

The Oyster River Protocol: A multi assembler and kmer approach for *de novo* transcriptome assembly

Matthew D MacManes ^{Corresp. 1}

¹ Department of Molecular, Cellular and Biomedical Sciences, University of New Hampshire, Durham, NH, United States

Corresponding Author: Matthew D MacManes

Email address: macmanes@gmail.com

Characterizing transcriptomes in non-model organisms has resulted in a massive increase in our understanding of biological phenomena. This boon, largely made possible via high-throughput sequencing, means that studies of functional, evolutionary and population genomics are now being done by hundreds or even thousands of labs around the world. For many, these studies begin with a *de novo* transcriptome assembly, which is a technically complicated process involving several discrete steps. The Oyster River Protocol (ORP), described here, implements a standardized and benchmarked set of bioinformatic processes, resulting in an assembly with enhanced qualities over other standard assembly methods. Specifically, ORP produced assemblies have higher Detonate and TransRate scores and mapping rates, which is largely a product of the fact that it leverages a multi-assembler and kmer assembly process, thereby bypassing the shortcomings of any one approach. These improvements are important, as previously unassembled transcripts are included in ORP assemblies, resulting in a significant enhancement of the power of downstream analysis. Further, as part of this study, I show that assembly quality is unrelated with the number of reads generated, above 30 million reads. Code Availability: The version controlled open-source code is available at https://github.com/macmanes-lab/Oyster_River_Protocol. Instructions for software installation and use, and other details are available at <http://oyster-river-protocol.rtfld.org/>.

The Oyster River Protocol: A Multi-Assembler and Kmer Approach For *de novo* Transcriptome Assembly

Matthew D. MacManes^{1,2}

¹Department of Molecular Cellular and Biomedical Sciences, University of New Hampshire, Durham NH, USA

²Hubbard Center for Genomic Studies, University of New Hampshire, Durham NH, USA

Corresponding author:

Matthew MacManes^{1,2}

Email address: macmanes@gmail.com

ABSTRACT

Characterizing transcriptomes in non-model organisms has resulted in a massive increase in the understanding of biological phenomena. This boon, largely made possible via high-throughput sequencing, means that studies of functional, evolutionary and population genomics are now being done by hundreds or even thousands of labs around the world. For many, these studies begin with a *de novo* transcriptome assembly, which is a technically complicated process involving several discrete steps. The Oyster River Protocol (ORP), described here, implements a standardized and benchmarked set of bioinformatic processes, resulting in an assembly with enhanced qualities over other standard assembly methods. Specifically, ORP produced assemblies have higher *Detonate* and *TransRate* scores and mapping rates, which is largely a product of the fact that it leverages a multi-assembler and kmer assembly process, thereby bypassing the shortcomings of any one approach. These improvements are important, as previously unassembled transcripts are included in ORP assemblies, resulting in a significant enhancement of the power of downstream analysis. Further, as part of this study, I show that assembly quality is unrelated with the number of reads generated, above 30 million reads. **Code Availability:** The version controlled open-source code is available at https://github.com/macmanes-lab/Oyster_River_Protocol. Instructions for software installation and use, and other details are available at <http://oyster-river-protocol.rtf.d.org/>.

1 INTRODUCTION

For all biology, modern sequencing technologies have provided for an unprecedented opportunity to gain a deep understanding of genome level processes that underlie a very wide array of natural phenomena, from intracellular metabolic processes to global patterns of population variability. Transcriptome sequencing has been influential (Mortazavi et al., 2008; Wang et al., 2009), particularly in functional genomics (Lappalainen et al., 2013; Cahoy et al., 2008), and has resulted in discoveries not possible even just a few years ago. This in large part is due to the scale at which these studies may be conducted (Li et al., 2017; Tan et al., 2017). Unlike studies of adaptation based on one or a small number of candidate genes (*e.g.*, (Fitzpatrick et al., 2005; Panhuis, 2006)), modern studies may assay the entire suite of expressed transcripts – the transcriptome – simultaneously. In addition to issues of scale, as a direct result of enhanced dynamic range, newer sequencing studies have increased ability to simultaneously reconstruct and quantitate lowly- and highly-expressed transcripts (Wolf, 2013; Vijay et al., 2013). Lastly, improved methods for the detection of differences in gene expression (*e.g.*, (Robinson et al., 2010; Love et al., 2014)) across experimental treatments have resulted in increased resolution for studies aimed at understanding changes in gene expression.

As a direct result of their widespread popularity, a diverse tool set for the assembly of transcriptome exists, with each potentially reconstructing transcripts others fail to reconstruct. Amongst the earliest of

specialized *de novo* transcriptome assemblers were the packages Trans-ABYSS (Robertson et al., 2010), Oases (Schulz et al., 2012), and SOAPdenovoTrans (Xie et al., 2014), which were fundamentally based on the popular *de Bruijn* graph-based genome assemblers ABySS (Simpson et al., 2009), Velvet (Zerbino and Birney, 2008), and SOAP (Li et al., 2008) respectively. These early efforts gave rise to a series of more specialized *de novo* transcriptome assemblers, namely Trinity (Haas et al., 2013), and IDBA-Tran (Peng et al., 2013). While the *de Bruijn* graph approach remains powerful, newly developed software explores novel parts of the algorithmic landscape, offering substantial benefits, assuming novel methods reconstruct different fractions of the transcriptome. BinPacker (Liu et al., 2016), for instance, abandons the *de Bruijn* graph approach to model the assembly problem after the classical bin packing problem, while Shannon (Kannan et al., 2016) uses information theory, rather than a set of software engineer-decided heuristics. These newer assemblers, by implementing fundamentally different assembly algorithms, may reconstruct fractions of the transcriptome that other assemblers fail to accurately assemble.

In addition to the variety of tools available for the *de novo* assembly of transcripts, several tools are available for pre-processing of reads via read trimming (e.g., Skewer (Jiang et al., 2014), Trimmomatic (Bolger et al., 2014), Cutadapt (Martin and 2011, 2011)), read normalization (khmer (Pell et al., 2012)), and read error correction (SEECER (Le et al., 2013) and RCorrector (Song and Florea, 2015), Reptile (Yang et al., 2010)). Similarly, benchmarking tools that evaluate the quality of assembled transcriptomes including TransRate (Smith-Unna et al., 2016), BUSCO (Benchmarking Universal Single-Copy Orthologs - (Simão et al., 2015)), and Detonate (Li et al., 2014) have been developed. Despite the development of these evaluative tools, this manuscript describes the first systematic effort coupling them with the development of a *de novo* transcriptome assembly pipeline.

The ease with which these tools may be used to produce and characterize transcriptome assemblies belies the true complexity underlying the overall process (Ungaro et al., 2017; Wang and Gribskov, 2017; Moreton et al., 2015; Yang and Smith, 2013). Indeed, the subtle (and not so subtle) methodological challenges associated with transcriptome reconstruction may result in highly variable assembly quality. In particular, while most tools run using default settings, these defaults may be sensible only for one specific (often unspecified) use case or data type. Because parameter optimization is both dataset-dependent and factorial in nature, an exhaustive optimization particularly of entire pipelines, is never possible. Given this, the production of a *de novo* transcriptome assembly requires a large investment in time and resources, with each step requiring careful consideration. Here, I propose an evidence-based protocol for assembly that results in the production of high quality transcriptome assemblies, across a variety of commonplace experimental conditions or taxonomic groups.

This manuscript describes the development of The Oyster River Protocol¹ for transcriptome assembly. It explicitly considers and attempts to address many of the shortcomings described in (Vijay et al., 2013), by leveraging a multi-kmer and multi-assembler strategy. This innovation is critical, as all assembly solutions treat the sequence read data in ways that bias transcript recovery. Specifically, with the development of assembly software comes the use of a set of heuristics that are necessary given the scope of the assembly problem itself. Given each software development team carries with it a unique set of ideas related to these heuristics while implementing various assembly algorithms, individual assemblers exhibit unique assembly behavior. By leveraging a multi-assembler approach, the strengths of one assembler may complement the weaknesses of another. In addition to biases related to assembly heuristics, it is well known that assembly kmer-length has important effects on transcript reconstruction, with shorter kmers more efficiently reconstructing lower-abundance transcripts relative to more highly abundant transcripts. Given this, assembling with multiple different kmer lengths, then merging the resultant assemblies may effectively reduce this type of bias. Recognizing these issue, I hypothesize that an assembly that results from the combination of multiple different assemblers and lengths of assembly-kmers will be better than each individual assembly, across a variety of metrics.

In addition to developing an enhanced pipeline, the work suggests an exhaustive way of characterizing assemblies while making available a set of fully-benchmarked reference assemblies that may be used by other researchers in developing new assembly algorithms and pipelines. Although many other researchers have published comparisons of assembly methods, up until now these have been limited to single datasets assembled a few different ways (Marchant et al., 2016; Finseth and Harrison, 2014), thereby failing to

100 provide more general insights.

101 2 METHODS

102 2.1 Datasets

103 In an effort at benchmarking the assembly and merging protocols, I downloaded a set of publicly available
104 RNAseq datasets (Table 1) that had been produced on the Illumina sequencing platform. These datasets
105 were chosen to represent a variety of taxonomic groups, so as to demonstrate the broad utility of the
106 developed methods. Because datasets were selected randomly with respect to sequencing center and read
107 number, they are likely to represent the typical quality of Illumina data circa 2014-2017.

