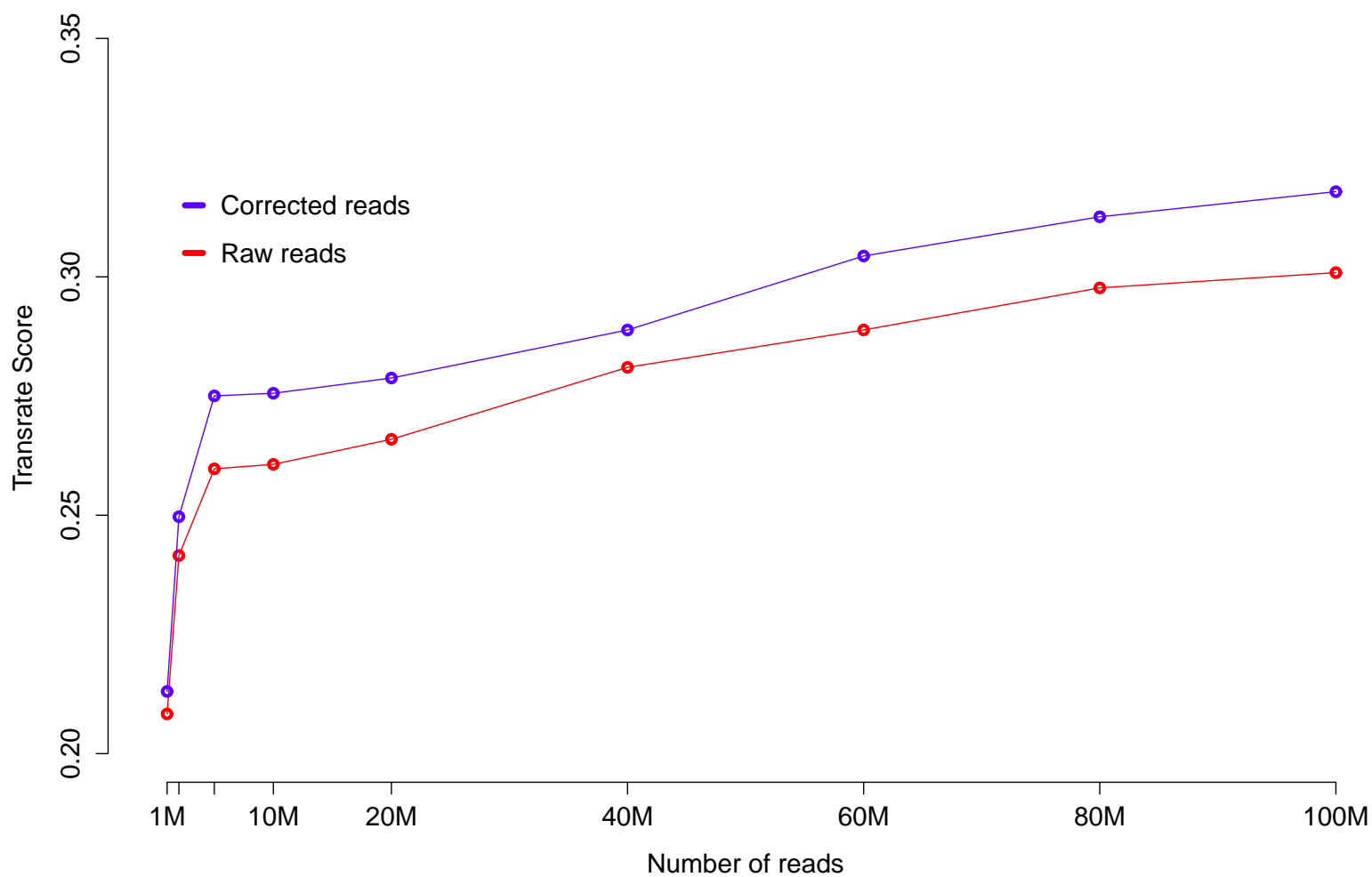


190 **Figure 2**



191 Figure 2. Error correction of reads results in a performance gain (defined as: (perfect error corrected  
192 reads - perfect raw reads) + reads made better - reads made worse). Perfect reads are reads that map to  
193 the reference without mismatch. Better and worse reads are those that map with fewer or more  
194 mismatches. Low coverage datasets are best corrected with **bfc**, which higher coverage datasets are  
195 optimally corrected with **RCorrector**. The best performing corrections improve the quality of more than  
196 15% of reads. To emphasize the patterns of performance of the two best-performing error correctors,  
197 their points are connected by lines.

