

# GigaScience

## Re-assembly, quality evaluation, and annotation of 678 microbial eukaryotic reference transcriptomes --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-18-00191R2	
<b>Full Title:</b>	Re-assembly, quality evaluation, and annotation of 678 microbial eukaryotic reference transcriptomes	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	Gordon and Betty Moore Foundation (GBMF4551)	Dr. C. Titus Brown
<b>Abstract:</b>	<p><b>Background</b> <i>De novo</i> transcriptome assemblies are required prior to analyzing RNAseq data from a species without an existing reference genome or transcriptome. Despite the prevalence of transcriptomic studies, the effects of using different workflows, or “pipelines”, on the resulting assemblies are poorly understood. Here, a pipeline was programmatically automated and used to assemble and annotate raw transcriptomic short read data collected by the Marine Microbial Eukaryotic Transcriptome Sequencing Project (MMETSP). The resulting transcriptome assemblies were evaluated and compared against assemblies that were previously generated with a different pipeline developed by the National Center for Genome Research (NCGR).</p> <p><b>Results</b> New transcriptome assemblies contained the majority of previous contigs as well as new content. On average, 7.8% of the annotated contigs in the new assemblies were novel gene names not found in the previous assemblies. Taxonomic trends were observed in the assembly metrics, with assemblies from the Dinoflagellata and Ciliophora phyla showing a higher percentage of open reading frames and number of contigs than transcriptomes from other phyla.</p> <p><b>Conclusions</b> Given current bioinformatics approaches, there is no single ‘best’ reference transcriptome for a particular set of raw data. As the optimum transcriptome is a moving target, improving (or not) with new tools and approaches, automated and programmable pipelines are invaluable for managing the computationally-intensive tasks required for re-processing large sets of samples with revised pipelines and ensuring a common evaluation workflow is applied to all samples. Thus, re-assembling existing data with new tools using automated and programmable pipelines may yield more accurate identification of taxon-specific trends across samples in addition to novel and useful products for the community.</p>	
<b>Corresponding Author:</b>	C. Titus Brown UC Davis Davis, CA UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	UC Davis	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Lisa Kristine Johnson	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Lisa Kristine Johnson Harriet Alexander C. Titus Brown	
<b>Order of Authors Secondary Information:</b>		
<b>Response to Reviewers:</b>	Minor changes as requested.	
<b>Additional Information:</b>		
<b>Question</b>	<b>Response</b>	

<p>Are you submitting this manuscript to a special series or article collection?</p>	<p>No</p>
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	<p>Yes</p>
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum</a></p>	<p>Yes</p>

[Standards Reporting Checklist?](#)

[Click here to view linked References](#)

GigaScience, 2018, 1–12

doi: [xx.xxxx/xxxx](#)Manuscript in Preparation  
Technical Note

## TECHNICAL NOTE

# Re-assembly, quality evaluation, and annotation of 678 microbial eukaryotic reference transcriptomes

Lisa K. Johnson<sup>1,2</sup>, Harriet Alexander<sup>1</sup> and C. Titus Brown<sup>1,2,3,\*</sup>

<sup>1</sup>Department of Population Health & Reproduction, School of Veterinary Medicine, University of California Davis and <sup>2</sup>Molecular, Cellular, and Integrative Physiology Graduate Group, University of California Davis and <sup>3</sup>Genome Center, University of California Davis

\*ctbrown@ucdavis.edu

## Abstract

**Background** *De novo* transcriptome assemblies are required prior to analyzing RNAseq data from a species without an existing reference genome or transcriptome. Despite the prevalence of transcriptomic studies, the effects of using different workflows, or “pipelines”, on the resulting assemblies are poorly understood. Here, a pipeline was programmatically automated and used to assemble and annotate raw transcriptomic short read data collected by the Marine Microbial Eukaryotic Transcriptome Sequencing Project (MMETSP). The resulting transcriptome assemblies were evaluated and compared against assemblies that were previously generated with a different pipeline developed by the National Center for Genome Research (NCGR). **Results** New transcriptome assemblies contained the majority of previous contigs as well as new content. On average, 7.8% of the annotated contigs in the new assemblies were novel gene names not found in the previous assemblies. Taxonomic trends were observed in the assembly metrics. Assemblies from the Dinoflagellata showed a higher number of contigs and unique k-mers than transcriptomes from other phyla while assemblies from Ciliophora had a lower percentage of open reading frames compared to other phyla. **Conclusions** Given current bioinformatics approaches, there is no single ‘best’ reference transcriptome for a particular set of raw data. As the optimum transcriptome is a moving target, improving (or not) with new tools and approaches, automated and programmable pipelines are invaluable for managing the computationally-intensive tasks required for re-processing large sets of samples with revised pipelines and ensuring a common evaluation workflow is applied to all samples. Thus, re-assembling existing data with new tools using automated and programmable pipelines may yield more accurate identification of taxon-specific trends across samples in addition to novel and useful products for the community.

**Key words:** marine microbial eukaryote; transcriptome assembly; automated pipeline; re-analysis

## Introduction

The analysis of gene expression from high-throughput nucleic acid sequence data relies on the presence of a high quality reference genome or transcriptome. When there is no reference genome or transcriptome for an organism of interest, raw RNA sequence data (RNAseq) must be assembled *de novo* into a transcriptome [1]. This type of analysis is ubiquitous across many fields, including: evolutionary developmental biology [2], cancer biology [3], agriculture [4, 5], ecological physiology [6, 7], and biological oceanography [8]. In recent years, substantial

investments have been made in data generation, primary data analysis, and development of downstream applications, such as biomarkers and diagnostic tools [9, 10, 11, 12, 13, 14, 15, 16]

Methods for *de novo* RNAseq assembly of the most common short read Illumina sequencing data continue to evolve rapidly, especially for non-model species [17]. At this time, there are several major *de novo* transcriptome assembly software tools available to choose from, including Trinity [18], SOAPdenovo-Trans [19], Trans-ABYSS [20], Oases [21], SPAdes [22], IDBA-tran [23], and Shannon [24]. The availability of these options stems from continued research into the unique computational

Compiled on: November 16, 2018.

Draft manuscript prepared by the author.

## Key Points

- Re-assembly with new tools can yield new results
- Automated and programmable pipelines can be used to process arbitrarily many samples.
- Analyzing many samples using a common pipeline identifies taxon-specific trends.

challenges associated with transcriptome assembly of short read Illumina RNAseq data, including large memory requirements, alternative splicing and allelic variants [18, 25],

The continuous development of new tools and workflows for RNAseq analysis combined with the vast amount of publicly available RNAseq data [26] raises the opportunity to re-analyze existing data with new tools. This, however, is rarely done systematically. To evaluate the performance impact of new tools on old data, we developed and applied a programmatically automated *de novo* transcriptome assembly workflow that is modularized and extensible based on the Eel Pond Protocol [27]. This workflow incorporates Trimmomatic [28], digital normalization with khmer software [29, 30], and the Trinity *de novo* transcriptome assembler [18].

To evaluate this pipeline, we re-analyzed RNAseq data from 678 samples generated as part of the Marine Microbial Eukaryotic Transcriptome Sequencing Project (MMETSP) [31]. The MMETSP data set was generated to broaden the diversity of sequenced marine protists to enhance our understanding of their evolution and roles in marine ecosystems and biogeochemical cycles [31, 32]. With data from species spanning more than 40 eukaryotic phyla, the MMETSP provides one of the largest publicly-available collections of RNAseq data from a diversity of species. Moreover, the MMETSP used a standardized library preparation procedure and all of the samples were sequenced at the same facility, making this data set unusually comparable.

Reference transcriptomes for the MMETSP were originally assembled by the National Center for Genome Research (NCGR) with a pipeline which used the Trans-ABYSS software program to assemble the short reads [31]. The transcriptomes generated from the NCGR pipeline have already facilitated discoveries in the evolutionary history of ecologically significant genes [33, 34], differential gene expression under shifting environmental conditions [8, 35], inter-group transcriptomic comparisons [36], unique transcriptional features [37, 38, 39], and meta-transcriptomic studies [34, 35, 36]

In re-assembling the MMETSP data, we sought to compare and improve the original MMETSP reference transcriptome and to create a platform which facilitates automated re-assembly and evaluation. Here, we show that our re-assemblies had better evaluation metrics and contained most of the NCGR contigs as well as adding new content.

## Methods

### Programmatically Automated Pipeline

An automated pipeline was developed to execute the steps of the Eel Pond mRNAseq Protocol [27], a lightweight protocol for assembling short Illumina RNAseq reads that uses the Trinity *de novo* transcriptome assembler. This protocol generates *de novo* transcriptome assemblies of acceptable quality [40]. The pipeline was used to assemble all of the data from the MMETSP (Figure 1). The code and instructions for running the pipeline are available at <https://doi.org/10.5281/zenodo.740440> [41].

The steps of the pipeline applied to the MMETSP are as follows:

### 1. Download the raw data

Raw RNA-seq data sets were obtained from the National Center for Biotechnology Information (NCBI) Sequence Read Archive (SRA) from [BioProject PRJNA231566](#). Data were paired-end (PE) Illumina reads with lengths of 50 bases for each read. A metadata (SraRunInfo.csv) file obtained from the SRA web interface was used to provide a list of samples to the `get_data.py` pipeline script, which was then used to download and extract fastq files from 719 records. The script uses the `fastq-dump` program from the SRA Toolkit to extract the SRA-formatted fastq files (version 2.5.4) [42]. There were 18 MMETSP samples with more than one SRA record (MMETSP0693, MMETSP1019, MMETSP0923, MMETSP0008, MMETSP1002, MMETSP1325, MMETSP1018, MMETSP1346, MMETSP0088, MMETSP0092, MMETSP0717, MMETSP0223, MMETSP0115, MMETSP0196, MMETSP0197, MMETSP0398, MMETSP0399, MMETSP0922). In these cases, reads from multiple SRA records were concatenated together per sample. Taking these redundancies into consideration, there were a total of 678 re-assemblies generated from the 719 records in PRJNA231566 (Supplemental Notebook 1 [43]). Assembly evaluation metrics were not calculated for MMETSP samples with more than one SRA record because these assemblies were different than the others, containing multiple samples, and thus not as comparable.

Initial transcriptomes that were assembled by the National Center for Genome Resources (NCGR), using methods and data described in the original publication [31], were downloaded from the iMicrobe repository to compare with our re-assemblies (<ftp://ftp.imicrobe.us/projects/104/>). There were two versions of each assembly, 'nt' and 'cds'. The version used for comparison is noted below in each evaluation step. To our knowledge, the NCGR took extra post-processing steps to filter content, leaving only coding sequences in the 'cds' versions of each assembly [31]

### 2. Perform quality control

Reads were analyzed with FastQC (v0.11.5) [44] and multiqc (version 1.2) [45] to confirm overall qualities before and after trimming. A conservative trimming approach [46] was used with Trimmomatic (version 0.33) [28] to remove residual Illumina adapters and cut bases off the start (LEADING) and end (TRAILING) of reads if they were below a threshold Phred quality score ( $Q < 2$ ).