108

Table 1

Type	Accession	Species	Num. Reads	Read Length
Animalia	ERR489297	<i>Anopheles gambiae</i>	206M	100bp
Animalia	DRR030368	<i>Echinococcus multilocularis</i>	73M	100bp
Animalia	ERR1016675	<i>Heterorhabditis indica</i>	51M	100bp
Animalia	SRR2086412	<i>Mus musculus</i>	54M	100bp
Animalia	DRR036858	<i>Mus musculus</i>	114M	100bp
Animalia	DRR046632	<i>Oncorhynchus mykiss</i>	82M	76bp
Animalia	SRR1789336	<i>Oryctolagus cuniculus</i>	31M	100bp
109 Animalia	SRR2016923	<i>Phyllodoce medipapillata</i>	86M	100bp
Animalia	ERR1674585	<i>Schistosoma mansoni</i>	39M	100bp
Plant	DRR082659	<i>Aeginetia indica</i>	69M	90bp
Plant	DRR053698	<i>Cephalotus follicularis</i>	126M	90bp
Plant	DRR069093	<i>Hevea brasiliensis</i>	103M	100bp
Plant	SRR3499127	<i>Nicotiana tabacum</i>	30M	150bp
Plant	DRR031870	<i>Vigna angularis</i>	60M	100bp
Protozoa	ERR058009	<i>Entamoeba histolytica</i>	68M	100bp

110 Table 1 lists the datasets used in this study. All datasets are publicly available for download by accession
111 number at the European Nucleotide Archive or NCBI Short Read Archive.

112 2.2 Software

113 The Oyster River Protocol can be installed on the Linux platform, and does not require superuser
114 privileges, assuming Linuxbrew (Jackman and Birol, 2016) is installed. The software is implemented
115 as a stand-alone makefile which coordinates all steps described below. All scripts are available at
116 https://github.com/macmanes-lab/Oyster_River_Protocol, and run on the Linux
117 platform. The software is version controlled and openly-licensed to promote sharing and reuse. A guide
118 for users is available at <http://oyster-river-protocol.rtf.d.io>.

119 2.3 Pre-assembly procedures

120 For all assemblies performed, Illumina sequencing adapters were removed from both ends of the se-
121 quencing reads, as were nucleotides with quality Phred ≤ 2 , using the program Trimmomatic version
122 0.36 (Bolger et al., 2014), following the recommendations from (MacManes, 2014). After trimming,
123 reads were error corrected using the software RCorrector version 1.0.2 (Song and Florea, 2015),
124 following recommendations from (MacManes and Eisen, 2013). The code for running this step of the
125 Oyster River protocols is available at https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/oyster.mk#L134. The trimmed and error corrected reads were then
126 subjected to *de novo* assembly.
127

¹Named the Oyster River Protocol because the ideas, and some of the code, was developed while overlooking the Oyster River, located in Durham, New Hampshire, NB, the naming assembly of protocols after bodies of water was, to the best of my knowledge, first done by C. Titus Brown (The Eel Pond Protocol: <http://khmer-protocols.readthedocs.io/en/latest/mrnaseq/index.html>), and may have subconsciously influenced me in naming this protocol.

2.4 Assembly

I assembled each trimmed and error corrected dataset using three different *de novo* transcriptome assemblers and three different kmer lengths, producing 4 unique assemblies. First, I assembled the reads using Trinity release 2.4.0 (Haas et al., 2013), and default settings (k=25), without read normalization. The decision to forgo normalization is based on previous work (MacManes, 2015) showing slightly worse performance of normalized datasets. Next, the SPAdes RNAseq assembler (version 3.10) (Chikhi and Medvedev, 2014) was used, in two distinct runs, using kmer sizes 55 and 75. Lastly, reads were assembled using the assembler Shannon version 0.0.2 (Kannan et al., 2016), using a kmer length of 75. These assemblers were chosen based on the fact that they [1] use an open-science development model, whereby end-users may contribute code, [2] are all actively maintained and are undergoing continuous development, and [3] occupy different parts of the algorithmic landscape.

This assembly process resulted in the production of four distinct assemblies. The code for running this step of the Oyster River protocols is available at https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/oyster.mk#L142.

2.5 Assembly Merging via OrthoFuse

To merge the four assemblies produced as part of the Oyster River Protocol, I developed new software that effectively merges transcriptome assemblies. Described in brief, OrthoFuse begins by concatenating all assemblies together, then forms groups of transcripts by running a version of OrthoFinder (Emms and Kelly, 2015) packaged with the ORP, modified to accept nucleotide sequences from the merged assembly. These groupings represent groups of homologous transcripts. While isoform reconstruction using short-read data is notoriously poor, by increasing the inflation parameter by default to I=4, it attempts to prevent the collapsing of transcript isoforms into single groups. After OrthoFinder has completed, a modified version of TransRate version 1.0.3 (Smith-Unna et al., 2016) which is packaged with the ORP, is run on the merged assembly, after which the best (= highest contig score) transcript is selected from each group and placed in a new assembly file to represent the entire group. The resultant file, which contains the highest scoring contig for each orthogroup, may be used for all downstream analyses. OrthoFuse is run automatically as part of the Oyster River Protocol, and additionally is available as a stand alone script, https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/orthofuser.mk.

2.6 Assembly Evaluation

All assemblies were evaluated using ORP-TransRate, Detonate version 1.11 (Li et al., 2014), shmlast version 1.2 (Scott, 2017), and BUSCO version 3.0.2 (Simão et al., 2015). TransRate evaluates transcriptome assembly contiguity by producing a score based on length-based and mapping metrics, while Detonate conducts an orthogonal analysis, producing a score that is maximized by an assembly that is representative of input sequence read data. BUSCO evaluates assembly content by searching the assemblies for conserved single copy orthologs found in all Eukaryotes. I report default BUSCO metrics as described in (Simão et al., 2015). Specifically, "complete orthologs", are defined as query transcripts that are within 2 standard deviations of the length of the BUSCO group mean, while contigs falling short of this metric are listed as "fragmented". Shmlast implements the conditional reciprocal best hits (CRBH) test (Aubry et al., 2014), conducted in this case against the Swiss-Prot protein database (downloaded October, 2017) using an e-value of 1E-10.

In addition to the generation of metrics to evaluation the quality of transcriptome assemblies, I generated a distance matrix of assemblies for each dataset using the sourmash package (Titus Brown and Irber, 2016), in an attempt at characterizing the algorithmic landscape of assemblers. Specifically, each assembly was characterized using the compute function using 5000 independent sketches. The distance between assemblies was calculated using the compare function and a kmer length of 51. These distance matrices were visualized using the isoMDS function of the MASS package (<https://CRAN.R-project.org/package=MASS>).

2.7 Statistics

All statistical analyses were conducted in R version 3.4.0 (R Core Development Team, 2011). Violin plots were constructed using the beanplot (Kampstra, 2008) and the beeswarm R packages (<https://>

180 //CRAN.R-project.org/package=beeswarm). Expression distributions were plotted using the
181 ggridges package (<https://CRAN.R-project.org/package=ggridges>).

182 3 RESULTS AND DISCUSSION

183 Fifteen RNAseq datasets, ranging in size from (30-206M paired end reads) were assembled using the
184 Oyster River Protocol and with Trinity. Each assembly was evaluated using the software BUSCO,
185 shmlast, Detonate, and TransRate. From these, several metrics were chosen to represent the
186 quality of the produced assemblies. Of note, all the assemblies produced as part of this work are available at
187 <https://www.dropbox.com/sh/ehxvd0ont9ge8id/AABZxRCwcpaxb7rXWclTBbJga>, and
188 will be moved to dataDryad after acceptance. A file containing the evaluative metrics is available at
189 [https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/](https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/orp.csv)
190 [orp.csv](https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/orp.csv), while the distance matrices are available within the folder [https://github.com/](https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/)
191 [macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/](https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/). R code used to
192 conduct analyses and make figures is found at [https://github.com/macmanes-lab/Oyster_](https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/R-analysis.Rmd)
193 [River_Protocol/blob/master/manuscript/R-analysis.Rmd](https://github.com/macmanes-lab/Oyster_River_Protocol/blob/master/manuscript/R-analysis.Rmd).

194 3.1 Assembled transcriptomes

195 The Trinity assembly of trimmed and error corrected reads generally completed on a standard Linux
196 server using 24 cores, in less than 24 hours. RAM requirement is estimated to be close to 0.5Gb per
197 million paired-end reads. The assemblies on average contained 176k transcripts (range 19k - 643k) and
198 97Mb (range 14MB - 198Mb). Other quality metrics will be discussed below, specifically in relation to
199 the ORP produced assemblies.

200 ORP assemblies generally completed on a standard Linux server using 24 cores in three days. Typically
201 Trinity was the longest running assembler, with the individual SPAdes assemblies being the shortest.
202 RAM requirement is estimated to be 1.5Gb - 2Gb per million paired-end reads, with SPAdes requiring
203 the most. The assemblies on average contained 153k transcripts (range 23k - 625k) and 64Mb (range
204 8MB - 181Mb).