198 **Figure 3**



199 Figure 3. Error correction (with the best performing correction software, described in Figure 2), results  
200 in a consistent increase in the **Transrate** score, which indicates a higher quality assembly across all  
201 coverage depths.

### 202 **0.3 Coverage normalization**

203 **Summary Statement: Normalize your data, only if you have to.**

204 Depending on the volume of input data, the availability of a high-memory workstation, and the rapidity  
205 with which the assembly is needed, coverage normalization may be employed. This process, which, using a  
206 streaming algorithm and measurement of the median kmer abundance of each read, aims to erode areas of  
207 high coverage while leaving untouched, reads spanning lower coverage areas. Normalization may be  
208 accomplished in the software package *khmer* (Pell et al., 2012), or within *Trinity* using a computational  
209 algorithm based on *khmer*. To evaluate *khmer* normalization, I generated 6 datasets, 3 using 20M paired  
210 reads and 3 using 100M paired reads. These datasets were assembled (available here) using *Trinity* and were  
211 evaluated via standard methods. These tests revealed that normalization did dramatically reduce RAM  
212 requirements and runtime, though it also decreased the number of complete BUSCO's found by an average  
213 of 1% for all assemblies. Normalization also decreased the *Transrate* score on average by 0.0248 for the  
214 20M read assemblies and 0.0272 for the 100M read assemblies. Interestingly, normalization *increased* the  
215 percent BUSCO duplication by an average of only 1% for the smaller assembly, but by over 14% for the  
216 larger assembly. Given this, our recommendation is to employ digital normalization when the assembly is  
217 otherwise impossible, or when results are urgently needed, but that it should not be used by default for the  
218 production of transcriptome assemblies.

### 219 **0.4 Assembly**

220 **Summary Statement: Assemble your data using *Trinity*, then remove poorly supported**  
221 **contigs.**