### 3. Apply digital normalization

To decrease the memory requirements for each assembly, digital normalization was applied with the khmer software package (v2.0) prior to assembly [47]. First, reads were interleaved, normalized to a k-mer ( $k = 20$ ) coverage of 20 and a memory size of  $4e9$ , then low-abundance k-mers from reads with a coverage above 18 were trimmed. Orphaned reads, where the mated pair was removed during normalization, were added to the normalized reads.

### 4. Assemble

Transcriptomes were assembled from normalized reads with Trinity 2.2.0 using default parameters ( $k = 25$ ). This version of Trinity (v2.2.0) did not include the "in silico normaliza-

tion" option as a default parameter. The digital normalization approach we used with khmer is the same algorithm implemented in Trinity, but it requires less memory and is faster [48].

The resulting assemblies are referred to below as the "Lab for Data Intensive Biology" assemblies, or DIB assemblies. The original assemblies are referred to as the NCGR assemblies.

### 5. Post-assembly assessment

Transcriptomes were annotated using the dammit pipeline (version v1.0.dev0) [49], which relies on the following databases as evidence: Pfam-A (version 28.0) [50], Rfam (version 12.1) [51], OrthoDB (version 8) [52]. In the case where there were multiple database hits, one gene name per contig was selected by choosing the name of the lowest e-value match ( $<1e-05$ ).

All assemblies were evaluated using metrics generated by the Transrate program (v1.0.3) [53]. Trimmed reads were used to calculate a Transrate score for each assembly, which represents the geometric mean of all contig scores multiplied by the proportion of input reads providing positive support for the assembly [50]. Comparative metrics were calculated using Transrate for each MMETSP sample between DIB and the NCGR assemblies using the Conditional Reciprocal Best BLAST hits (CRBB) algorithm [54]. A forward comparison was made with the NCGR assembly used as the reference and each DIB re-assembly as the query. Reverse comparative metrics were calculated with each DIB re-assembly as the reference and the NCGR assembly as the query. Transrate scores were calculated for each assembly using the Trimmomatic quality-trimmed reads, prior to digital normalization.

Benchmarking Universal Single-Copy Orthologs (BUSCO) software (version 3) was used with a database of 215 orthologous genes specific to protistans and 303 genes specific to eukaryota with open reading frames in the assemblies. BUSCO scores are frequently used as one measure of assembly completeness [55]

To assess the occurrences of fixed-length words in the assemblies, unique 25-mers were measured in each assembly using the HyperLogLog (HLL) estimator of cardinality built into the khmer software package [56]. We used the HLL function to digest each assembly and count the number of distinct fixed-length substrings of DNA (k-mers).

**Table 1.** Number of assemblies with higher values in NCGR or DIB for each quality metric. The "cds" or "nt" indicate the version of the NCGR assembly compared with. The NCGR "cds" assemblies were filtered for ORF content.

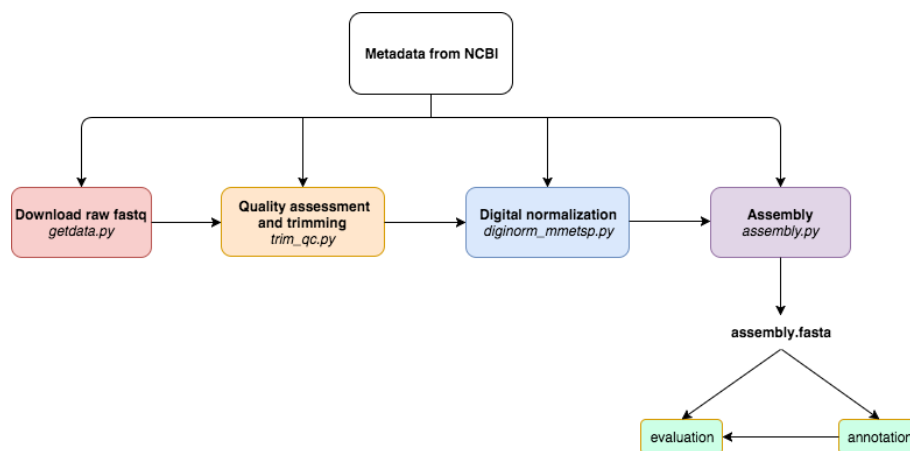
Quality Metric	Higher in NCGR	Higher in DIB
Transrate score, "cds"	44	583
Transrate score, "nt"	495	143
Mean ORF %, "cds"	592	35
Mean ORF %, "nt"	42	596
% References with CRBB, "nt"	100	538
Number of contigs, "nt"	12	626
% Complete BUSCO, Eukaryota, "nt"	381	235

Unique gene names were compared from a random subset of 296 samples using the dammit annotation pipeline [49]. If a gene name was annotated in NCGR but not in DIB, this was considered a gene uniquely annotated in NCGR. Unique gene names were normalized to the total number of annotated genes in each assembly.

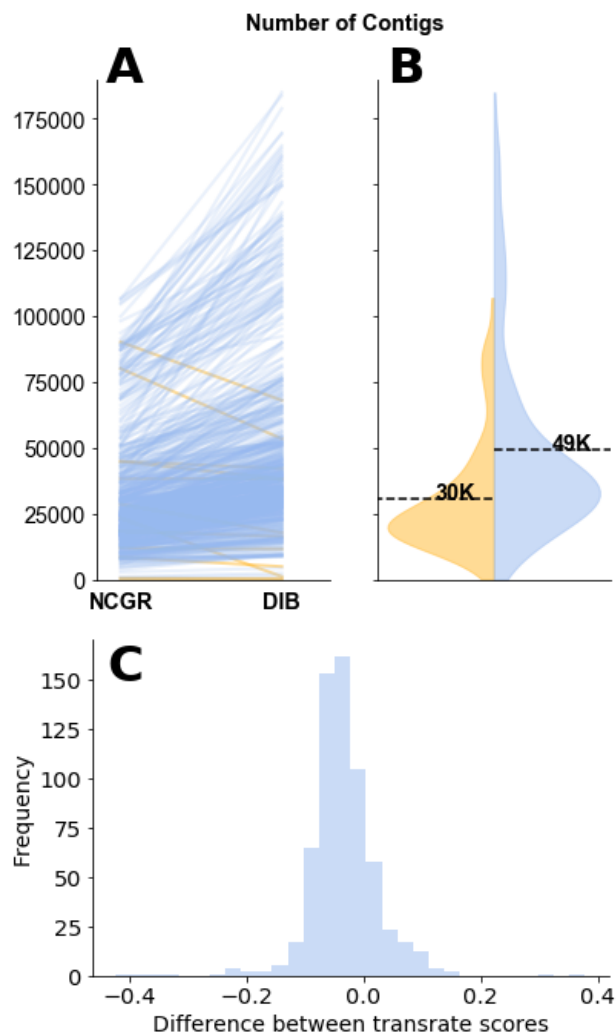
A Tukey's honest significant different (HSD) post-hoc range test of multiple pairwise comparisons was used in conjunction with an ANOVA to measure differences between distributions of data from the top eight most-represented phyla ("Bacillariophyta", "Dinophyta", "Ochrophyta", "Haptophyta", "Ciliophora", "Chlorophyta", "Cryptophyta", "Others") using the 'agricolae' package version 1.2-8 in R version 3.4.2 (2017-09-28). Margins sharing a letter in the group label are not significantly different at the 5% level (8). Averages are reported  $\pm$  standard deviation.

## Results

After assemblies and annotations were completed, files were uploaded to Figshare and Zenodo and are available for download [57, 58, 59]. Due to obstacles encountered uploading and maintaining 678 assemblies on Figshare, Zenodo will be the long-term archive for these re-assemblies <https://doi.org/10.5281/zenodo.740440>. Assembly quality metrics were summarized and are available (Supplemental Tables 1 and 2 [43]).



**Figure 1.** A programmatically automated *de novo* transcriptome assembly pipeline was developed for this study. Metadata in the SraRunInfo.csv file downloaded from NCBI was used as input for each step of the pipeline to indicate which samples were processed. The steps of the pipeline are as follows: download raw fastq data with the fastq-dump script in the SRA Toolkit, perform quality assessment with FastQC and trim residual Illumina adapters and low quality bases (Q<2) with Trimmomatic, do digital normalization with khmer version 2.0, and perform *de novo* transcriptome assembly with Trinity. If a process was terminated, the automated nature of this pipeline allowed for the last process to be run again without starting the pipeline over. In the future, if a new sample is added, the pipeline can be run from beginning to end with just new samples, without having to repeat the processing of all samples in the dataset as one batch. If a new tool becomes available, for example a new assembler, it can be substituted in lieu of the original tool used by this pipeline.



**Figure 2.** The number of contigs and Transrate quality score for each data set varied between DIB and NCGR assemblies. (A) Slopegraphs show shifts in the number of contigs for each individual sample between the DIB and the NCGR assembly pipelines. Negative slope (yellow) lines represent values where NCGR was higher than DIB and positive slope (blue) lines represent values where DIB was higher than NCGR. (B) Split violin plots show the distribution of the number of contigs in each assembly with the original assemblies from NCGR in yellow (left) and the DIB re-assemblies and in blue (right side of B). (C) The difference in Transrate score between the DIB and NCGR assemblies is shown as a histogram. Negative values on the x-axis indicate that the NCGR assembly had a higher Transrate score and positive values indicate that the DIB assembly had a higher Transrate score.

#### Differences in available evaluation metrics between NCGR and DIB were variable.

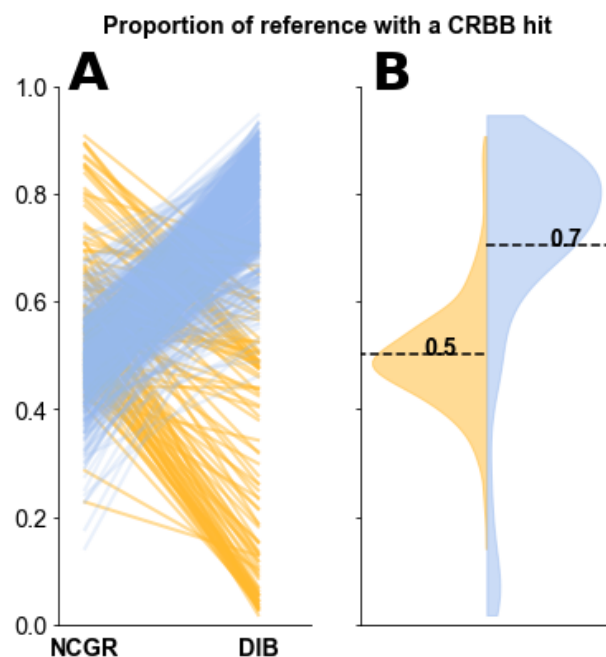
The majority of transcriptome evaluation metrics collected for each sample were higher in Trinity-based DIB re-assemblies than for the Trans-ABYSS-based NCGR assemblies, 'cds' versions (1). The Transrate score from the "nt" version of the assemblies were higher in NCGR vs. DIB, whereas compared to the 'cds' version, the DIB re-assemblies were higher (Supplemental Figure 1 [43]). Since the NCGR 'cds' assemblies were filtered for open reading frame (ORF) content, and the DIB re-assemblies were not filtered, the unfiltered NCGR 'nt' assemblies are more comparable to the DIB re-assemblies.