205 MinHash sketch signatures (Ondov et al., 2016) of each assemblies of a given dataset were calculated
206 using sourmash (Titus Brown and Irber, 2016), and a MDS plot was generated (Figure 1) from their
207 distances. Interestingly, each assembler tends to produce a specific signature which is relatively consistent
208 between the fifteen datasets. Shannon differentiates itself from the other assemblers on the first (x) MDS
209 axis, while the other assemblers (SPAdes and Trinity) are separated on the second (y) MDS axis.

Figure 1

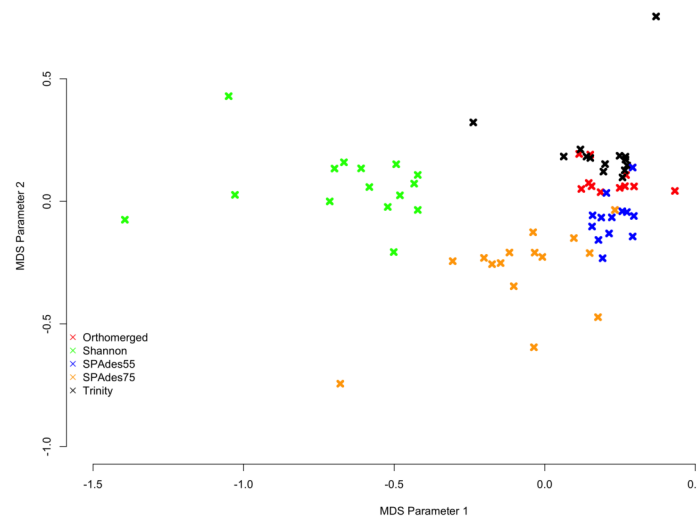


Figure 1. MDS plot describing the similarity within and between assemblers. Colored x's mark individual assemblies, with red marks corresponding to the ORP assemblies, green marks corresponding to the Shannon assemblies, blue marks corresponding to the SPAdes55 assemblies, orange marks corresponding to the SPAdes75 assemblies, and the black marks corresponding to the Trinity assemblies. In general assemblies produced by a given assembler tend to cluster together.

3.1.1 Assembly Structure

The structural integrity of each assembly was evaluated using the *TransRate* and *Detonate* software packages. As many downstream applications depend critically on accurate read mapping, assembly quality is correlated with increased mapping rates. The split violin plot presented in figure 2A visually represents the mapping rates of each assembly, with lines connecting the mapping rates of datasets assembled with *Trinity* and with the ORP, respectively. The average mapping rate of the *Trinity* assembled datasets was 87% (sd = 8%), while the average mapping rates of the ORP assembled datasets was 93% (sd=4%). This test is statistically significant (one-sided Wilcoxon rank sum test, $p = 2E-2$). Mapping rates of the other assemblies are less than that of the ORP assembly, but in most cases, greater than that of the *Trinity* assembly. This aspect of assembly quality is critical. Specifically mapping rates measure how representative the assembly is of the reads. If I assume that the vast majority of generated reads come from the biological sample under study, when reads fail to map, that fraction of the biology is lost from all downstream analysis and inference. This study demonstrates that across a wide variety of taxa, assembling RNAseq reads with any single assembler alone may result in a decrease in mapping rate and in turn, the lost ability to draw conclusions from that fraction of the sample.

Figure 2B describes the distribution of *TransRate* assembly scores, which is a synthetic metric taking into account the quality of read mapping and coverage-based statistics. The *Trinity* assemblies had an average optimal score of 0.35 (sd = .14), while the ORP assembled datasets had an average score of 0.46 (sd = .07). This test is statistically significant (one-sided Wilcoxon rank sum test, p -value = $1.8E-2$). Optimal scores of the other assemblies are less than that of the ORP assembly, but in most cases, greater than that of the *Trinity* assembly. Figure 2C describes the distribution of *Detonate* scores. The *Trinity* assemblies had an average score of $-6.9E9$ (sd = $5.2E9$), while the ORP assembled datasets had an average score of $-5.3E9$ (sd = $3.5E9$). This test not is statistically significant, though in all cases, relative to all other assemblies, scores of the ORP assemblies are improved (become less negative), indicating that the ORP produced assemblies of higher quality.

In addition to reporting synthetic metrics related to assembly structure, *TransRate* reports individual metrics related to specific elements of assembly quality. One such metric estimates the rate of chimerism, a phenomenon which is known to be problematic in *de novo* assembly (Ungaro et al., 2017; Singhal, 2013). Rates of chimerism are relatively constant between all assemblers, ranging from 10% for the Shannon assembly, to 12% for the SPAdes75 assembly. The chimerism rate for the ORP assemblies

averaged 10.5% (\pm 4.7%). While the new method would ideally improve this metric by exclusively selecting non-chimeric transcripts, this does not seem to be the case, and may be related to the inherent shortcomings of short-read transcriptome assembly.

Of note, consistent with all short-read assemblers (Ungaro et al., 2017), the ORP assemblies may not accurately reflect the true isoform complexity. Specifically, because of the way that single representative transcripts are chosen from a cluster of related sequences, some transcriptional complexity may be lost. Consider the cluster containing contigs {AB, A, B} where AB is a false-chimera, selecting a single representative transcript with the best score could yield either A or B, thereby excluding an important transcript in the final output. I believe this type of transcript loss is not common, based on how contigs are scored (Table 1, Figure 3, (Smith-Unna et al., 2016)), though strict demonstration of this is not possible, given the lack of high-quality reference genomes for the majority of the datasets. More generally, mapping rates, Detonate and TransRate score improvements suggest that this type of loss is not widespread.

Figure 2

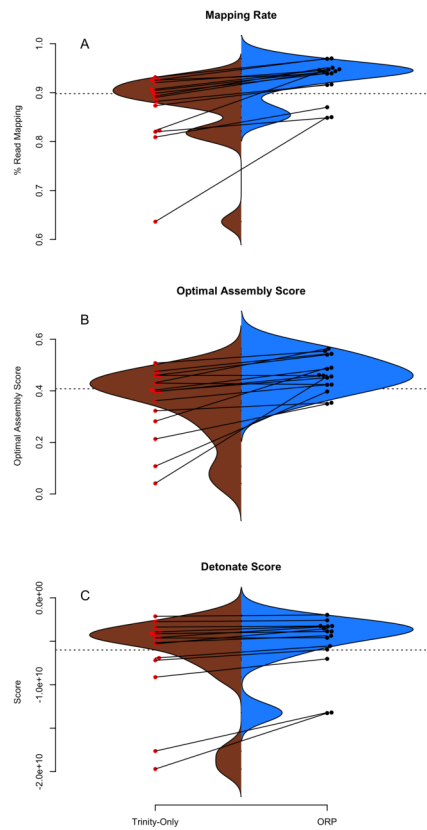


Figure 2. TransRate and Detonate generated statistics. Split violin plots depict the relationship between Trinity assemblies (brown color) and ORP produced assemblies (blue color). Red and black dots indicate the value of a given metric for each assembly. Lines connecting the red and black dots connect datasets assembled via the two methods.

3.1.2 Assembly Content

The genic content of assemblies was measured using the software package Shmblast, which implements the conditional reciprocal blast test against the Swiss-prot database. Presented in Table 2 and in Figure 3A, ORP assemblies recovered on average 13364 (sd=3391) blast hits, while all other assemblies recovered fewer (minimum Shannon, mean=10299). In every case across all assemblers, the ORP assembler retained more reciprocal blast hits, though only the comparison between the ORP assembly and Shannon was significant (one-sided Wilcoxon rank sum test, $p = 4E-3$). Notably, in all cases, each assembler was both missing transcripts contained in other assemblies, and contributed unique transcripts to the final merged assembly (Table 2), highlighting the utility of using multiple assemblers.

Table 2

276

Assembly	Genes	Delta	Unique
Concatenated	14674 ± 3590		
SPAdes55		-1739 ± 758	570 ± 266
SPAdes75		-2711 ± 2047	301 ± 195
Shannon		-4375 ± 3508	302 ± 241
Trinity		-1952 ± 803	520 ± 301

277
278
279
280
281
282

Table 2 describes the number of genes contained in the assemblies, with the row labeled concatenated representing the combined average (\pm standard deviation) number of genes contained in all assemblies of a given dataset. The other rows contain information about each assembly. The column labeled delta contains the average number (\pm standard deviation) of genes missing, relative to the concatenated number. The unique column contains the average number of genes (\pm standard deviation) unique to that assembly.