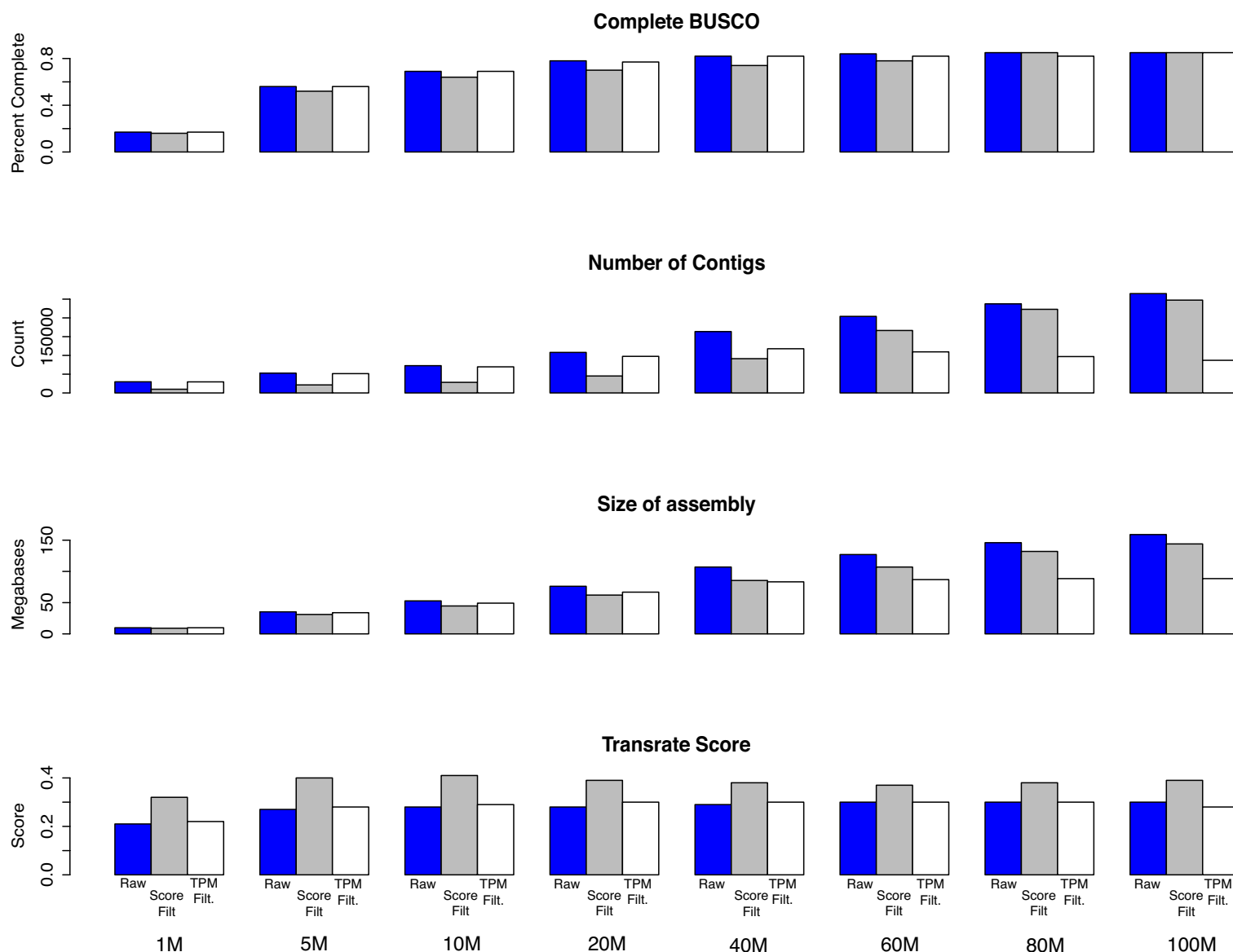
222 For non-model organisms lacking reference genomic resources, the error correction, adapter and quality  
223 trimming reads should be assembled *de novo* into transcripts. Currently, the assembly package *Trinity*  
224 (Haas et al., 2013) is thought to currently be the most accurate (Li et al., 2014), and therefore is  
225 recommended over other assemblers. While attempting a merged assembly with multiple assemblers may  
226 *ultimately* result in the highest quality assembly, options for merging assemblies are currently limited, and  
227 therefore is not recommended.

228 *Trinity*'s underlying algorithm has been pre-optimized to recover large numbers of alternative isoforms,

229 including many that are minimally supported by read data. As a result, in many cases, the raw assembly  
230 will require filtration to remove these assembly artifacts. Reference dependent and independent evaluative  
231 tools (*e.g.*, **Transrate**, **BUSCO**) allow for evidence-based post-assembly filtration. Typically, an initial  
232 quality-evaluation and filtration step is implemented using **Transrate**. This process assigns a score to the  
233 assembly, and creates an alternative assembly by removing contigs based on read-mapping metrics. This  
234 filtration step may result in the removal of a large proportion (as much as 67%) of the transcripts.  
235 Reference-based metrics are generated before and after this filtration step to ensure that filtration has not  
236 been too aggressive - that a significant number of known transcripts have not been removed. After  
237 **Transrate** filtration, or alternative to it, it is often helpful to employ a filtration step based on TPM.  
238 Because underlying assumptions of gene expression estimation software vary, which may result in variation  
239 of the actual estimates, gene expression is typically estimated using two different packages, **Salmon** and  
240 **Kallisto**. Transcripts whose abundance is less than either 1 or 0.5 are removed. Again, reference-based  
241 metrics are generated to ensure that a significant number of known transcripts are not removed.

242 The results of filtration on several datasets of varying size are presented in Figure 4. The reads used in  
243 the 1M,5M,10M,20M subset assemblies were corrected with **bfc**, while the reads for the larger assemblies  
244 were corrected with **RCorrector**. Each dataset was trimmed to a quality of Phred <2, and assembled with  
245 **Trinity**. The raw assembly was filtered by **Transrate** and by gene expression. **BUSCO** evaluation was  
246 performed before and after these filtration steps. In general, for low coverage datasets (less than 20 million  
247 reads), filtering based on expression, using TPM=1 as a threshold performs well, with **Transrate** filtering  
248 being too aggressive. With higher coverage data (more than 60 million reads) **Transrate** filtering may be  
249 optimal, as may gene expression filtering using a threshold of TPM=0.5. Again, the strength of this process  
250 is that it is guided by evidence, with filtering thresholds chosen based upon objective metrics.