The DIB re-assemblies had more contigs than the NCGR assemblies in 83.5% of the samples (1). The mean number of contigs in the DIB re-assemblies was  $48,361 \pm 35,703$  while the mean number of contigs in the NCGR 'nt' assemblies was  $30,532 \pm 21,353$  (2). A two-sample Kolmogorov-Smirnov test

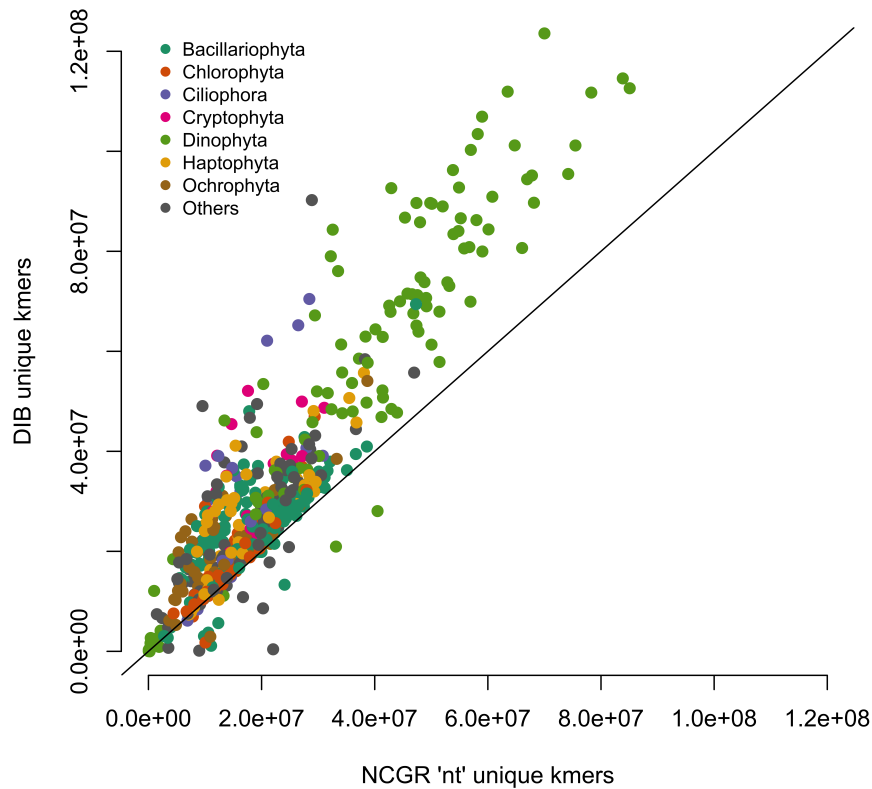
comparing distributions indicated that the number of contigs were significantly different between DIB and NCGR assemblies ( $p < 0.001$ ,  $D = 0.35715$ ). Transrate scores [53], which calculate the overall quality of the assembly based on the original reads, were significantly higher in the DIB re-assemblies ( $0.31 \pm 0.1$ ) compared to the 'cds' versions of the NCGR assemblies ( $0.22 \pm 0.09$ ) ( $p < 0.001$ ,  $D = 0.49899$ ). The Transrate scores in the NCGR 'nt' assemblies ( $0.35 \pm 0.09$ ) were significantly higher than the DIB assemblies ( $0.22 \pm 0.09$ ) ( $p < 0.001$ ,  $D = 0.22475$ ) (Supplemental Figure 1 [43]). The frequency of the differences between Transrate scores in the NCGR 'nt' assemblies and the DIB re-assemblies is centered around zero (Figure 2C). Transrate scores from the DIB assemblies relative to the NCGR 'nt' assemblies did not appear to have taxonomic trends (Supplemental Figure 2 [43]).

#### The DIB re-assemblies contained most of the NCGR contigs as well as new content.

We applied CRBB to evaluate overlap between the assemblies. A positive CRBB result indicates that one assembly contains the same contig information as the other. Thus, the proportion of positive CRBB hits can be used as a scoring metric to compare the relative similarity of content between two assemblies. For example, MMETSP0949 (*Chattonella subsalsa*) had 39,051 contigs and a CRBB score of 0.71 in the DIB re-assembly whereas in the NCGR assembly of the same sample had 18,873 contigs and a CRBB score of 0.34. This indicated that 71% of the reference of DIB was covered by the NCGR assembly, whereas in the reverse alignment, the NCGR reference assembly was only covered by 34% of the DIB re-assembly. The mean CRBB score in DIB when queried against NCGR 'nt' as a reference was  $0.70 \pm 0.22$ , while the mean proportion for NCGR 'nt' assemblies queried against DIB re-assemblies was  $0.49 \pm 0.10$  ( $p < 0.001$ ,  $D = 0.71121$ ) (3). This indicates that more content from the NCGR assemblies was included in the DIB re-assemblies than vice versa and also suggests that the DIB re-assemblies overall have additional content. This finding is reinforced by higher



**Figure 3.** (A) Slopegraphs comparing the proportion of CRBB hits between NCGR 'nt' assemblies and DIB assemblies between the same samples. (B) Violin plots showing the distribution of the proportion of NCGR transcripts with reciprocal BLAST hits to DIB (blue) and the proportion of DIB transcripts with reciprocal BLAST hits to NCGR (yellow).



**Figure 4.** Unique numbers of k-mers ( $k=25$ ) in seven most represented phyla, calculated with the HyperLogLog function in the khmer software package. DIB re-assemblies were compared to the NCGR 'nt' assemblies along a 1:1 line. Samples are colored based on their phylum level affiliation. More than 95% of the DIB re-assemblies had more unique k-mers than to the NCGR assembly of the same sample.

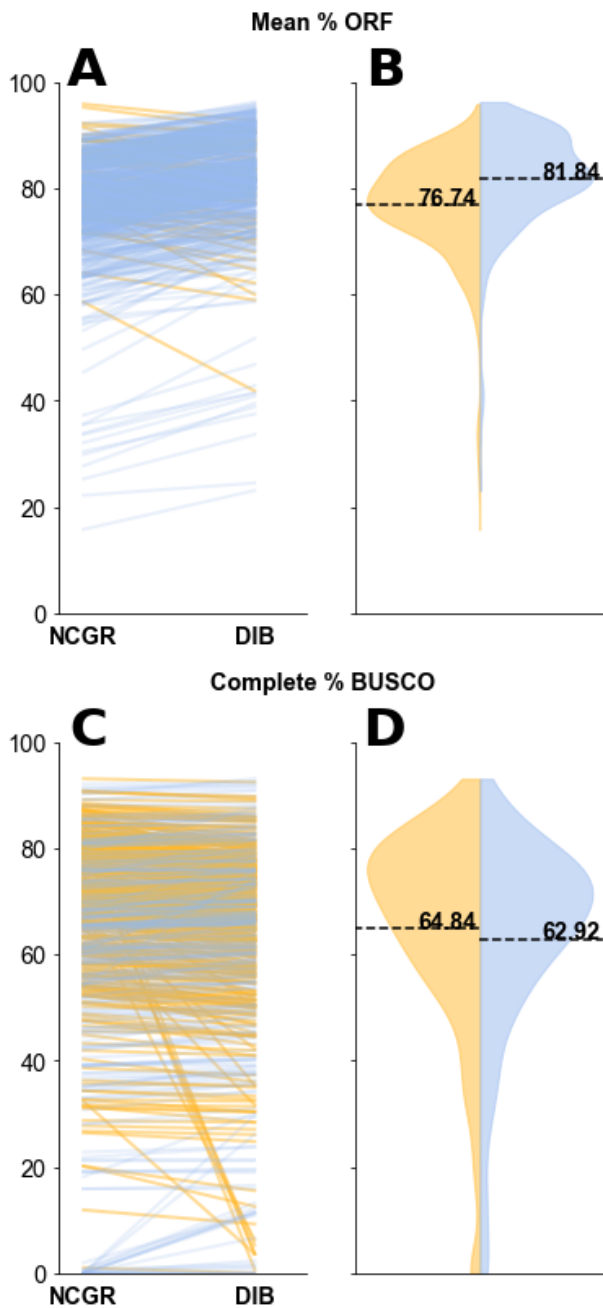
unique k-mer content found in the DIB re-assemblies compared to NCGR, where more than 95% of the samples had more unique k-mers in the DIB re-assemblies compared to NCGR assemblies (4).

To investigate whether the new sequence content was genuine, we examined two different metrics that take into account the biological quality of the assemblies. First, the estimated content of open reading frames (ORFs), or coding regions, across contigs was quantified. Though DIB re-assemblies had more contigs, the ORF content is similar to the original assemblies, with a mean of  $81.8 \pm 9.9\%$  ORF content in DIB re-assemblies and  $76.7 \pm 10.1\%$  ORF content in the NCGR assemblies. Nonetheless, ORF content in DIB re-assemblies was higher than NCGR assemblies for 95% of the samples (5), although DIB re-assemblies had significantly higher ORF content ( $p < 0.001$ ,  $D = 2681$ ). Second, when the assemblies were queried against the eukaryotic BUSCO database [55], the percentages of BUSCO eukaryotic matches in the DIB re-assemblies ( $61.8 \pm 19.9\%$ ) were similar to the original NCGR assemblies ( $63.8 \pm 20.3\%$ ) (5). However, the DIB re-assemblies were significantly different compared to the NCGR assemblies ( $p = 0.002408$ ,  $D = 0.099645$ ). Therefore, although the number of contigs and amount of CRBB content were dramatically increased in the DIB re-assemblies compared to the NCGR assemblies, the differences in ORF content and BUSCO matches compared to eukaryotic (5) and protistan (Supplemental Figure 3 [43]) databases, while they were significantly different, were less dramatic. This suggests that content was not lost by gaining extra contigs. Since the extra content contained roughly similar proportions of ORFs and BUSCO annotations, it is likely

that the re-assemblies contribute more biologically meaningful information.