283
284
285
286
287
288
289
290
291
292
293

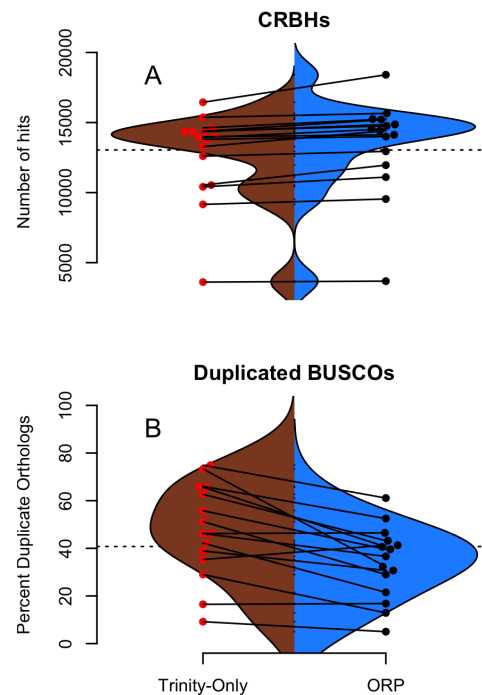
Regarding BUSCO scores, Trinity assemblies contained on average 86% (sd = 21%) of the full-length orthologs as defined by the BUSCO developers, while the ORP assembled datasets contained on average 86% (sd = 13%) of the full length transcripts. Other assemblers contained fewer full-length orthologs. The Trinity and ORP assemblies were missing, on average 4.5% (sd = 8.7%) of orthologs. The Trinity assembled datasets contained 9.5% (sd = 17%) of fragmented transcripts while the ORP assemblies each contained on average 9.4% (sd = 9%) of fragmented orthologs. The other assemblers in all cases contained more fragmentation. The rate of transcript duplication, depicted in figure 3B is 47% (sd = 20%) for Trinity assemblies, and 34% (sd = 15%) for ORP assemblies. This result is statistically significant (One sided Wilcoxon rank sum test, p-value = 0.02). Of note, all other assemblers produce less transcript duplication than does the ORP assembly, but none of these differences arise to the level of statistical significance.

294
295
296
297
298
299
300

While the majority of the BUSCO metrics were unchanged, the number of orthologs recovered in duplicate (>1 copy), was decreased when using the ORP. This difference is important, given that the relative frequency of transcript duplication may have important implications for downstream abundance estimation, with less duplication potentially resulting in more accurate estimation. Although gene expression quantitation software (Patro et al., 2017; Bray et al., 2016) probabilistically assigns reads to transcripts in an attempt at mitigating this issue, a primary solution related to decreasing artificial transcript duplication could offer significant advantages.

301

Figure 3



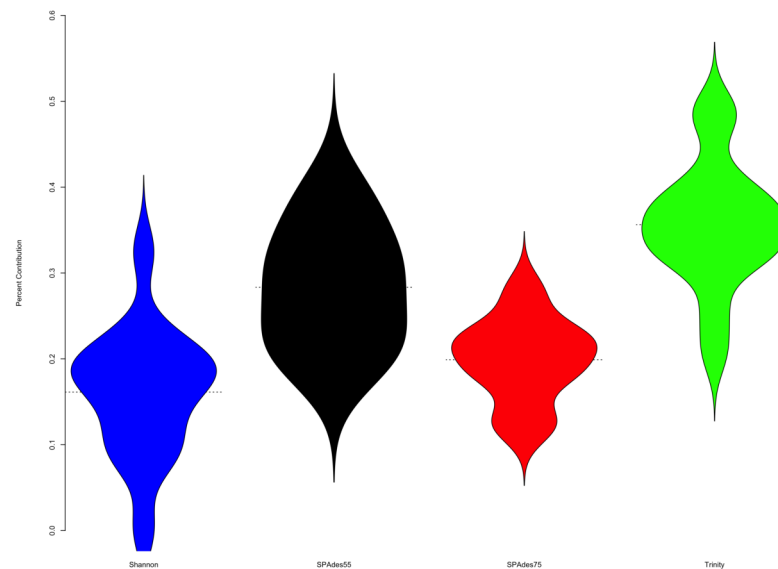
302

303 Figure 3. Shmblast and BUSCO generated statistics. Split violin plots depict the relationship between
 304 Trinity assemblies (brown color) and ORP produced assemblies (blue color). Red and black dots
 305 indicate the value of a given metric for each assembly. Lines connecting the red and black dots connect
 306 datasets assembled via the two methods.

307 3.1.3 Assembler Contributions

308 To understand the relative contribution of each assembler to the final merged assembly produced by the
 309 Oyster River Protocol, I counted the number of transcripts in the final merged assembly that originated
 310 from a given assembler (Figure 4). On average, 36% of transcripts in the merged assembly were produced
 311 by the Trinity assembler. 16% were produced by Shannon. SPAdes run with a kmer value of
 312 length=55 produced 28% of transcripts, while SPAdes run with a kmer value of length=75 produced
 313 20% of transcripts

314 Figure 4

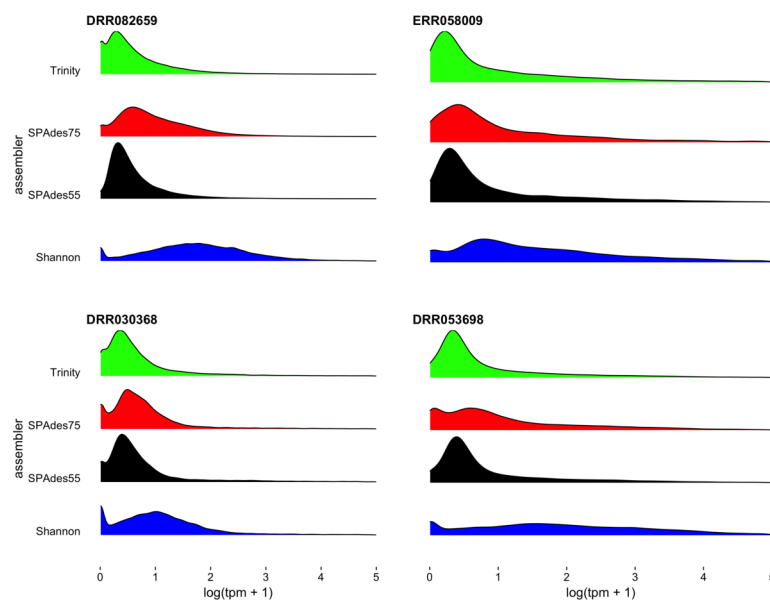


315

316 Figure 4 depicts the proportion of the final transcripts contained in the merged assembly that are a
 317 product of each assembler. Violin plots illustrate that Shannon contributes on average the fewest number
 318 of transcripts (<20 % of transcripts) to the final merged assembly, while Trinity contributes on average
 319 the most. Small dashed lines on each side of the plot mark the median of the distribution.

320 To further understand the potential biases intrinsic to each assembler, I plotted the distribution of
 321 gene expression estimates for each merged assembly, broken down by the assembler of origin (Figure 5,
 322 depicting four randomly selected representative assemblies). As is evident, most transcripts are lowly
 323 expressed, with SPAdes and Trinity both doing a sufficient job in reconstructing these transcripts.
 324 Of note, the SPAdes assemblies using kmer-length=75 is biased, as expected, towards more highly
 325 expressed transcripts relative to kmer-length 55 assemblies. Shannon demonstrates a unique profile,
 326 consisting of, almost exclusively high-expression transcripts, showing a previously undescribed bias
 327 against low-abundance transcripts. These differences may reflect a set of assembler-specific heuristics
 328 which translate into differential recovery of distinct fractions of the transcript community. Figure 5 and
 329 Table 2 describe the outcomes of these processes in terms of transcript recovery. Taken together, these
 330 expression profiles suggest a mechanism by which the ORP outperforms single-assembler assemblies.
 331 While there is substantial overlap in transcript recovery, each assembler recovers unique transcripts (Table
 332 2 and Figure 5) based on expression (and potentially other properties), which when merged together into
 333 a final assembly, increases the completeness

334 **Figure 5**



335

336 Figure 5 depicts the density distribution of gene expression ($\log(\text{TPM}+1)$), broken down by individual
 337 assembly, for four representative datasets. As predicted, the use of a higher kmer value with the SPAdes
 338 assembler resulted in biasing reconstruction towards more highly expressed transcripts. Interestingly,
 339 Shannon uniquely exhibits a bias towards the reconstruction of high-expression transcripts (or away
 340 from low-abundance transcripts).

341 3.2 Quality is independent of read depth

342 This study included read datasets of a variety of sizes. Because of this, I was interested in understanding
 343 if the number of reads used in assembly was strongly related to the quality of the resultant assembly.
 344 Conclusively, this study demonstrates that between 30 million paired-end reads and 200 million paired-end
 345 reads, no strong patterns in quality are evident (Figure 6). This finding is in line with previous work,
 346 (MacManes, 2015) suggesting that assembly metrics plateau at between 20M and 40M read pairs, with
 347 sequencing beyond this level resulting in minimal gain in performance.

348 **Figure 6**

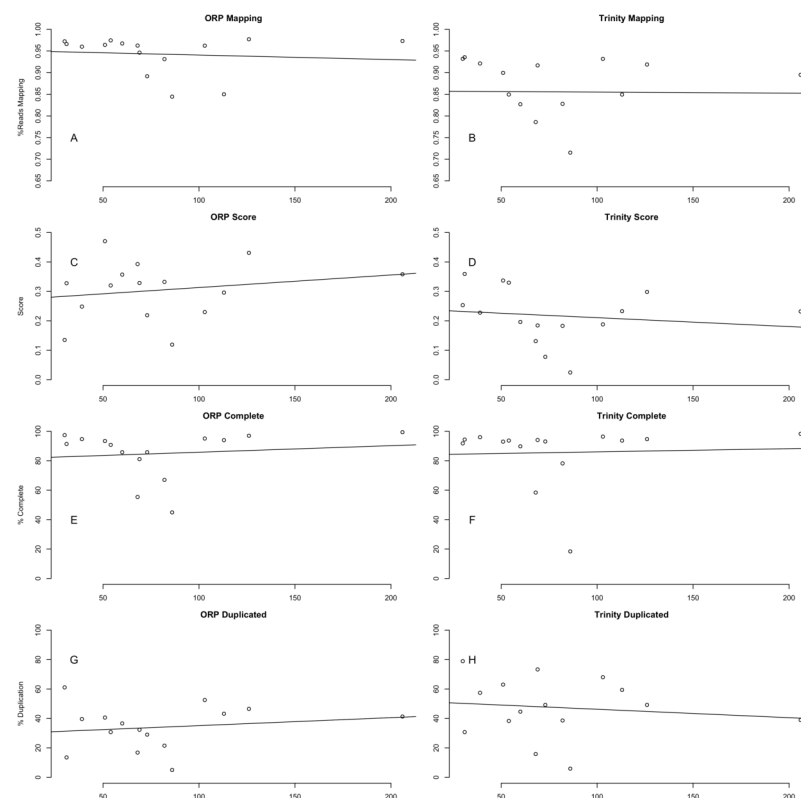


Figure 6 depicts the relationship between a subset of assembly metrics and the number of read pairs. There is no significant relationship. In all cases the x-axis is millions of paired-end reads.