251 **Figure 4**



252 Figure 4. Post-assembly filtration. Using assemblies from the 1M,5M,10M,20M,40M,60M,80M,100M  
 253 read subsets, I evaluated the effects of **Transrate** and TPM filtration using a threshold of TPM=1. Both  
 254 **Transrate** and TPM filtering reduced the number of contigs and assembly size, though the magnitudes  
 255 were dependent on the depth of sequencing. BUSCO scores were either decreased in some cases, or stable  
 256 in others, representing the differential effects of filtering on different sized assemblies. In general, for low  
 257 coverage datasets (less than 20 million reads), filtering based on expression, using TPM=1 as a threshold  
 258 performs well, with **Transrate** filtering being too aggressive. With higher coverage data (more than 60  
 259 million reads) **Transrate** filtering may perform better, as mat expression filtering with a lower threshold.

## 260 **0.5 Annotation, post-assembly quality verification, & reporting**

261 **Summary Statement: Verify the quality of your assembly using content based metrics.**

262 **Annotate using dammit Report Transrate score, BUSCO statistics, number of unique transcripts,**  
263 **etc. Do not report meaningless statistics such as N50.**

264 Annotation is a critically important step in transcriptome assembly. Much like other steps, numerous  
265 options exist. Though the research requirements may drive the annotation process, I propose that a core set  
266 of annotations be provided with all *de novo* transcriptome assembly projects. The process through which  
267 these core annotations are accomplished is coordinated by the software package **dammit**. This software takes  
268 as input a fasta file and outputs a standard gff3 containing annotations. After annotation, but before  
269 downstream use, it is important to assess the quality of a transcriptome. Many authors have attempted to  
270 use typical genome assembly quality metrics for this purpose. In particular, N50 and other length-based  
271 summary statistic are often reported (e.g. (Hiz et al., 2014; Liang et al., 2013; Shinzato et al., 2014)).  
272 However, in addition to being a poor proxy for quality in genome assembly (Bradnam et al., 2013), N50 in  
273 the context of a transcriptome assembly carries very little information because the optimal contig length is  
274 not known (Li et al., 2014) - real transcripts vary greatly in length, ranging from tens of nucleotides to tens  
275 of thousands of nucleotides. Reportable metrics should be chosen based on their relevance for assembly  
276 optimization given the biological question at hand. In most cases, this means maximizing the number of  
277 transcripts that can be confidently attributed to the organism, while minimizing the number of technical  
278 artifacts related to the process of sequencing, quality control, and assembly. For many researchers, this  
279 means evaluation with both **BUSCO** and **Transrate**. The statistics found in Table 1 should be presented for  
280 all assemblies, with additional information supplementing these core vital statistics as needed.

## 281 **Testing the Oyster River Protocol**

282 To evaluate the Oyster River Protocol for Transcriptome Assembly, I selected three publicly available  
283 Illumina RNAseq datasets and their corresponding assembled transcriptomes. These three assemblies  
284 included the Nile Tilapia, *Oreochromis niloticus* ((Zhang et al., 2013), SRR797490), an unpublished study  
285 of the Mediterranean black widow, *Latrodectus tredecimguttatus* (SRR954929), and lastly a work on *Delia*  
286 *antiqua* ((Guo et al., 2015), SRR916227). I analyzed the original transcriptomes using both **BUSCO** and  
287 **Transrate**, then followed the protocol as described here. Code for data analysis of the *Oreochromis* is

288 available here. The other samples were processed in an identical fashion. The application of the Oyster  
289 River Protocol on these datasets resulted universally in a substantial (as much as 22%) improvement in the  
290 completeness of assemblies. Given a major goal of these types of studies includes reconstruction all  
291 expressed genes, this improvement may have substantial improvement on downstream work. The  
292 **Transrate** score was dramatically improved as well, particularly in the *Oreochromis* and *Delia* assemblies.  
293 This improvement speaks to the improvement of the structure of the assembly.