Looking through the results for missing BUSCOs in the eight samples where NCGR had  $>30\%$  higher complete BUSCO evaluation score (MMETSP0121, MMETSP0932, MMETSP0045, MMETSP0169, MMETSP0232, MMETSP0439, MMETSP0329, MMETSP0717), in some cases a particular orthogroup in the BUSCO database does not produce output for reasons that we don't understand. For example, the Trinity-based pipeline only produced 342 contigs for sample MMETSP0232 while the NCGR 'nt' assembly had 4234 and the 'cds' version had 2736. BUSCO did not recognize any of the DIB contigs but it did recognize the NCGR contigs. For other samples, MMETSP0169 (*Corethron pennatum*, Phylum: Bacillariophyta), the BUSCO software recognized several DIB contigs but the BUSCO group was still considered "missing", even though there were lengths of the contig identified in the output as being similar. For example, the BUSCO orthogroup "EOG0937060I" is a "DNARNA helicase, ATP-dependent, DEAH-box type, conserved site". The BUSCO output indicates the DIB contig, "TRINITY\_DN13758\_c3\_g2\_i1" with a length 974 bases is related to this orthogroup. When we look for this gene in the gff annotation file for MMETSP0169, there are no annotation results for this contig. Another DIB contig, "TRINITY\_DN3716\_co\_g1\_i1" (length 154) is also identified as similar to this same orthogroup. This contig does have annotation results, but it matches with a BUSCO orthogroup, "EOG090C08EI", which is a different gene, Abl-interactor, homeo-domain homologous domain (ABI family, member 3a). The top results comparing the contig





**Figure 5.** The percentage of contigs with a predicted open reading frame (ORF) (A, B) and the percentage of complete protistan universal single-copy orthologs (BUSCO) recovered in each assembly (C, D). In the blue (right side B, D) are the “DIB” re-assemblies and in yellow (left side B, D) are the original ‘nt’ assemblies from NCGR. Slopegraphs (A,C) compare values between the DIB and the NCGR ‘nt’ assemblies. Yellow lines represent negative slope values where NCGR was higher than DIB and blue lines represent positive slope values where DIB was higher than NCGR.

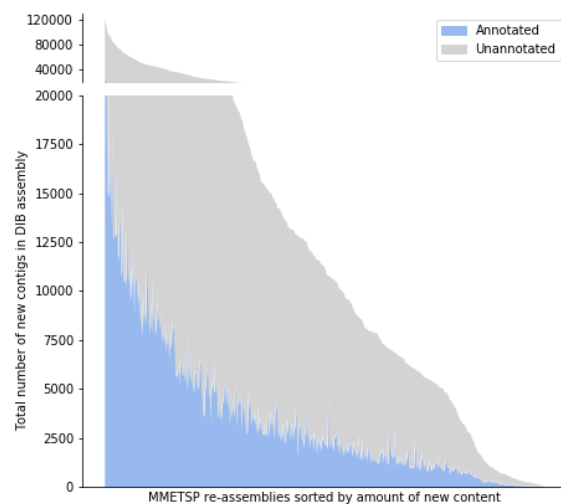
sequence, "TRINITY\_DN13758\_c3\_g2\_i1" against the NCBI blastn database matches with small, several hundred bp regions of the EOG0937060I gene sequence (XM\_021257656.1, XM\_004843976.2, XM\_010604294.2, XM\_010604293.2, XM\_010604291.2, XR\_776390.2). Even though this contig was assembled, it did not successfully annotate. We do not know whether there are errors associated with this assembled contig, or if the contig sequence is unique to this MMETSP0169 organism. Since the BUSCO database and corresponding orthogroups were constructed from multiple sequence alignments with existing individuals in the

databases, it is possible that the transcriptome from the newly sequenced, MMETSP0169 (*Corethron pennatum*) may naturally fall outside the hmm scoring cutoffs for matching with the BUSCO orthogroups. Since the corresponding NCGR assembly had a “Duplicated” result from this particular BUSCO, it is also possible that there is a particular oddity within this ortholog.

There are many examples that can be picked over in these results, which suggests that there is more to learn about the evaluation tools within the context of the organisms in this data set. For now, we conclude that our assemblies are differently fragmented in some regions relative to the NCGR assemblies. We have assembled additional sequences that were not assembled by NCGR. Some NCGR assemblies had different and more complete content than the DIB assemblies. As far as we can tell, there does not appear to be a pattern in the samples that fared well with this pipeline vs. NCGR. This could be a future avenue to explore.

Following annotation by the dammit pipeline [49],  $91 \pm 1.6\%$  of the contigs in the DIB re-assemblies had positive matches with sequence content in the databases queried (Pfam, Rfam, and OrthoDB), with  $48 \pm 0.9\%$  of those containing unique gene names (the remaining are fragments of the same gene). Of those annotations,  $7.8 \pm 0.2\%$  were identified as novel compared to the NCGR ‘nt’ assemblies, determined by a “false” CRBB result (6). Additionally, the number of unique gene names in DIB re-assemblies were higher in 97% of the samples compared to NCGR assemblies, suggesting an increase in genic content (7).

Novel contigs in the DIB re-assemblies likely represent a combination of unique annotations, allelic variants and alternatively spliced isoforms. For example, "FoXV46\_GROCL", "Helicase\_C", "ODR4-like", "PsaA\_PsaB", and "Metazoa\_SRP" are novel gene names found annotated in the DIB re-assembly of the sample MMETSP1473 (*Stichococcus* sp.) that were absent in the NCGR assembly of this same sample. Other gene names,



**Figure 6.** A histogram across MMETSP samples depicting the number of contigs identified as novel in DIB assemblies. These contigs were absent in the NCGR assemblies, based on negative conditional reciprocal best BLAST (CRBB) results. Samples are sorted from highest to lowest number of ‘new’ contigs. The region in gray indicates the number of unannotated contigs present in the DIB re-assemblies, absent from NCGR ‘nt’ assemblies. Highlighted in blue are contigs that were annotated with dammit [49] to a gene name in the Pfam, Rfam, or OrthoDB databases, representing the number of contigs unique to the DIB re-assemblies with an annotation.

for example "Pkinase\_Tyr", "Bromodomain", and "DnaJ", are found in both the NCGR and DIB assemblies, but are identified as novel contigs based on negative CRBB results in the DIB re-assembly of sample MMETSP1473 compared to the NCGR reference.

#### **Assembly metrics varied by taxonomic group being assembled.**

To examine systematic taxonomic differences in the assemblies, metrics for content and assembly quality were assessed (8). Metrics were grouped by the top eight most represented phyla in the MMETSP data set as follows: Bacillariophyta (N=173), Dinophyta (N=114), Ochrophyta (N=73), Chlorophyta (N=62), Haptophyta (N=61), Ciliophora (N=25), Cryptophyta (N=22) and Others (N=130).

While there were no major differences between the phyla in the number of input reads (Figure 8 A), the Dinoflagellates (Dinophyta) had significantly different (higher) number of contigs ( $p < 0.01$ ), unique k-mers ( $p < 0.001$ ), and % ORF ( $p < 0.001$ ) compared to other groups (8), and assemblies from Ciliates (Ciliophora) had lower % ORF ( $p < 0.001$ ) (8).

## **Discussion**

### ***DIB re-assemblies contained the majority of the previously-assembled contigs.***

We used a different pipeline than the original one used to create the NCGR assemblies, in part because new software was available [18] and in part because of new trimming guidelines [46]. The general genome assembler ABySS [20] was used in conjunction with a *de novo* transcriptome assembly pipeline described by Keeling et al. [31]. We had no a priori expectation for the similarity of the results, yet we found that the majority of new DIB re-assemblies included substantial portions of the previous NCGR assemblies seen in the CRBB results. Given this, it may seem surprising that the Transrate and BUSCO scores are lower in the DIB re-assemblies relative to the NCGR counterparts. However, given that the number of contigs and the k-mer content were both dramatically increased in the DIB re-assemblies, it is interesting that the ORFs and annotations were similar between the two assemblies. If the extra content observed was due to assembly artifact, we would not expect these content-based results to be similar. The two metrics, Transrate and BUSCO, which estimate "completeness" of the transcriptomes, may not be telling the whole story. Our results suggest that both pipelines yielded similarly valid contigs, even though the NCGR assemblies appeared to be less sensitive.

The relative increase in number of unique k-mers from the NCGR assemblies to the DIB re-assemblies could be due to the higher number of contigs generated by Trinity. Within the data, the Trinity assembler found evidence for building alternative isoforms. The ABySS assembler and transcriptome pipeline that NCGR used [31] appears to not have preserved that variation, perhaps in an attempt to narrow down the contigs to a consensus transcript sequence.

### ***Re-assembly with new tools can yield new results***

Evaluation with quality metrics suggested that the DIB re-assemblies were more inclusive than the NCGR assemblies. The Transrate scores in the DIB re-assemblies compared to the NCGR 'nt' assemblies were significantly lower, indicating that the NCGR 'nt' assemblies had better overall read inclusion in the assembled contigs whereas the DIB assemblies had higher Transrate scores than the NCGR 'cds' version. This suggests that the NCGR 'cds' version, which was post-processed to only include coding sequence content, was missing information originally in the quality-trimmed reads. As we also saw with % ORF, when filtration steps select only for ORF content

in the NCGR 'cds' versions, potentially useful content is lost. The Transrate score [53] is one of the few metrics available for evaluating the 'quality' of a *de novo* transcriptome. It is similar to the DETONATE RSEM-EVAL score in that it returns a metric indicating how well the assembly is supported by the read data [13]. It does not directly evaluate the underlying de Bruijn graph data structure used to produce the assembled contigs. In the future, metrics directly evaluating the underlying de Bruijn graph data structure may better evaluate assembly quality. Here, the DIB re-assemblies, which used the Trinity *de novo* assembly software, typically contained more k-mers, more annotated transcripts, and more unique gene names than the NCGR assemblies.

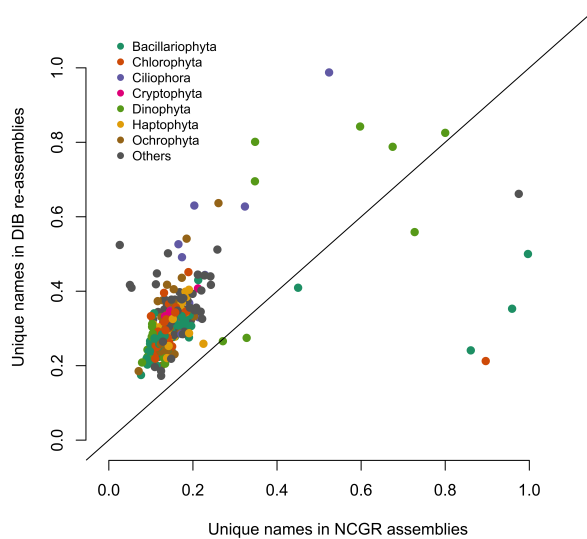
These points all suggest that additional content in these re-assemblies might be biologically relevant and that these re-assemblies provide new content not available in the previous NCGR assemblies. Since contigs are probabilistic predictions of full-length transcripts made by assembly software [18], 'final' reference assemblies are approximations of the full set of transcripts in the transcriptome. Results from this study suggest that the 'ideal' reference transcriptome is a moving target and that these predictions may continue to improve given updated tools in the future.