REFERENCES

- Aubry, S., Kelly, S., Kümpers, B. M. C., Smith-Unna, R. D., and Hibberd, J. M. (2014). Deep Evolutionary Comparison of Gene Expression Identifies Parallel Recruitment of Trans-Factors in Two Independent Origins of C4 Photosynthesis. *PLOS Genetics*, 10(6):e1004365.
- Bolger, A. M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*, 30(15):btu170–2120.
- Bray, N. L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic RNA-seq quantification. *Nature Biotechnology*, 34(5):525–527.
- Cahoy, J. D., Emery, B., Kaushal, A., Foo, L. C., Zamanian, J. L., Christopherson, K. S., Xing, Y., Lubischer, J. L., Krieg, P. A., Krupenko, S. A., Thompson, W. J., and Barres, B. A. (2008). A transcriptome database for astrocytes, neurons, and oligodendrocytes: a new resource for understanding brain development and function. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 28(1):264–278.
- Chikhi, R. and Medvedev, P. (2014). Informed and automated k-mer size selection for genome assembly. *Bioinformatics*, 30(1):31–37.
- Emms, D. M. and Kelly, S. (2015). OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biology*, 16(1):157.
- Finseth, F. R. and Harrison, R. G. (2014). A comparison of next generation sequencing technologies for transcriptome assembly and utility for RNA-Seq in a non-model bird. *PloS one*, 9(10):e108550.
- Fitzpatrick, M., Ben-Shahar, Y., Vet, L., Smid, H., Robinson, G. E., and Sokolowski, M. (2005). Candidate genes for behavioural ecology. *Trends In Ecology & Evolution*, 20(2):96–104.
- Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., Couger, M. B., Eccles, D., Li, B., Lieber, M., MacManes, M. D., Ott, M., Orvis, J., Pochet, N., Strozzi, F., Weeks, N., Westerman, R., William, T., Dewey, C. N., Henschel, R., Leduc, R. D., Friedman, N., and Regev, A. (2013). *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nature Protocols*, 8(8):1494–1512.

- 378 Jackman, S. D. and Birol, I. (2016). Linuxbrew and Homebrew for cross-platform package manage-
379 ment [version 1; not peer reviewed]. In *F1000*.
- 380 Jiang, H., Lei, R., Ding, S.-W., and Zhu, S. (2014). Skewer: a fast and accurate adapter trimmer for
381 next-generation sequencing paired-end reads. *BMC Bioinformatics*, 15(1):182.
- 382 Kampstra, P. (2008). Beanplot: A boxplot alternative for visual comparison of distributions.
- 383 Kannan, S., Hui, J., Mazooji, K., Pachter, L., and Tse, D. (2016). Shannon: An Information-Optimal de
384 Novo RNA-Seq Assembler. *bioRxiv*.
- 385 Lappalainen, T., Sammeth, M., Friedländer, M. R., t Hoen, P. A. C., Monlong, J., Rivas, M. A., González-
386 Porta, M., Kurbatova, N., Griebel, T., Ferreira, P. G., Barann, M., Wieland, T., Greger, L., van Iterson,
387 M., Almlöf, J., Ribeca, P., Pulyakhina, I., Esser, D., Giger, T., Tikhonov, a., Sultan, M., Bertier, G.,
388 MacArthur, D. G., Lek, M., Lizano, E., Buermans, H. P. J., Padiou, I., Schwarzmayr, T., Karlberg,
389 O., Ongen, H., Kilpinen, H., Beltran, S., Gut, M., Kahlem, K., Amstislavskiy, V., Stegle, O., Pirinen,
390 M., Montgomery, S. B., Donnelly, P., McCarthy, M. I., Flicek, P., Strom, T. M., Geuvadis Consortium,
391 Lehrach, H., Schreiber, S., Sudbrak, R., Carracedo, A., Antonarakis, S. E., Häslér, R., Syvänen, A.-C.,
392 van Ommen, G.-J., Brazma, A., Meitinger, T., Rosenstiel, P., Guigo, R., Gut, I. G., Estivill, X., and
393 Dermitzakis, E. T. (2013). Transcriptome and genome sequencing uncovers functional variation in
394 humans. *Nature*, 501(7468):506–511.
- 395 Le, H. S., Schulz, M. H., McCauley, B. M., Hinman, V. F., and Bar-Joseph, Z. (2013). Probabilistic error
396 correction for RNA sequencing. *Nucleic Acids Research*, 41(10):1–11.
- 397 Li, B., Fillmore, N., Bai, Y., Collins, M., Thomson, J. A., Stewart, R., and Dewey, C. N. (2014). Evaluation
398 of *de novo* transcriptome assemblies from RNA-Seq data. *Genome Biology*, 15(12):663–21.
- 399 Li, R., Li, Y., Kristiansen, K., and Wang, J. (2008). SOAP: short oligonucleotide alignment program.
400 *Bioinformatics*, 24(5):713–714.
- 401 Li, X., Kim, Y., Tsang, E. K., Davis, J. R., Damani, F. N., Chiang, C., Hess, G. T., Zappala, Z.,
402 Strober, B. J., Scott, A. J., Li, A., Ganna, a., Bassik, M. C., Merker, J. D., GTEx Consortium,
403 Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis Working Group, Statistical
404 Methods groups—Analysis Working Group, Enhancing GTEx (eGTEx) groups, NIH Common
405 Fund, NIH/NCI, NIH/NHGRI, NIH/NIMH, NIH/NIDA, Biospecimen Collection Source Site—NDRI,
406 Biospecimen Collection Source Site—RPCI, Biospecimen Core Resource—VARI, Brain Bank Reposi-
407 tory—University of Miami Brain Endowment Bank, Leidos Biomedical—Project Management, ELSI
408 Study, Genome Browser Data Integration & Visualization—EBI, Genome Browser Data Integration
409 & Visualization—UCSC Genomics Institute, University of California Santa Cruz, Hall, I. M., Battle, A.,
410 and Montgomery, S. B. (2017). The impact of rare variation on gene expression across tissues. *Nature*,
411 550(7675):239–243.
- 412 Liu, J., Li, G., Chang, Z., Yu, T., Liu, B., McMullen, R., Chen, P., and Huang, X. (2016). BinPacker:
413 Packing-Based De Novo Transcriptome Assembly from RNA-seq Data. *PLOS Computational Biology*,
414 12(2):e1004772.
- 415 Love, M. I., Huber, W., and anders, S. (2014). Moderated estimation of fold change and dispersion for
416 RNA-seq data with DESeq2. *Genome Biology*, 15(12):550.
- 417 MacManes, M. D. (2014). On the optimal trimming of high-throughput mRNA sequence data. *Frontiers*
418 *in Genetics*, 5:13.
- 419 MacManes, M. D. (2015). Establishing evidenced-based best practice for the *de novo* assembly and
420 evaluation of transcriptomes from non-model organisms. Technical report.
- 421 MacManes, M. D. and Eisen, M. B. (2013). Improving transcriptome assembly through error correction
422 of high-throughput sequence reads. *PeerJ*, 1:e113.
- 423 Marchant, A., Mougél, F., Mendonça, V., Quartier, M., Jacquin-Joly, E., da Rosa, J. A., Petit, E., and Harry,
424 M. (2016). Comparing *de novo* and reference-based transcriptome assembly strategies by applying them
425 to the blood-sucking bug *Rhodnius prolixus*. *Insect Biochemistry and Molecular Biology*, 69:25–33.
- 426 Martin, M. and 2011 (2011). Cutadapt removes adapter sequences from high-throughput sequencing
427 reads. *journal.embnet.org*, 17(1):10.
- 428 Moreton, J., Izquierdo, A., and Emes, R. D. (2015). Assembly, Assessment, and Availability of *De novo*
429 Generated Eukaryotic Transcriptomes. *Frontiers in Genetics*, 6(217):361.
- 430 Mortazavi, A., Williams, B. A., Mccue, K., Schaeffer, L., and Wold, B. (2008). Mapping and quantifying
431 mammalian transcriptomes by RNA-Seq. *Nature Methods*, 5(7):621–628.
- 432 Ondov, B. D., Treangen, T. J., Melsted, P., Mallonee, A. B., Bergman, N. H., Koren, S., and Phillippy,