294 The filtering process through which these more optimal assemblies were is key. Evaluating both the  
295 BUSCO and **Transrate** scores before and after, allows for an objective way to decide if filtering has been too  
296 restrictive or not. Indeed, for the *Latrodectus* assembly, both **Transrate** and TPM filtering reduced the  
297 BUSCO score, while substantially increasing the **Transrate** score. Depending on the goals of the experiment,  
298 it may be determined that the structural integrity of the assembly outweighs improved content. In contrast  
299 to how post-assembly filtering is typically done, this method allow for the researcher to make an informed  
300 decision about these processes.

301

Name	Number Reads	Number Contigs	Assembly Size (Mb)	Transrate Score	BUSCO Score
<b>Oreochromis</b>	25.2M	79198/140035/100376/116038/88456	32.0/75.1/69.5/58.6/57.7	0.1103/0.2173/0.4778/0.2595/0.4479	C:39%,M:46%/C:58%,M:28%/C
<b>Latrodectus</b>	27.6M	10259/36394/30932/27973/NA	10.6/13.5/13.1/10.9/NA	0.43673/0.2795/0.4968/0.338/NA	C:48%,M:38%/C:58%,M:28%/C
<b>Delia</b>	25.8M	29451/49099/38614/46145/32689	12.4/19.3/18.8/17.9/15.8	0.393/0.2036/0.4572/0.2305/0.4341	C:40%,M:48%/C:62%,M:21%/C

303 Table 2. The results of the application of the Oyster River Protocol to three available transcriptomes.  
304 Within each column, the 5 metrics, separated by forward slashes are: 1. The original assembly 2. The  
305 raw **Trinity** assembly 3. The **Transrate** filtered assembly 4. The TPM=1 filtered assembly, and 5. The  
306 **Transrate** filtered assembly that has been further filtered by expression. In all cases the assembly  
307 content, as evaluated by the BUSCO score is dramatically improved over the original assembly. These  
308 content-improved assemblies have acceptable **Transrate** scores, which in 2 of 3 cases are vastly superior  
309 to the scores of the original assembly.

## 310 Conclusions

311 With the rapid adoption of high-throughput sequencing, studies of functional, evolutionary and population  
312 genomics are now being done by hundreds or even thousands of labs around the world. These studies  
313 typically begin with a *de novo* transcriptome assembly. Assembly may be accomplished in one of several  
314 different ways, using different software packages, with each method producing different results. This  
315 complexity begs the question – *Which method(s) are optimal?* Using reference and non-reference based  
316 evaluative methods, I have proposed a set of guidelines **The Oyster River Protocol for Transcriptome**  
317 **Assembly** that aim to standardize and facilitate the process of transcriptome assembly. These

318 recommendations include limiting assembly to between 20 million and 40 million sequencing reads from  
319 single individual where possible, error correction of reads, gently quality trimming, assembly filtering using  
320 **Transrate** or gene expression, annotation using **dammit**, and appropriate reporting. The processes result in  
321 a high quality transcriptome assembly appropriate for downstream usage.

## 322 **Acknowledgments**

323 This work was significantly improved by discussions with Richard Smith-Unna, Brian Haas and many  
324 others. More generally, the work and it's presentation has been influenced by supporters of the Open Access  
325 and Science movements.