For some samples, complete BUSCO scores were lower than over half of DIB vs. NCGR. This could be an effect of the BUSCO metric, given that these samples did not perform poorly with other metrics such as % ORF and number of contigs compared to the NCGR. For other samples, MMETSP0252 (*Prorocentrum lima*) in particular, assemblies required several tries and only four contigs were assembled from 30 million reads. The fastq reports were unremarkable, compared to the other samples. In such a large dataset with a diversity of species with no prior sequencing data it is challenging to speculate why each anomaly occurred. However, further investigation into the reasons for failures and peculiarities in the evaluation metrics may lead to interesting discoveries about how we should be effectively assembling and evaluating nucleic acid sequencing data from a diversity of species.

We predict that assembly metrics could have been further improved with longer read lengths of the original data since MMETSP data had only 50 bp read lengths, although this would have presented Keeling et al. [31] with a more expensive data collection endeavor. A study by Chang et al. [25] reported a consistent increase in the percentage of full-length transcript reconstruction and a decrease in the false positive rate moving from 50 to 100 bp read lengths with the Trinity assembler. However, regardless of length, the conclusions we draw would likely remain the same that assembling data with new tools can yield new results.

The DIB re-assemblies, including the additional biologically relevant information, are likely to be meta-transcriptomes. RNA sequences generated from the MMETSP experiments are likely to contain genetic information from more than the target species, as many were not or could not be cultured axenically. Thus, both the NCGR assemblies and the DIB re-assemblies, including the additional biologically relevant information, might be considered meta-transcriptomes. Sequencing data and unique k-mer content likely include bacteria, viruses, or other protists that occurred within the sequenced sample. We did not make an attempt to de-contaminate the assemblies.

The evaluation metrics described here generally serve as a framework for better contextualizing the quality of protistan transcriptomes. For some species and strains in the MMETSP data set, these data represent the first nucleic acid sequence information available [31].



**Figure 7.** Unique gene names found in a subset (296 samples) of either NCGR 'nt' assemblies or DIB re-assemblies but not found in the other assembly, normalized to the number of annotated contigs in each assembly. The line indicates a 1:1 relationship between the unique gene names in DIB and NCGR. More than 97% of the DIB assemblies had more unique gene names than in NCGR assemblies of the same sample.

#### Automated and programmable pipelines can be used to process arbitrarily many RNAseq samples.

The automated and programmable nature of this pipeline was useful for processing large data sets like the MMETSP as it allowed for batch processing of the entire collection, including re-analysis when new tools or new samples become available (see op-ed Alexander et al. 2018). During the course of this project, we ran two re-assemblies of the complete MMETSP data set and one subset as new versions Trinity were released (Supplemental Notebook 2 [43]). Each re-analysis of the complete dataset required only a single command and approximately half a CPU-year of compute. The value of automation is clear when new data from samples become available to expand the data set, tools are updated, or many tools are compared in benchmark studies. Despite this, few assembly efforts completely automate their process, perhaps because the up-front cost of doing so is high compared to the size of the dataset typically being analyzed.

For the purposes of future benchmarking studies, a subset of 12 "High" and 15 "Low" performing samples were identified based on the evaluation metrics: number of contigs, longest contig length, unique k-mers (k=25), and % Complete BUSCO (eukaryota) (Supplemental Figure 4 [43]).

#### Analyzing many samples using a common pipeline identifies taxon-specific trends.

The MMETSP dataset presents an opportunity to examine transcriptome qualities for hundreds of taxonomically diverse species spanning a wide array of protistan lineages. This is among the largest set of diverse RNAseq data to be sequenced. In comparison, the Assemblathon2 project compared genome assembly pipelines using data from three vertebrate species [60]. The BUSCO paper assessed 70 genomes and 96 transcriptomes representing groups of diverse species (vertebrates, arthropods, other metazoans, fungi) [55]. Other benchmarking studies have examined transcriptome qualities for samples representing dozens of species from different taxonomic groups

[61, 62]. A study with a more restricted evolutionary analysis of 15 plant and animals species [62] found no evidence of taxonomic trends in assembly quality but did find evidence of differences between assembly software packages [59].

With the MMETSP data set, we show that comparison of assembly evaluation metrics across this diversity provides not only a baseline for assembly performance, but also highlights particular metrics which are unique within some taxonomic groups. For example, the phyla Ciliophora had a significantly lower percentage of ORFs compared to other phyla. This is supported by recent work which has found that ciliates have an alternative triplet codon dictionary, with codons normally encoding STOP serving a different purpose [37, 38, 39], thus application of typical ORF finding tools fail to identify ORFs accurately in Ciliophora. Additionally, Dinophyta data sets had a significantly higher number of unique k-mers and total contigs in assemblies compared to the assemblies from other data sets, despite having the same number of input reads. Such a finding supports previous evidence from studies showing that large gene families are constitutively expressed in Dinophyta [63].

In future development of *de novo* transcriptome assembly software, the incorporation of phylum-specific information may be useful in improving the overall quality of assemblies for different taxa. Phylogenetic trends are important to consider in the assessment of transcriptome quality, given that the assemblies from Dinophyta and Ciliophora are distinguished from other assemblies by some metrics. Applying domain-specific knowledge, such as specialized transcriptional features in a given phyla, in combination with other evaluation metrics can help to evaluate whether a transcriptome is of good quality or "finished" enough to serve as a high quality reference to answer the biological questions of interest.

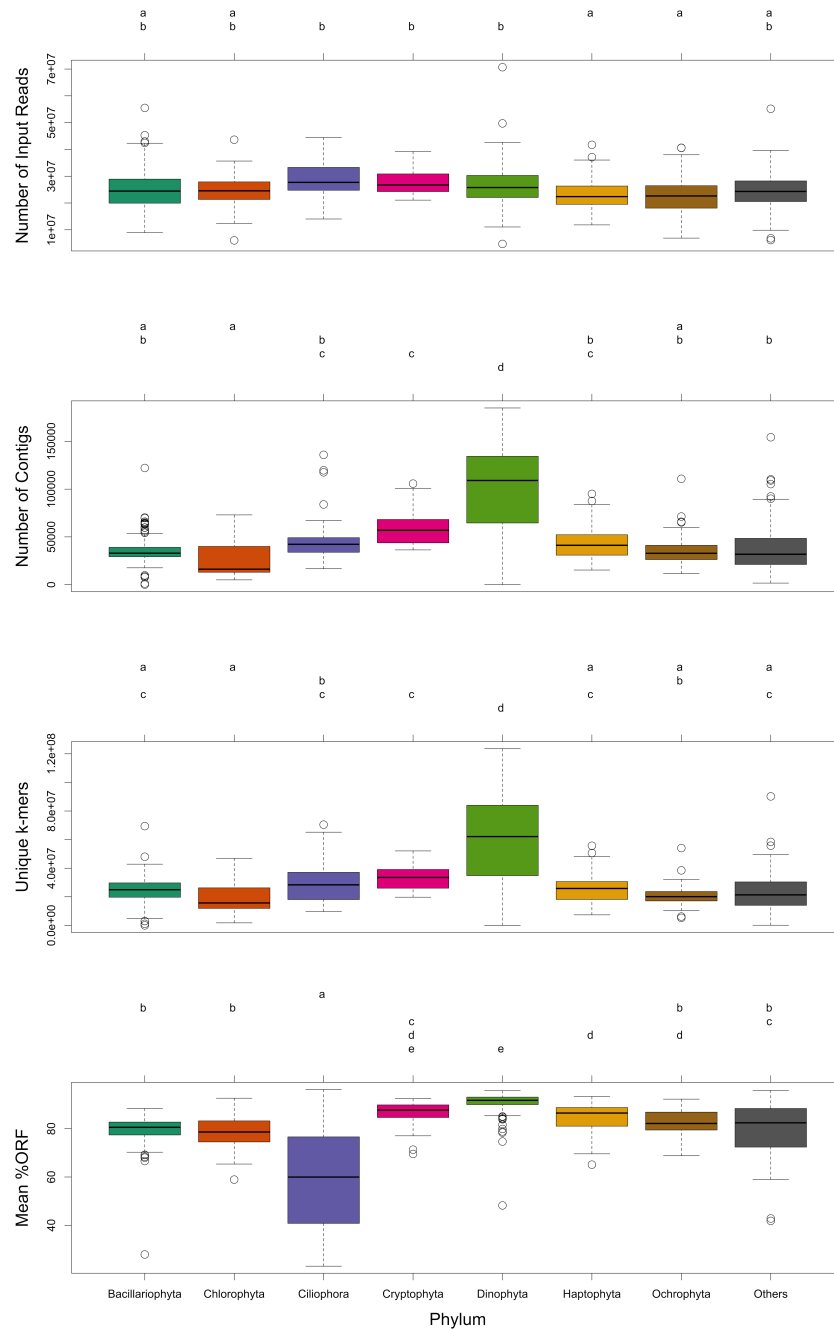
## Conclusion

As the rate of sequencing data generation continues to increase, efforts to programmatically automate the processing and evaluation of sequence data will become increasingly important. Ultimately, the goal in generating *de novo* transcriptomes is to create the best possible reference against which downstream analyses can be accurately based. This study demonstrated that re-analysis of old data with new tools and methods improved the quality of the reference assembly through an expansion of the gene catalog of the dataset. Notably, these improvements arose without further experimentation or sequencing.

With the growing volume of nucleic acid data in centralized and decentralized repositories, streamlining methods into pipelines will not only enhance the reproducibility of future analyses, but will facilitate inter-comparisons amongst from both similar and diverse datasets. Automation tools were key in successfully processing and analyzing this large collection of 678 samples.

## Acknowledgements

The authors would like to gratefully acknowledge the contributions of Camille Scott (0000-0001-8822-8779), Luiz Irber (0000-0003-4371-9659), Taylor Reiter (0000-0002-7388-421X), Daniel Standage (0000-0003-0342-8531) and other members of the Data Intensive Biology lab at UC Davis for providing their time towards helpful assistance with troubleshooting the assembly, annotation and evaluation pipeline and visualizations presented in this paper.



**Figure 8.** Box-and-whisker plots for the seven most common phyla in the MMETSP dataset, (A) number of input reads, (B) number of contigs in the assembly, (C) unique k-mers ( $k = 25$ ) in the assembly, (D) mean percentage open reading frames (ORF). Groups sharing a letter in the top margin were generated from Tukey's HSD post-hoc range test of multiple pairwise comparisons used in conjunction with an ANOVA.