- 433 A. M. (2016). Mash: fast genome and metagenome distance estimation using MinHash. *Genome*
434 *Biology*, 17(1):132.
- 435 Panhuis, T. M. (2006). Molecular evolution and population genetic analysis of candidate female reproduc-
436 tive genes in *Drosophila*. *Genetics*, 173(4):2039–2047.
- 437 Patro, R., Duggal, G., Love, M. I., Irizarry, R. A., and Kingsford, C. (2017). Salmon provides fast and
438 bias-aware quantification of transcript expression. *Nature Methods*, 14(4):417–419.
- 439 Pell, J., Hintze, A., Canino-Koning, R., Howe, A., Tiedje, J. M., and Brown, C. T. (2012). Scaling
440 metagenome sequence assembly with probabilistic *de Bruijn* graphs. *Proceedings of the National*
441 *Academy of Sciences*, 109(33):13272–13277.
- 442 Peng, Y., Peng, Y., Leung, H. C. M., Leung, H. C. M., Yiu, S.-M., Lv, M.-J., Lv, M.-J., Zhu, X.-G., Zhu,
443 X.-G., Chin, F. Y. L., and Chin, F. Y. L. (2013). IDBA-tran: a more robust *de novo de Bruijn* graph
444 assembler for transcriptomes with uneven expression levels. *Bioinformatics*, 29(13):i326–i334.
- 445 R Core Development Team, F. (2011). R: A Language and Environment for Statistical Computing.
- 446 Robertson, G., Schein, J., Chiu, R., Corbett, R., Field, M., Jackman, S. D., Mungall, K., Lee, S., Okada,
447 H. M., Qian, J. Q., Griffith, M., Raymond, A., Thiessen, N., Cezard, T., Butterfield, Y. S., Newsome,
448 R., Chan, S. K., She, R., Varhol, R., Kamoh, B., Prabhu, A.-L., Tam, A., Zhao, Y., Moore, R. A., Hirst,
449 M., Marra, M. A., Jones, S. J. M., Hoodless, P. A., and Birol, I. (2010). *De novo* assembly and analysis
450 of RNA-seq data. *Nature Methods*, 7(11):909–912.
- 451 Robinson, M. D., McCarthy, D. J., and Smyth, G. K. (2010). edgeR: a Bioconductor package for
452 differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140.
- 453 Schulz, M. H., Zerbino, D. R., Vingron, M., and Birney, E. (2012). Oases: robust *de novo* RNA-seq
454 assembly across the dynamic range of expression levels. *Bioinformatics*, 28(8):1086–1092.
- 455 Scott, C. (2017). shmlast: An improved implementation of Conditional Reciprocal Best Hits with LAST
456 and Python. *The Journal of Open Source Software*, 2(9):1–4.
- 457 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., and Zdobnov, E. M. (2015). BUSCO:
458 assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*,
459 31(19):3210–3212.
- 460 Simpson, J. T., Wong, K., Jackman, S. D., Schein, J. E., Jones, S. J. M., and Birol, I. (2009). ABySS: A
461 parallel assembler for short read sequence data. *Genome Research*, 19(6):1117–1123.
- 462 Singhal, S. (2013). *De novotranscriptomic* analyses for non-model organisms: an evaluation of methods
463 across a multi-species data set. *Molecular Ecology Resources*, 13(3):n/a–n/a.
- 464 Smith-Unna, R., Bournsnel, C., Patro, R., Hibberd, J. M., and Kelly, S. (2016). TransRate: reference-free
465 quality assessment of *de novo* transcriptome assemblies. *Genome Research*, 26(8):1134–1144.
- 466 Song, L. and Florea, L. (2015). Rcorrector: efficient and accurate error correction for Illumina RNA-seq
467 reads. *GigaScience*, 4(1):48.
- 468 Tan, M. H., Li, Q., Shanmugam, R., Piskol, R., Kohler, J., Young, A. N., Liu, K. I., Zhang, R., Ra-
469 maswami, G., Ariyoshi, K., Gupte, A., Keegan, L. P., George, C. X., Ramu, A., Huang, N., Pollina,
470 E. A., Leeman, D. S., Rustighi, A., Goh, Y. P. S., GTEx Consortium, Laboratory, Data Analysis
471 & Coordinating Center (LDACC)—Analysis Working Group, Statistical Methods groups—Analysis
472 Working Group, Enhancing GTEx (eGTEx) groups, NIH Common Fund, NIH/NCI, NIH/NHGRI,
473 NIH/NIMH, NIH/NIDA, Biospecimen Collection Source Site—NDRI, Biospecimen Collection Source
474 Site—RPCI, Biospecimen Core Resource—VARI, Brain Bank Repository—University of Miami Brain
475 Endowment Bank, Leidos Biomedical—Project Management, ELSI Study, Genome Browser Data Inte-
476 gration & Visualization—EBI, Genome Browser Data Integration & Visualization—UCSC Genomics
477 Institute, University of California Santa Cruz, Chawla, A., Del Sal, G., Peltz, G., Brunet, A., Conrad,
478 D. F., Samuel, C. E., O’Connell, M. A., Walkley, C. R., Nishikura, K., and Li, J. B. (2017). Dynamic
479 landscape and regulation of RNA editing in mammals. *Nature*, 550(7675):249–254.
- 480 Titus Brown, C. and Irber, L. (2016). sourmash: a library for MinHash sketching of DNA. *The Journal of*
481 *Open Source Software*, 1(5):27–1.
- 482 Ungaro, A., Pech, N., Martin, J.-F., McCairns, R. J. S., Mévy, J.-P., Chappaz, R., and Gilles, a. (2017).
483 Challenges and advances for transcriptome assembly in non-model species. *PloS one*, 12(9):e0185020–
484 21.
- 485 Vijay, N., Poelstra, J. W., Künstner, A., and Wolf, J. B. W. (2013). Challenges and strategies in
486 transcriptome assembly and differential gene expression quantification. A comprehensive *in silico*
487 assessment of RNA-seq experiments. *Molecular Ecology*, 22(3):620–634.

488 Wang, S. and Gribskov, M. (2017). Comprehensive evaluation of *de novo* transcriptome assembly
489 programs and their effects on differential gene expression analysis. *Bioinformatics*, 33(3):327–333.
490 Wang, Z., Wang, Z., Gerstein, M., and Snyder, M. (2009). RNA-Seq: a revolutionary tool for transcrip-
491 tomics. *Nature Reviews Genetics*, 10(1):57–63.
492 Wolf, J. B. W. (2013). Principles of transcriptome analysis and gene expression quantification: an
493 RNA-seq tutorial. *Molecular Ecology Resources*, 13(4):559–572.
494 Xie, Y., Wu, G., Tang, J., Luo, R., Patterson, J., Liu, S., Huang, W., He, G., Gu, S., Li, S., Zhou, X., Li, Y.,
495 Xu, X., Wong, G. K.-S., and Wang, J. (2014). SOAPdenovo-Trans: *de novo* transcriptome assembly
496 with short RNA-Seq reads. *Bioinformatics*, 30(12):1660–1666.
497 Yang, X., Dorman, K. S., and Aluru, S. (2010). Reptile: representative tiling for short read error correction.
498 *Bioinformatics*, 26(20):2526–2533.
499 Yang, Y. and Smith, S. A. (2013). Optimizing *de novo* assembly of short-read RNA-seq data for
500 phylogenomics. *BMC Genomics*, 14:328.
501 Zerbino, D. R. and Birney, E. (2008). Velvet: Algorithms for *de novo* short read assembly using *de Bruijn*
502 graphs. *Genome Research*, 18(5):821–829.

Figure 1

MDS plot describing the similarity within and between assemblers.

Colored x's mark individual assemblies, with red marks corresponding to the ORP assemblies, green marks corresponding to the Shannon assemblies, blue marks corresponding to the SPAdes55 assemblies, orange marks corresponding to the SPAdes75 assemblies, and the black marks corresponding to the Trinity assemblies. In general assemblies produced by a given assembler tend to cluster together.