## 326 **References**

- 327 Alameldin, H.F., Awad, S., Boucher, E., Caldwell, A., Cartwright, R., Charbonneau, A., Constantinides, B.,  
328 Fay, S., Fenton, J., Fenzl, T., Fish, J., Garcia-Gutierrez, L., Garland, P., Gluck, J., González, I.,  
329 Guermond, S., Guo, J., Gupta, A., Herr, J.R., Howe, A., Hyer, A., Härpfer, a., Irber, L., Kidd, R., Lin, D.,  
330 Lippi, J., Mansour, T., McA’Nulty, P., Mizzi, J., Murray, K.D., Nahum, J.R., Nanlohy, K., Nederbragt,  
331 A.J., Ortiz-Zuazaga, H., Ory, J., Pell, J., Pepe-Ranney, C., Russ, Z.N., Schwarz, E., Seaman, J., Sievert,  
332 S., Simpson, J., Skennerton, C.T., Spencer, J., Srinivasan, R., Standage, D., Stapleton, J.A., Steinman,  
333 S.R., Stein, J., Taylor, B., Trimble, W., Wiencko, H.L., Wright, M., Wyss, B., Zhang, Q., Zyme, E., 2015.  
334 The khmer software package: enabling efficient nucleotide sequence analysis. *F1000Research* 4, 900.
- 335 Au, K.F., Sebastiano, V., Afshar, P.T., Durruthy, J.D., Lee, L., Williams, B.A., van Bakel, H., Schadt, E.E.,  
336 Reijo-Pera, R.A., Underwood, J.G., Wong, W.H., 2013. Characterization of the human ESC  
337 transcriptome by hybrid sequencing. *PNAS* 110, 201320101–30.
- 338 Bolger, A.M., Lohse, M., Usadel, B., 2014. Trimmomatic: a flexible trimmer for Illumina sequence data.  
339 *Bioinformatics* 30, btu170–2120.
- 340 Bradnam, K.R., Fass, J.N., Alexandrov, A., Baranay, P., Bechner, M., Birol, I., Boisvert, S., Chapman, J.A.,  
341 Chapuis, G., Chikhi, R., Chitsaz, H., Chou, W.C., Corbeil, J., Del Fabbro, C., Docking, T.R., Durbin, R.,  
342 Earl, D., Emrich, S., Fedotov, P., Fonseca, N.A., Ganapathy, G., Gibbs, R.A., Gnerre, S., Godzaridis, E.,  
343 Goldstein, S., Haimel, M., Hall, G., Haussler, D., Hiatt, J.B., Ho, I.Y., Howard, J., Hunt, M., Jackman,  
344 S.D., Jaffe, D.B., Jarvis, E.D., Jiang, H., Kazakov, S., Kersey, P.J., Kitzman, J.O., Knight, J.R., Koren,



- 345 S., Lam, T.W., Lavenier, D., Laviolette, F., Li, Y., Li, Z., Liu, B., Liu, Y., Luo, R., Maccallum, I.,  
346 MacManes, M.D., Maillet, N., Melnikov, S., Naquin, D., Ning, Z., Otto, T.D., Paten, B., Paulo, O.S.,  
347 Phillippy, A.M., Pina-Martins, F., Place, M., Przybylski, D., Qin, X., Qu, C., Ribeiro, F.J., Richards, S.,  
348 Rokhsar, D.S., Ruby, J.G., Scalabrin, S., Schatz, M.C., Schwartz, D.C., Sergushichev, A., Sharpe, T.,  
349 Shaw, T.I., Shendure, J., Shi, Y., Simpson, J.T., Song, H., Tsarev, F., Vezzi, F., Vicedomini, R., Vieira,  
350 B.M., Wang, J., Worley, K.C., Yin, S., Yiu, S.M., Yuan, J., Zhang, G., Zhang, H., Zhou, S., Korf, I.F.,  
351 2013. Assemblathon 2: evaluating *de novo* methods of genome assembly in three vertebrate species.  
352 GigaScience 2, 10.
- 353 Bray, N., Pimentel, H., Melsted, P., Pachter, L., 2015. Near-optimal RNA-Seq quantification. arXiv.org  
354 arXiv:14054600747265555258.
- 355 Brown, C.T., Howe, A., Zhang, Q., Pyrkosz, A.B., Brom, T.H., 2012. A Reference-Free Algorithm for  
356 Computational Normalization of Shotgun Sequencing Data. arXiv.org arXiv:1203.4802v2.
- 357 Cox, M.P., Peterson, D.A., Biggs, P.J., 2010. SolexaQA: At-a-glance quality assessment of Illumina  
358 second-generation sequencing data. BMC Bioinformatics 11, 485.
- 359 Fitzpatrick, M., Ben-Shahar, Y., Vet, L., Smid, H., Robinson, G.E., Sokolowski, M., 2005. Candidate genes  
360 for behavioural ecology. Trends In Ecology & Evolution 20, 96–104.
- 361 Garber, M., Grabherr, M.G., Guttman, M., Trapnell, C., 2011. Computational methods for transcriptome  
362 annotation and quantification using RNA-seq. Nature Methods 8, 469–477.
- 363 Glenn, T.C., 2011. Field guide to next-generation DNA sequencers. Molecular Ecology Resources 11,  
364 759–769.
- 365 Guo, Q., Hao, Y.J., Li, Y., Zhang, Y.J., Ren, S., Si, F.L., Chen, B., 2015. Gene cloning, characterization  
366 and expression and enzymatic activities related to trehalose metabolism during diapause of the onion  
367 maggot *Delia antiqua* (Diptera: Anthomyiidae). Gene 565, 106–115.
- 368 Haas, B.J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P.D., Bowden, J., Couger, M.B., Eccles,  
369 D., Li, B., Lieber, M., MacManes, M.D., Ott, M., Orvis, J., Pochet, N., Strozzi, F., Weeks, N.,  
370 Westerman, R., William, T., Dewey, C.N., Henschel, R., Leduc, R.D., Friedman, N., Regev, A., 2013. *De*  
371 *novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation  
372 and analysis. Nature Protocols 8, 1494–1512.