### Availability of supporting data and materials

The data sets supporting the results of this article are available in zenodo and the GigaDB [58, 59, 64]. Pipeline code generated and used for this re-analysis are in zenodo and GigaDB [41, 64]. Original raw RNAseq data generated by Keeling et al. (2014) [31] can be found at [NCBI BioProject PRJNA231566](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA231566).

<https://github.com/dib-lab/dib-MMETSP> Operating system: Linux Programming language: Python Other requirements: fastqc (v0.11.5), trimmomatic (v0.33), khmer (v2.0), trinity, dammit (v1.0.dev0), transrate (v1.0.3), busco (v3) License: 3-clause BSD

### Availability of supporting source code and requirements

Project name: MMETSP: Marine Microbial Eukaryote Transcriptome Sequencing Project code Project home page:

### Declarations

#### List of abbreviations

- BLAST = Basic Local Alignment Search Tool
- BUSCO = Benchmarking Universal Single Copy Orthologs
- CRBB = Conditional Reciprocal Best BLAST

- DIB = Data Intensive Biology Lab at the University of California Davis
- HLL = HyperLogLog
- MMETSP = Marine Microbial Eukaryotic Transcriptome Sequencing Project
- NCGR = National Center for Genome Research
- ORF = Open Reading Frame
- NCBI = National Center for Biotechnology Information
- SRA = Sequence Read Archive

## Ethical Approval

Data were downloaded from public repositories, provided by [31] and NCBI BioProject PRJNA231566 as cited in the text.

## Consent for publication

Not applicable.

## Competing Interests

The authors declare that they have no competing interests.

## Funding

Funding was provided from the Gordon and Betty Moore Foundation under award number GBMF4551 to CTB. Scripts were tested and run on the MSU HPCC and NSF-XSEDE Jetstream cloud platform with allocation TG-BIO160028 [65, 66].

## Author's Contributions

Re-assembly, automated analysis pipeline, and data curation were done by LKJ. Writing, review & editing were done by LKJ, HA, CTB. Visualizations were developed by HA and LKJ.

## Editor's Note

Commentary for this article is available by Alexander et al. [67].

## References

- Geniza M, Jaiswal P. Tools for building de novo transcriptome assembly. *Current Plant Biology* 2017 sep;11-12:41-45. <https://doi.org/10.1016/2Fj.cpb.2017.12.004>.
- Tulin S, Aguiar D, Istrail S, Smith J. A quantitative reference transcriptome for *Nematostella vectensis* early embryonic development: a pipeline for de novo assembly in emerging model systems. *EvoDevo* 2013;4(1):16. <https://doi.org/10.1186/2F2041-9139-4-16>.
- Mittal VK, McDonald JF. De novo assembly and characterization of breast cancer transcriptomes identifies large numbers of novel fusion-gene transcripts of potential functional significance. *BMC Medical Genomics* 2017 aug;10(1). <https://doi.org/10.1186/2F12920-017-0289-7>.
- SONG Y, di LIU H, ZHOU Q, jun ZHANG H, dong ZHANG Z, dong LI Y, et al. High-throughput sequencing of highbush blueberry transcriptome and analysis of basic helix-loop-helix transcription factors. *Journal of Integrative Agriculture* 2017 mar;16(3):591-604. <https://doi.org/10.1016/2Fs2095-3119/2816/2961461-2>.
- Suárez-Vega A, Gutiérrez-Gil B, Klopp C, Tosser-Klopp G, Arranz JJ. Comprehensive RNA-Seq profiling to evaluate lactating sheep mammary gland transcriptome. *Scientific Data* 2016 jul;3:160051. <https://doi.org/10.1038/2Fsdata.2016.51>.
- Carruthers M, Yurchenko AA, Augley JJ, Adams CE, Herzyk P, Elmer KR. De novo transcriptome assembly, annotation and comparison of four ecological and evolutionary model salmonid fish species. *BMC Genomics* 2018 jan;19(1). <https://doi.org/10.1186/2Fs12864-017-4379-x>.
- Mansour TA, Rosenthal JJC, Brown CT, Roberson LM. Transcriptome of the Caribbean stony coral *Porites astreoides* from three developmental stages. *GigaScience* 2016 aug;5(1). <https://doi.org/10.1186/2Fs13742-016-0138-1>.
- Frischkorn KR, Harke MJ, Gobler CJ, Dyhrman ST. De novo assembly of *Aureococcus anophagefferens* transcriptomes reveals diverse responses to the low nutrient and low light conditions present during blooms. *Frontiers in Microbiology* 2014 jul;5. <https://doi.org/10.3389/2Ffmicb.2014.00375>.
- Mansour TA, Scott EY, Finno CJ, Bellone RR, Mienaltowski MJ, Penedo MC, et al. Tissue resolved, gene structure refined equine transcriptome. *BMC Genomics* 2017 jan;18(1). <https://doi.org/10.1186/2Fs12864-016-3451-2>.
- Gonzalez VL, Andrade SCS, Bieler R, Collins TM, Dunn CW, Mikkelsen PM, et al. A phylogenetic backbone for *Bivalvia*: an RNA-seq approach. *Proceedings of the Royal Society B: Biological Sciences* 2015 jan;282(1801):20142332-20142332. <https://doi.org/10.1098/2Frsob.2014.2332>.
- Müller M, Seifert S, Lübke T, Leuschner C, Finkeldey R. De novo transcriptome assembly and analysis of differential gene expression in response to drought in European beech. *PLOS ONE* 2017 sep;12(9):e0184167. <https://doi.org/10.1371/2Fjournal.pone.0184167>.
- Heikkinen LK, Kesäniemi JE, Knott KE. De novo transcriptome assembly and developmental mode specific gene expression of *Pygospio elegans*. *Evolution & Development* 2017 jul;19(4-5):205-217. <https://doi.org/10.1111/2Fede.12230>.
- Li F, Wang L, Lan Q, Yang H, Li Y, Liu X, et al. RNA-Seq Analysis and Gene Discovery of *Andrias davidianus* Using Illumina Short Read Sequencing. *PLOS ONE* 2015 apr;10(4):e0123730. <https://doi.org/10.1371/2Fjournal.pone.0123730>.
- Yu J, Lou Y, Zhao A. Transcriptome analysis of follicles reveals the importance of autophagy and hormones in regulating broodiness of Zhedong white goose. *Scientific Reports* 2016 nov;6(1). <https://doi.org/10.1038/2Fsrep36877>.
- Seo M, Kim K, Yoon J, Jeong JY, Lee HJ, Cho S, et al. RNA-seq analysis for detecting quantitative trait-associated genes. *Scientific Reports* 2016 apr;6(1). <https://doi.org/10.1038/2Fsrep24375>.
- Pedrotty DM, Morley MP, Cappola TP. Transcriptomic Biomarkers of Cardiovascular Disease. *Progress in Cardiovascular Diseases* 2012 jul;55(1):64-69. <https://doi.org/10.1016/2Fj.pcad.2012.06.003>.
- Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, et al. A survey of best practices for RNA-seq data analysis. *Genome Biology* 2016 jan;17(1). <https://doi.org/10.1186/2Fs13059-016-0881-8>.
- Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology* 2011 may;29(7):644-652. <https://doi.org/10.1038/2Fnbt.1883>.
- Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, et al. SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-Seq reads. *Bioinformatics* 2014 feb;30(12):1660-1666. <https://doi.org/10.1093/2Fbioinformatics/2Fbtu077>.