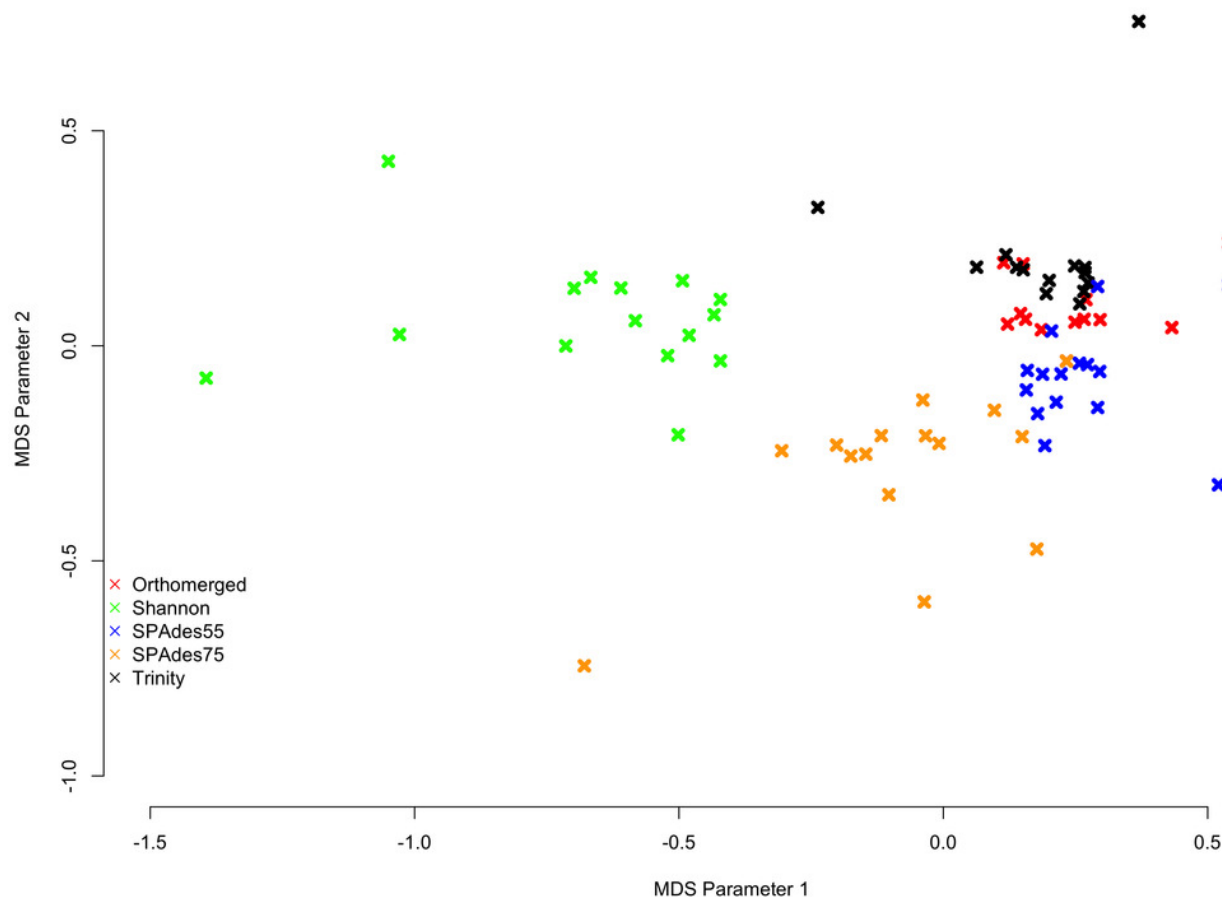


Figure 2

TransRate and Detonate generated statistics.

Split violin plots depict the relationship between Trinity assemblies (brown color) and ORP produced assemblies (blue color). Red and black dots indicate the value of a given metric for each assembly. Lines connecting the red and black dots connect datasets assembled via the two methods.

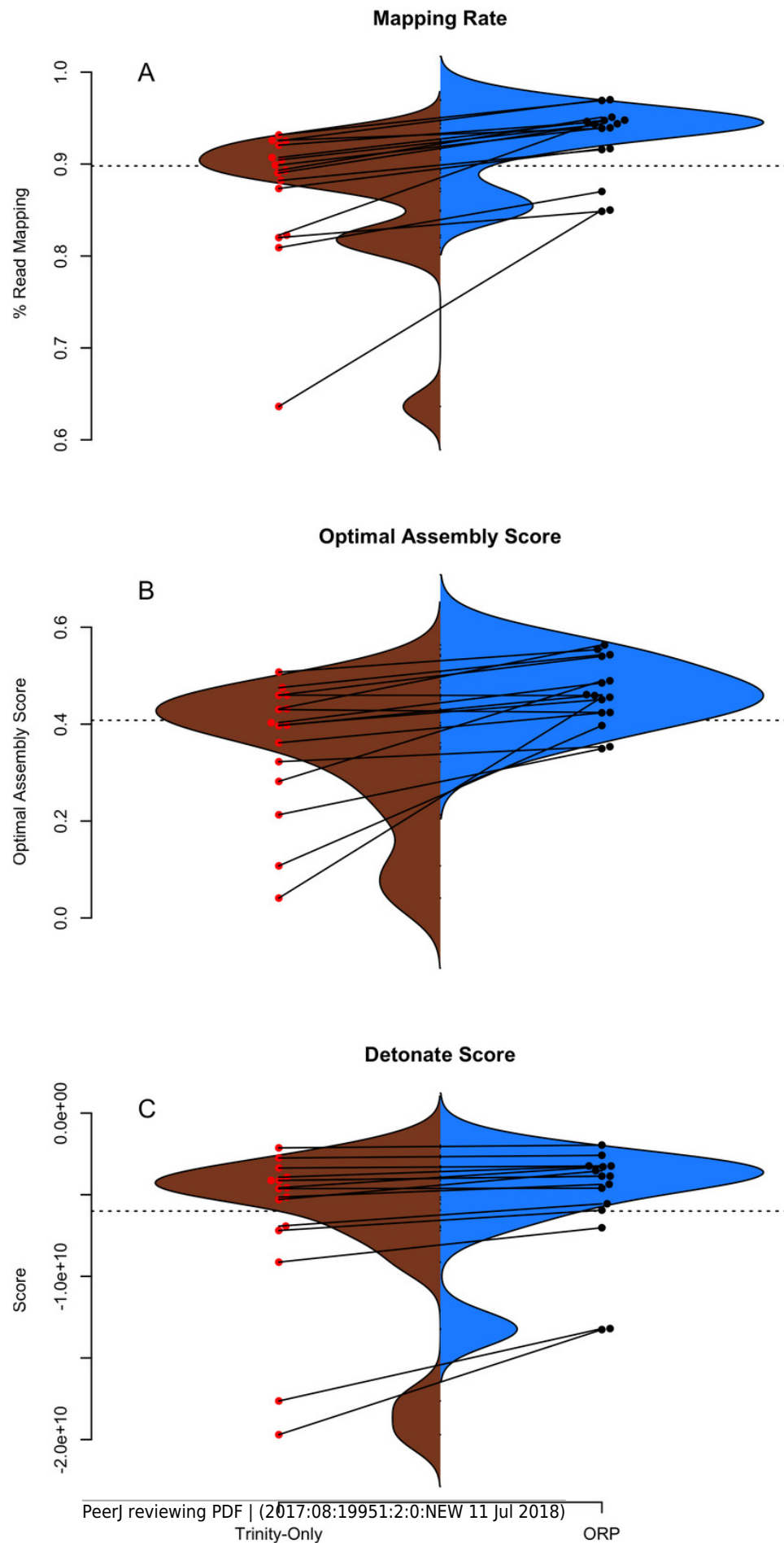


Figure 3

Shmlast and BUSCO generated statistics.

Split violin plots depict the relationship between Trinity assemblies (brown color) and ORP produced assemblies (blue color). Red and black dots indicate the value of a given metric for each assembly. Lines connecting the red and black dots connect datasets assembled via the two methods.

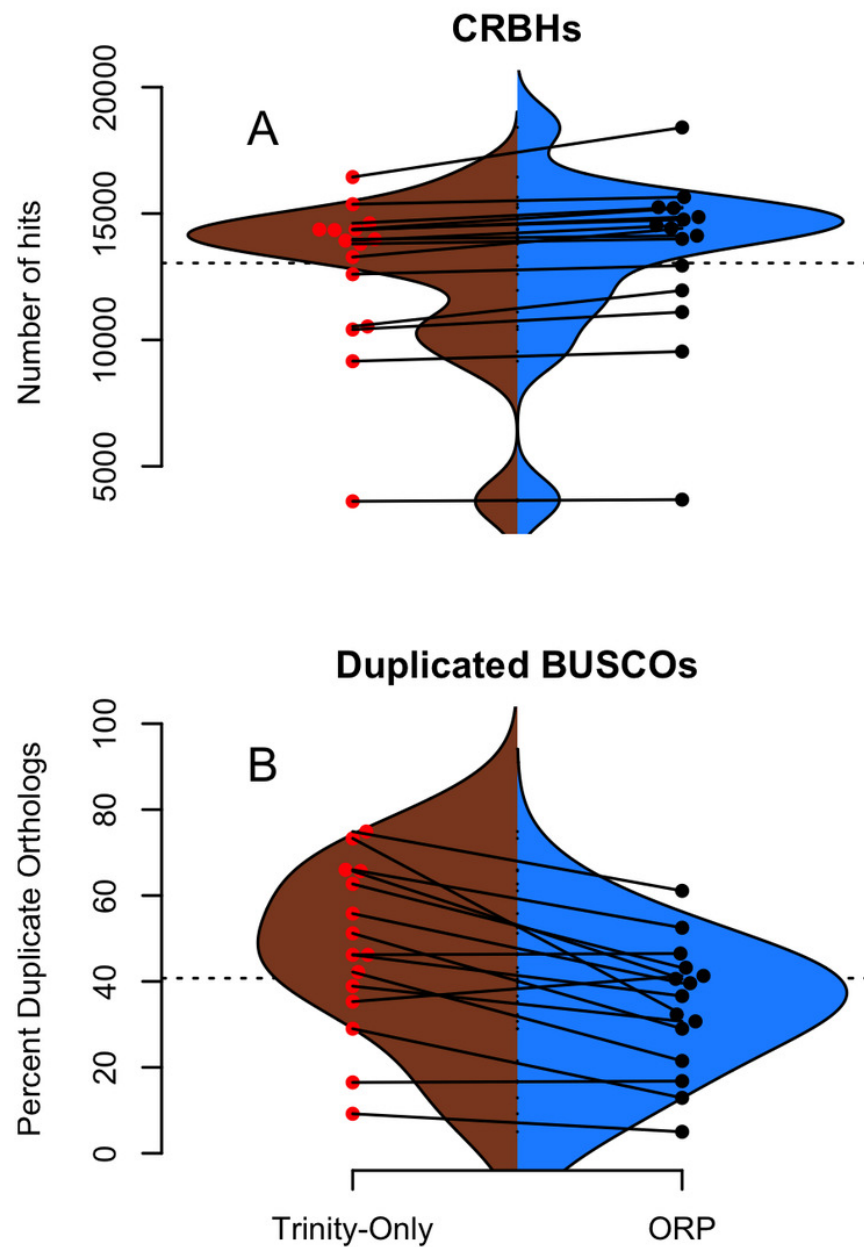


Figure 4

Plot describes the percent contribution of each assembler to the final ORP assembly.