- 373 Han, H., Irimia, M., Ross, P.J., Sung, H.K., Alipanahi, B., David, L., Golipour, A., Gabut, M., Michael, I.P.,  
374 Nachman, E.N., Wang, E., Trcka, D., Thompson, T., O'Hanlon, D., Slobodeniuc, V., Barbosa-Morais,  
375 N.L., Burge, C.B., Moffat, J., Frey, B.J., Nagy, a., Ellis, J., Wrana, J.L., Blencowe, B.J., 2013. MBNL  
376 proteins repress ES-cell-specific alternative splicing and reprogramming. *Nature* 498, 241–245.
- 377 Heo, Y., Wu, X.L., Chen, D., Ma, J., Hwu, W.M., 2014. BLESS: Bloom filter-based error correction  
378 solution for high-throughput sequencing reads. *Bioinformatics* 30, 1354–1362.
- 379 Hiz, M.C., Canher, B., Niron, H., Turet, M., 2014. Transcriptome analysis of salt tolerant common bean  
380 (*Phaseolus vulgaris* L.) under saline conditions. *PLOS ONE* 9, e92598.
- 381 Iqbal, Z., Caccamo, M., Turner, I., Flicek, P., McVean, G., 2012. *De novo* assembly and genotyping of  
382 variants using colored *de Bruijn* graphs. *Nature Publishing Group* 44, 226–232.
- 383 Jiang, H., Lei, R., Ding, S.W., Zhu, S., 2014. Skewer: a fast and accurate adapter trimmer for  
384 next-generation sequencing paired-end reads. *BMC Bioinformatics* 15, 182.
- 385 Johnson, J.M., Castle, J., Garrett-Engle, P., Kan, Z., Loerch, P.M., Armour, C.D., Santos, R., Schadt,  
386 E.E., Stoughton, Roland Shoemaker, D.D., 2003. Genome-wide survey of human alternative pre-mRNA  
387 splicing with exon junction microarrays. *Science* 302, 2141–2144.
- 388 Le, H.S., Schulz, M.H., McCauley, B.M., Hinman, V.F., Bar-Joseph, Z., 2013. Probabilistic error correction  
389 for RNA sequencing. *Nucleic Acids Research* 41, 1–11.
- 390 Li, B., Fillmore, N., Bai, Y., Collins, M., Thomson, J.A., Stewart, R., Dewey, C.N., 2014. Evaluation of *de*  
391 *novo* transcriptome assemblies from RNA-Seq data. *Genome Biology* 15, 663–21.
- 392 Li, H., 2015. BFC: correcting Illumina sequencing errors. *Bioinformatics* 31, 2885–2887.
- 393 Liang, C., Liu, X., Yiu, S.M., Lim, B.L., 2013. *De novo* assembly and characterization of *Camelina sativa*  
394 transcriptome by paired-end sequencing. *BMC Genomics* 14, 146.
- 395 Love, M.I., Huber, W., anders, S., 2014. Moderated estimation of fold change and dispersion for RNA-seq  
396 data with DESeq2. *Genome Biology* 15, 550.
- 397 MacManes, M.D., 2014. On the optimal trimming of high-throughput mRNA sequence data. *Frontiers in*  
398 *Genetics* 5.
- 399 MacManes, M.D., Eisen, M.B., 2013. Improving transcriptome assembly through error correction of  
400 high-throughput sequence reads. *PeerJ* 1, e113.