20. Robertson G, Schein J, Chiu R, Corbett R, Field M, Jackman SD, et al. De novo assembly and analysis of RNA-seq data. *Nature Methods* 2010 oct;7(11):909–912. <https://doi.org/10.1038%2Fnmeth.1517>.
21. Schulz MH, Zerbino DR, Vingron M, Birney E. Oases: robust de novo RNA-seq assembly across the dynamic range of expression levels. *Bioinformatics* 2012 feb;28(8):1086–1092. <https://doi.org/10.1093%2Fbioinformatics%2Fbts094>.
22. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *Journal of Computational Biology* 2012 may;19(5):455–477. <https://doi.org/10.1089%2Fcmb.2012.0021>.
23. Peng Y, Leung HCM, Yiu SM, Lv MJ, Zhu XG, Chin FYL. IDBA-tran: a more robust de novo de Bruijn graph assembler for transcriptomes with uneven expression levels. *Bioinformatics* 2013 jun;29(13):i326–i334. <https://doi.org/10.1093%2Fbioinformatics%2Fbtt219>.
24. Kannan S, Hui J, Mazooji K, Pachter L, Tse D. Shannon: An Information-Optimal de Novo RNA-Seq Assembler. *bioRxiv* 2016; <https://www.biorxiv.org/content/early/2016/02/09/039230>.
25. Chang Z, Wang Z, Li G. The Impacts of Read Length and Transcriptome Complexity for De Novo Assembly: A Simulation Study. *PLoS ONE* 2014 apr;9(4):e94825. <https://doi.org/10.1371%2Fjournal.pone.0094825>.
26. Solomon B, Kingsford C. Fast search of thousands of short-read sequencing experiments. *Nature Biotechnology* 2016 feb;34(3):300–302. <https://doi.org/10.1038%2Fbnbt.3442>.
27. Brown CT, Scott C, Crusoe MR, Sheneman L, Rosenthal J, Howe A, khmer-protocols 0.8.4 documentation; 2013. [https://figshare.com/articles/khmer\\_protocols\\_0\\_8\\_3\\_documentation/878460](https://figshare.com/articles/khmer_protocols_0_8_3_documentation/878460).
28. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 2014 apr;30(15):2114–2120. <https://doi.org/10.1093%2Fbioinformatics%2Fbtu170>.
29. Crusoe MR, Alameldin HF, Awad S, Boucher E, Caldwell A, Cartwright R, et al. The khmer software package: enabling efficient nucleotide sequence analysis. *F1000Research* 2015 sep; <https://doi.org/10.12688%2Ff1000research.6924.1>.
30. Zhang Q, Awad S, Brown CT. Crossing the streams: a framework for streaming analysis of short DNA sequencing reads 2015 mar; <https://doi.org/10.7287%2Fpeerj.preprints.890v1>.
31. Keeling PJ, Burki F, Wilcox HM, Allam B, Allen EE, Amaral-Zettler LA, et al. The Marine Microbial Eukaryote Transcriptome Sequencing Project (MMETSP): Illuminating the Functional Diversity of Eukaryotic Life in the Oceans through Transcriptome Sequencing. *PLoS Biology* 2014 jun;12(6):e1001889. <https://doi.org/10.1371%2Fjournal.pbio.1001889>.
32. Caron DA, Alexander H, Allen AE, Archibald JM, Armbrust EV, Bachy C, et al. Probing the evolution, ecology and physiology of marine protists using transcriptomics. *Nature Reviews Microbiology* 2016 nov;15(11):6–20. <https://doi.org/10.1038%2Fnrmicro.2016.160>.
33. Durkin CA, Koester JA, Bender SJ, Armbrust EV. The evolution of silicon transporters in diatoms. *Journal of Phycology* 2016 aug;52(5):716–731. <https://doi.org/10.1111%2Fjpy.12441>.
34. Groussman RD, Parker MS, Armbrust EV. Diversity and Evolutionary History of Iron Metabolism Genes in Diatoms. *PLOS ONE* 2015 jun;10(6):e0129081. <https://doi.org/10.1371%2Fjournal.pone.0129081>.
35. Harke MJ, Juhl AR, Haley ST, Alexander H, Dyhrman ST. Conserved Transcriptional Responses to Nutrient Stress in Bloom-Forming Algae. *Frontiers in Microbiology* 2017 jul;8. <https://doi.org/10.3389%2Ffmicb.2017.01279>.
36. Koid AE, Liu Z, Terrado R, Jones AC, Caron DA, Heidelberg KB. Comparative Transcriptome Analysis of Four Prymnesiophyte Algae. *PLoS ONE* 2014 jun;9(6):e97801. <https://doi.org/10.1371%2Fjournal.pone.0097801>.
37. Alkalaeva E, Mikhailova T. Reassigning stop codons via translation termination: How a few eukaryotes broke the dogma. *BioEssays* 2016 dec;39(3):1600213. <https://doi.org/10.1002%2Fbies.201600213>.
38. Heaphy SM, Mariotti M, Gladyshev VN, Atkins JF, Baranov PV. Novel Ciliate Genetic Code Variants Including the Reassignment of All Three Stop Codons to Sense Codons in *Condylostoma magnum*. *Molecular Biology and Evolution* 2016 aug;33(11):2885–2889. <https://doi.org/10.1093%2Fmolbev%2Fmsw166>.
39. Swart EC, Serra V, Petroni G, Nowacki M. Genetic Codes with No Dedicated Stop Codon: Context-Dependent Translation Termination. *Cell* 2016 jul;166(3):691–702. <https://doi.org/10.1016%2Fj.cell.2016.06.020>.
40. Lowe EK, Swalla BJ, Brown CT. Evaluating a lightweight transcriptome assembly pipeline on two closely related ascidian species 2014 sep; <https://doi.org/10.7287%2Fpeerj.preprints.505v1>.
41. Johnson LK, Alexander H, dib-lab/dib-MMETSP: v2; 2018. <https://doi.org/10.5281/zenodo.594854>.
42. Leinonen R, Sugawara H, and MS. The Sequence Read Archive. *Nucleic Acids Research* 2010 nov;39(Database):D19–D21. <https://doi.org/10.1093%2Fnar%2Fgkq1019>.
43. Johnson LK, Alexander H, Brown CT. Supplemental Information for MMETSP article: 'Re-assembly, quality evaluation, and annotation of 678 microbial eukaryotic reference transcriptomes' 2018; <https://doi.org/10.6084/m9.figshare.7091003>.
44. Andrews S, FastQC: A quality control tool for high throughput sequence data.; 2016. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
45. Ewels P, Magnusson M, Lundin S, Käller M. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* 2016 jun;32(19):3047–3048. <https://doi.org/10.1093%2Fbioinformatics%2Fbtw354>.
46. MacManes MD. On the optimal trimming of high-throughput mRNA sequence data. *Frontiers in Genetics* 2014;5. <https://doi.org/10.3389%2Ffgene.2014.00013>.
47. Brown CT, Howe A, Zhang Q, Pyrkosz AB, Brom TH. A Reference-Free Algorithm for Computational Normalization of Shotgun Sequencing Data 2012 mar; <http://arxiv.org/abs/1203.4802>.
48. Brown CT, What does Trinity's In Silico normalization do?; 2012. <https://doi.org/10.6084/m9.figshare.98198.v1>.
49. Scott C, dammit: an open and accessible de novo transcriptome annotator; 2016. [www.camillescott.org/dammit](http://www.camillescott.org/dammit).
50. Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, et al. The Pfam protein families database: towards a more sustainable future. *Nucleic Acids Research* 2015 dec;44(D1):D279–D285. <https://doi.org/10.1093%2Fnar%2Fgkv1344>.
51. Gardner PP, Daub J, Tate JG, Nawrocki EP, Kolbe DL, Lindgreen S, et al. Rfam: updates to the RNA families database. *Nucleic Acids Research* 2009 jan;37(Database):D136–D140. <https://doi.org/10.1093%2Fnar%2Fgkn766>.
52. Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis P, et al. OrthoDB v9.1: cataloging evolutionary and functional annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Research* 2016 nov;45(D1):D744–D749. <https://doi.org/10.1093%2Fnar%2Fgkz1234>.

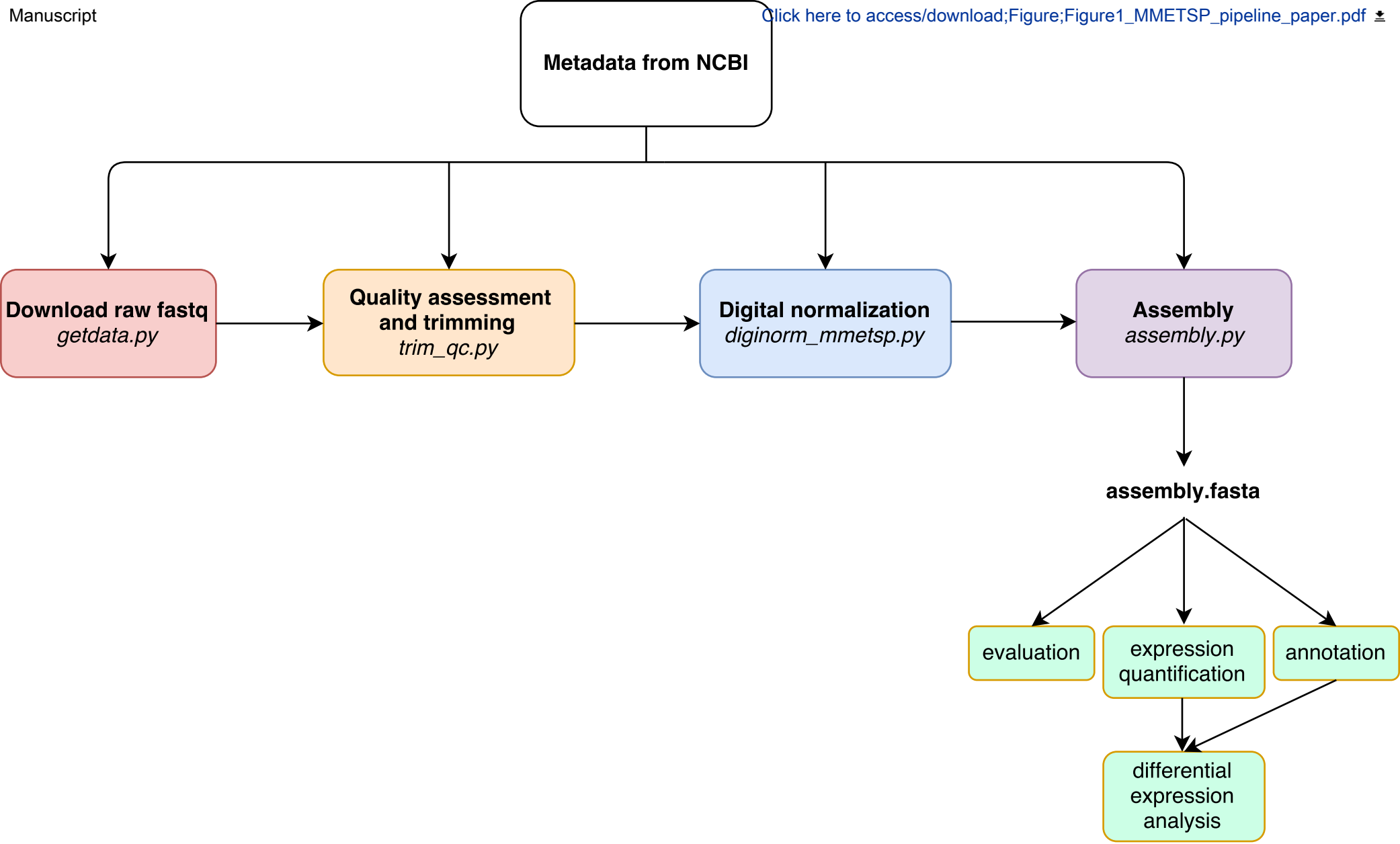
[//doi.org/10.1093%2Fnar%2Fgkw1119](https://doi.org/10.1093%2Fnar%2Fgkw1119).