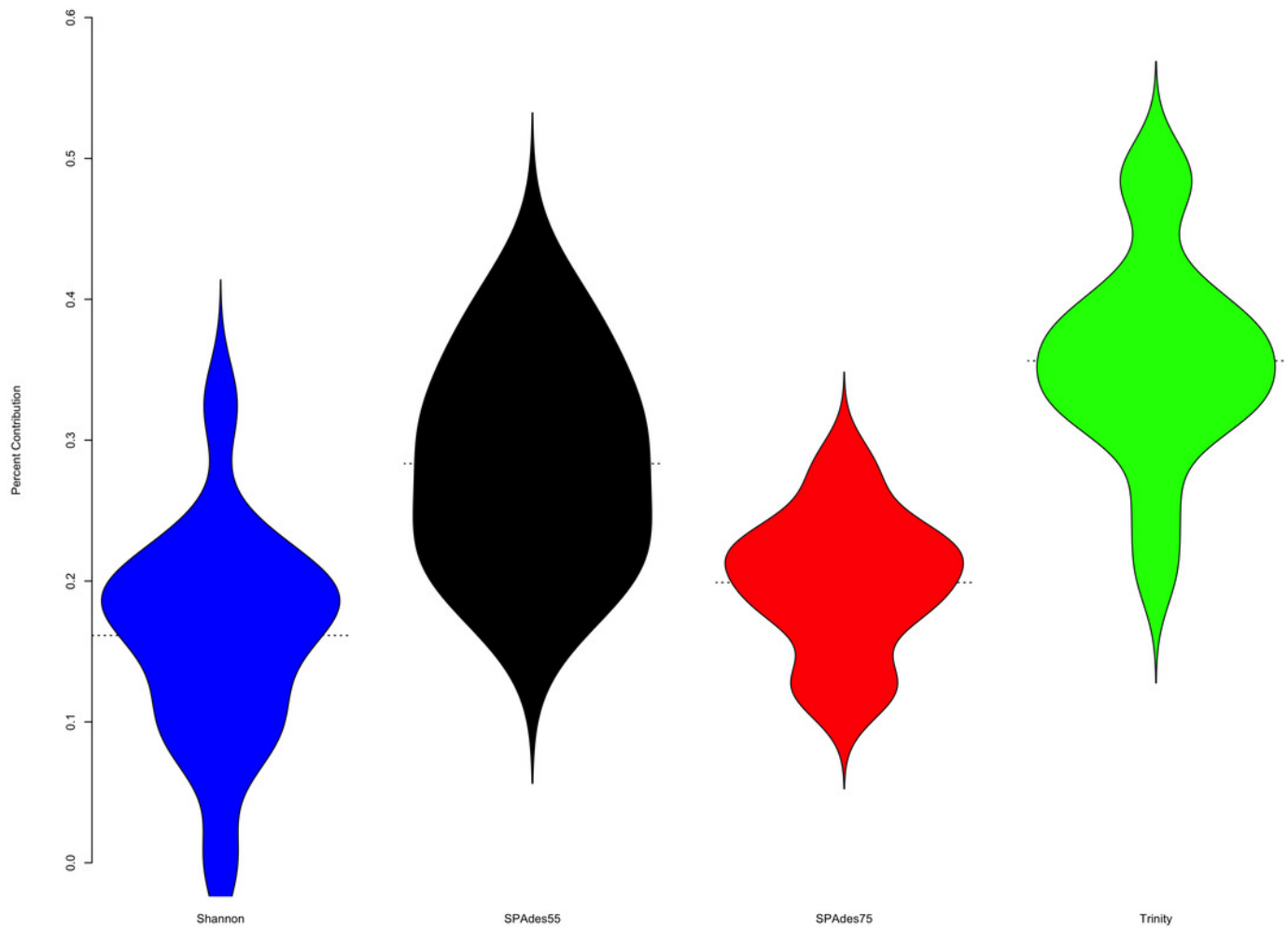


Figure 5

Distribution of gene expression for each assembler

Figure 5 depicts the distribution of gene expression ($\log(\text{TPM}+1)$), broken down by individual assembly, for four representative datasets. As predicted, the use of a higher kmer value with the SPAdes assembler resulted in biasing reconstruction towards more highly expressed transcripts. Interestingly, Shannon uniquely exhibits a bias towards the reconstruction of high-expression transcripts (or away from low-abundance transcripts).

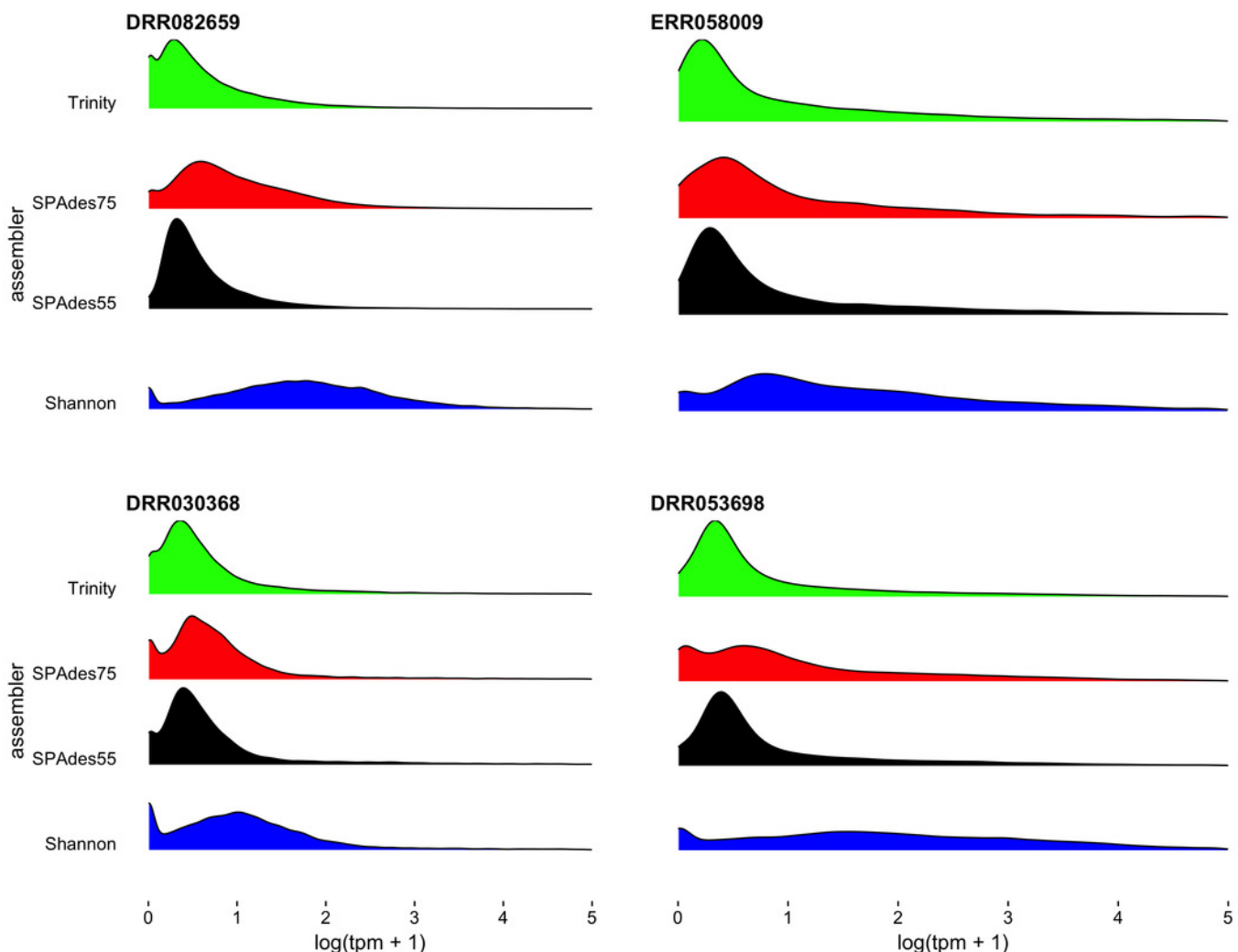


Figure 6

No relationship between metrics and dataset size

Figure 6 depicts the relationship between a subset of assembly metrics and the number of read pairs. There is no significant relationship. In all cases the x-axis is millions of paired-end reads.