- 401 MacManes, M.D., Lacey, E.A., 2012. The Social Brain: Transcriptome Assembly and Characterization of  
402 the Hippocampus from a Social Subterranean Rodent, the Colonial Tuco-Tuco (*Ctenomys sociabilis*).  
403 PLOS ONE 7, e45524.
- 404 Martin, M., 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
405 EMBnet.journal 17, pp. 10–12.
- 406 McDonald, E., Brown, C.T., 2013. khmer: Working with Big Data in Bioinformatics. arXiv.org  
407 arXiv:1303.2223v1.
- 408 Modrek, B., Resch, A., Grasso, C., Lee, C., 2001. Genome-wide detection of alternative splicing in  
409 expressed sequences of human genes. Nucleic Acids Research 29, 2850–2859.
- 410 Panhuis, T.M., 2006. Molecular evolution and population genetic analysis of candidate female reproductive  
411 genes in *Drosophila*. Genetics 173, 2039–2047.
- 412 Patro, R., Duggal, G., Kingsford, C., 2015. Accurate, fast, and model-aware transcript expression  
413 quantification with Salmon. biorxiv.org , 1–35.
- 414 Pell, J., Hintze, A., Canino-Koning, R., Howe, A., Tiedje, J.M., Brown, C.T., 2012. Scaling metagenome  
415 sequence assembly with probabilistic *de Bruijn* graphs. Proceedings of the National Academy of Sciences  
416 109, 13272–13277.
- 417 Robinson, M.D., McCarthy, D.J., Smyth, G.K., 2010. edgeR: a Bioconductor package for differential  
418 expression analysis of digital gene expression data. Bioinformatics 26, 139–140.
- 419 Shinzato, C., Inoue, M., Kusakabe, M., 2014. A snapshot of a coral "holobiont": a transcriptome assembly  
420 of the scleractinian coral, porites, captures a wide variety of genes from both the host and symbiotic  
421 zooxanthellae. PLOS ONE 9, e85182.
- 422 Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., Zdobnov, E.M., 2015. BUSCO: assessing  
423 genome assembly and annotation completeness with single-copy orthologs. Bioinformatics 31, 3210–3212.
- 424 Simpson, J.T., Durbin, R., 2012. Efficient *de novo* assembly of large genomes using compressed data  
425 structures. Genome Research 22, 549–556.
- 426 Smith-Unna, R.D., Bournnell, C., Patro, R., Hibberd, J.M., Kelly, S., 2015. TransRate: reference free  
427 quality assessment of *de novo* transcriptome assemblies. bioRxiv URL:  
428 <http://biorxiv.org/content/early/2015/06/27/021626>.

- 429 Song, L., Florea, L., 2015. Rcorrector: efficient and accurate error correction for Illumina RNA-seq reads.  
430 GigaScience 4, 48.
- 431 Song, L., Florea, L., Langmead, B., 2014. Lighter: fast and memory-efficient sequencing error correction  
432 without counting. Genome Biology 15, 509.
- 433 Studholme, D.J., 2010. *De novo* assembly of short sequence reads. Briefings In Bioinformatics 11, 457–472.
- 434 Vijay, N., Poelstra, J.W., Künstner, A., Wolf, J.B.W., 2013. Challenges and strategies in transcriptome  
435 assembly and differential gene expression quantification. A comprehensive *in silico* assessment of  
436 RNA-seq experiments. Molecular Ecology 22, 620–634.
- 437 Wolf, J.B.W., 2013. Principles of transcriptome analysis and gene expression quantification: an RNA-seq  
438 tutorial. Molecular Ecology Resources 13, 559–572.
- 439 Xie, Y., Wu, G., Tang, J., Luo, R., Patterson, J., Liu, S., Huang, W., He, G., Gu, S., Li, S., Zhou, X., Li,  
440 Y., Xu, X., Wong, G.K.S., Wang, J., 2014. SOAPdenovo-Trans: *de novo* transcriptome assembly with  
441 short RNA-Seq reads. Bioinformatics 30, 1660–1666.
- 442 Zhang, R., Zhang, L.l., Ye, X., Tian, Y.y., Sun, C.f., Lu, M.x., Bai, J.j., 2013. Transcriptome profiling and  
443 digital gene expression analysis of Nile tilapia (*Oreochromis niloticus*) infected by *Streptococcus*  
444 *agalactiae*. Molecular biology reports 40, 5657–5668.