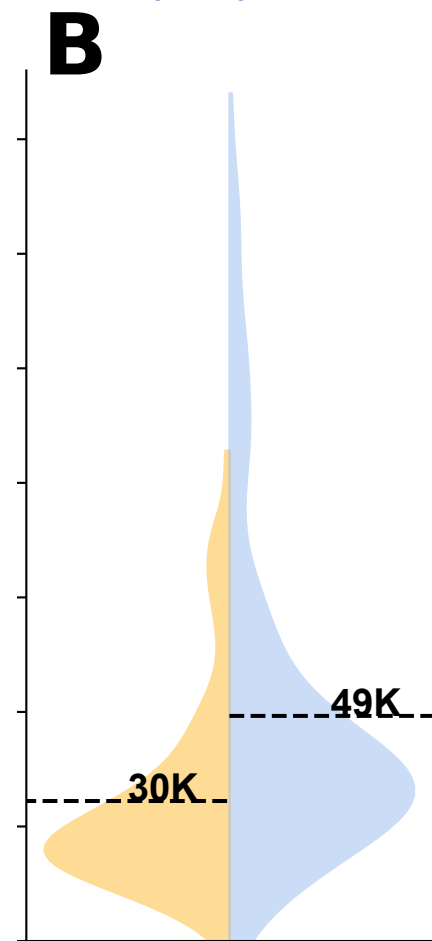
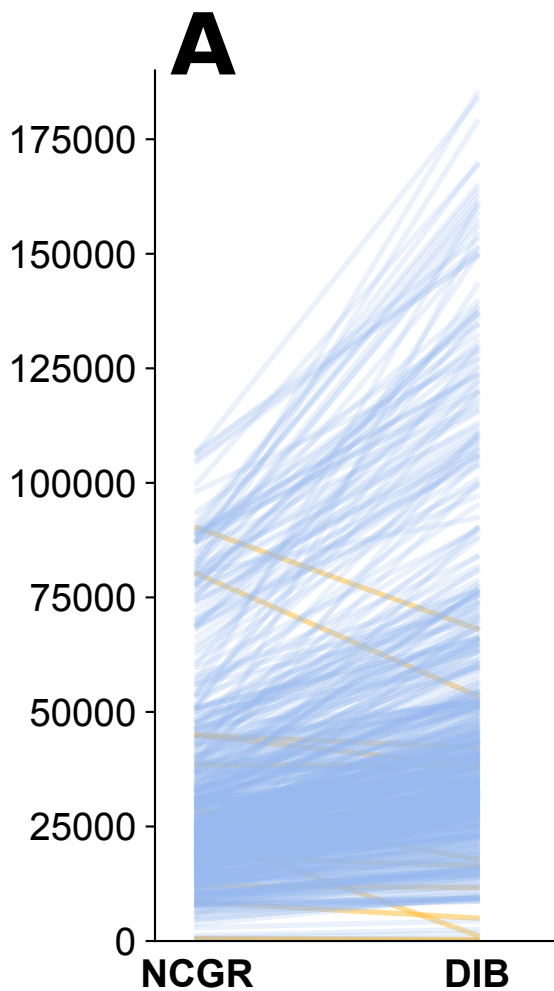
53. Smith-Unna R, Bournnell C, Patro R, Hibberd JM, Kelly S. TransRate: reference-free quality assessment of de novo transcriptome assemblies. *Genome Research* 2016 jun;26(8):1134–1144. <https://doi.org/10.1101%2Fgr.196469.115>.
54. Aubry S, Kelly S, Kümpers BMC, Smith-Unna RD, Hibberd JM. Deep Evolutionary Comparison of Gene Expression Identifies Parallel Recruitment of Trans-Factors in Two Independent Origins of C<sub>4</sub> Photosynthesis. *PLoS Genetics* 2014 jun;10(6):e1004365. <https://doi.org/10.1371%2Fjournal.pgen.1004365>.
55. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 2015 jun;31(19):3210–3212. <https://doi.org/10.1093%2Fbioinformatics%2Fbtv351>.
56. Junior LCI, Brown CT. Efficient cardinality estimation for k-mers in large DNA sequencing data sets 2016 jun;<https://doi.org/10.1101%2F056846>.
57. Johnson L, Alexander H, Brown CT, Marine Microbial Eukaryotic Transcriptome Sequencing Project, re-assemblies; 2018. <https://doi.org/10.6084/m9.figshare.3840153>.
58. Johnson LK, Alexander H, Brown CT. (all datasets) MMETSP re-assemblies 2017 jan;<https://doi.org/10.5281/zenodo.746048>.
59. Johnson LK, Alexander H, Brown CT. MMETSP re-assemblies 2018 apr;<https://doi.org/10.5281/zenodo.740440>.
60. Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, et al. Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. *GigaScience* 2013 jul;2(1). <https://doi.org/10.1186%2F2047-217x-2-10>.
61. Li B, Fillmore N, Bai Y, Collins M, Thomson JA, Stewart R, et al. Evaluation of de novo transcriptome assemblies from RNA-Seq data. *Genome Biology* 2014 dec;15(12). <https://doi.org/10.1186%2Fs13059-014-0553-5>.
62. MacManes MD. The Oyster River Protocol: a multi-assembler and kmer approach for de novo transcriptome assembly. *PeerJ* 2018 aug;6:e5428. <https://doi.org/10.7717%2Fpeerj.5428>.
63. Aranda M, Li Y, Liew YJ, Baumgarten S, Simakov O, Wilson MC, et al. Genomes of coral dinoflagellate symbionts highlight evolutionary adaptations conducive to a symbiotic lifestyle. *Scientific Reports* 2016 dec;6(1). <https://doi.org/10.1038%2Fsrep39734>.
64. Johnson LK, Alexander H, Brown CT, Supporting data for "Re-assembly, quality evaluation, and annotation of 678 microbial eukaryotic reference transcriptomes"; 2018. <http://dx.doi.org/10.5524/100522>.
65. Towns J, Cockerill T, Dahan M, Foster I, Gauthier K, Grimshaw A, et al. XSEDE: Accelerating Scientific Discovery. *Computing in Science & Engineering* 2014 sep;16(5):62–74. <https://doi.org/10.1109%2Fmcse.2014.80>.
66. Stewart CA, Turner G, Vaughn M, Gaffney NI, Cockerill TM, Foster I, et al. Jetstream. In: *Proceedings of the 2015 XSEDE Conference on Scientific Advancements Enabled by Enhanced Cyberinfrastructure - XSEDE '15 ACM Press*; 2015. <https://doi.org/10.1145%2F2792745.2792774>.
67. Alexander H, Johnson LK, Brown CT. Keeping it light: (Re)analyzing community-wide datasets without major infrastructure. *GigaScience* In press;.

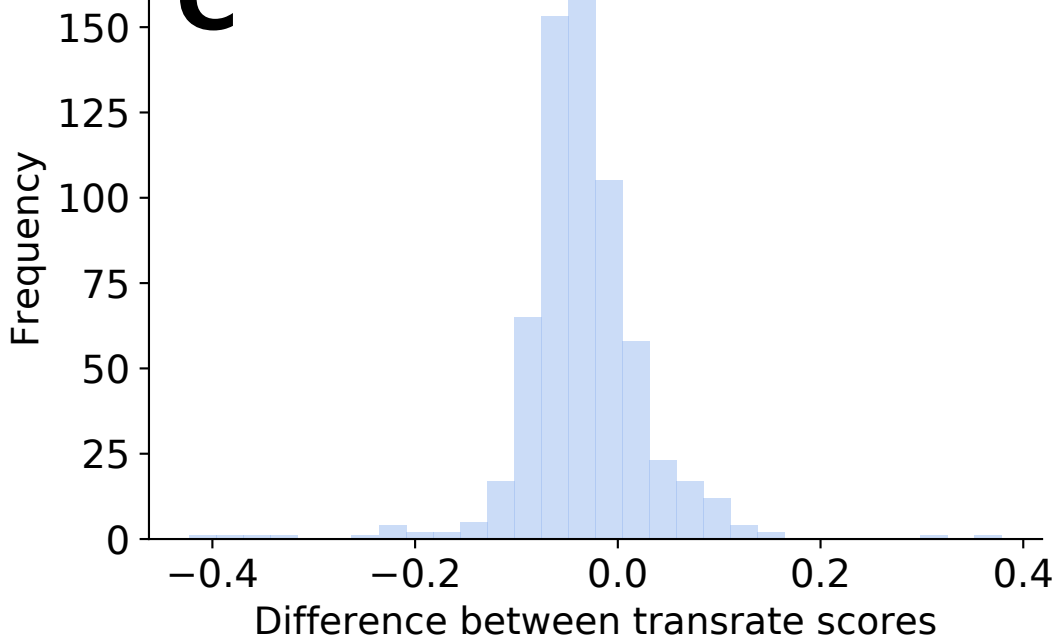
[Click here to view linked References](#)

```
1
2
3
4   This is pdfTeX, Version 3.14159265-2.6-1.40.19 (TeX Live 2018/W32TeX)
5   (preloaded format=pdflatex 2018.7.12)  28 NOV 2018 18:45
6   entering extended mode
7     restricted \writel8 enabled.
8     %&-line parsing enabled.
9   **main.tex
10  (./main.tex
11  LaTeX2e <2018-04-01> patch level 5
12
13
14  ! LaTeX Error: File `oup-contemporary.cls' not found.
15
16  Type X to quit or <RETURN> to proceed,
17  or enter new name. (Default extension: cls)
18
19  Enter file name:
20  ! Emergency stop.
21  <read *>
22
23  l.11 ^^M
24
25  *** (cannot \read from terminal in nonstop modes)
26
27
28  Here is how much of TeX's memory you used:
29    10 strings out of 492646
30    215 string characters out of 6133325
31    56709 words of memory out of 5000000
32    3994 multiletter control sequences out of 15000+600000
33    3640 words of font info for 14 fonts, out of 8000000 for 9000
34    1141 hyphenation exceptions out of 8191
35    10i,0n,8p,86b,8s stack positions out of 5000i,500n,10000p,200000b,80000s
36  ! ==> Fatal error occurred, no output PDF file produced!
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
```

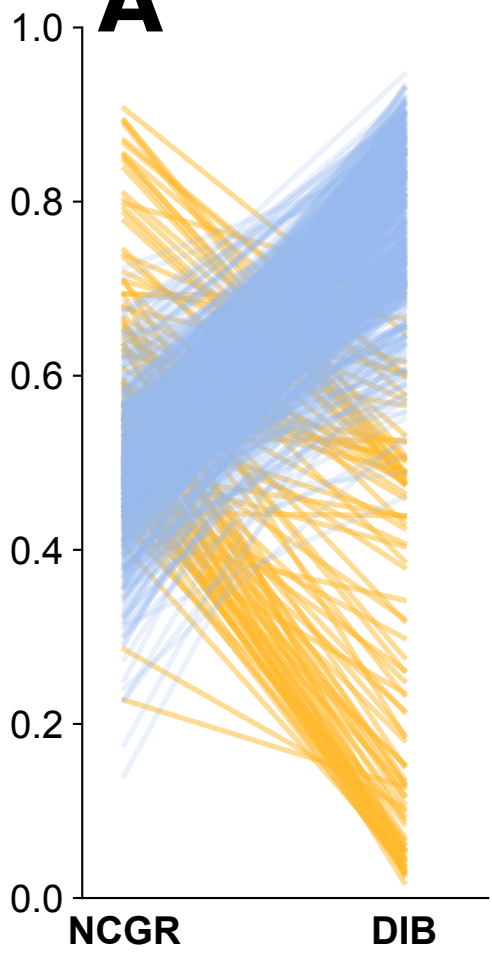




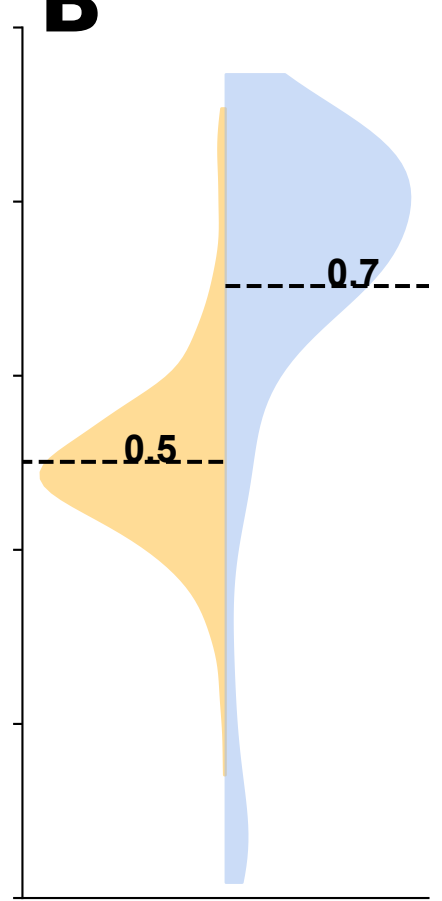


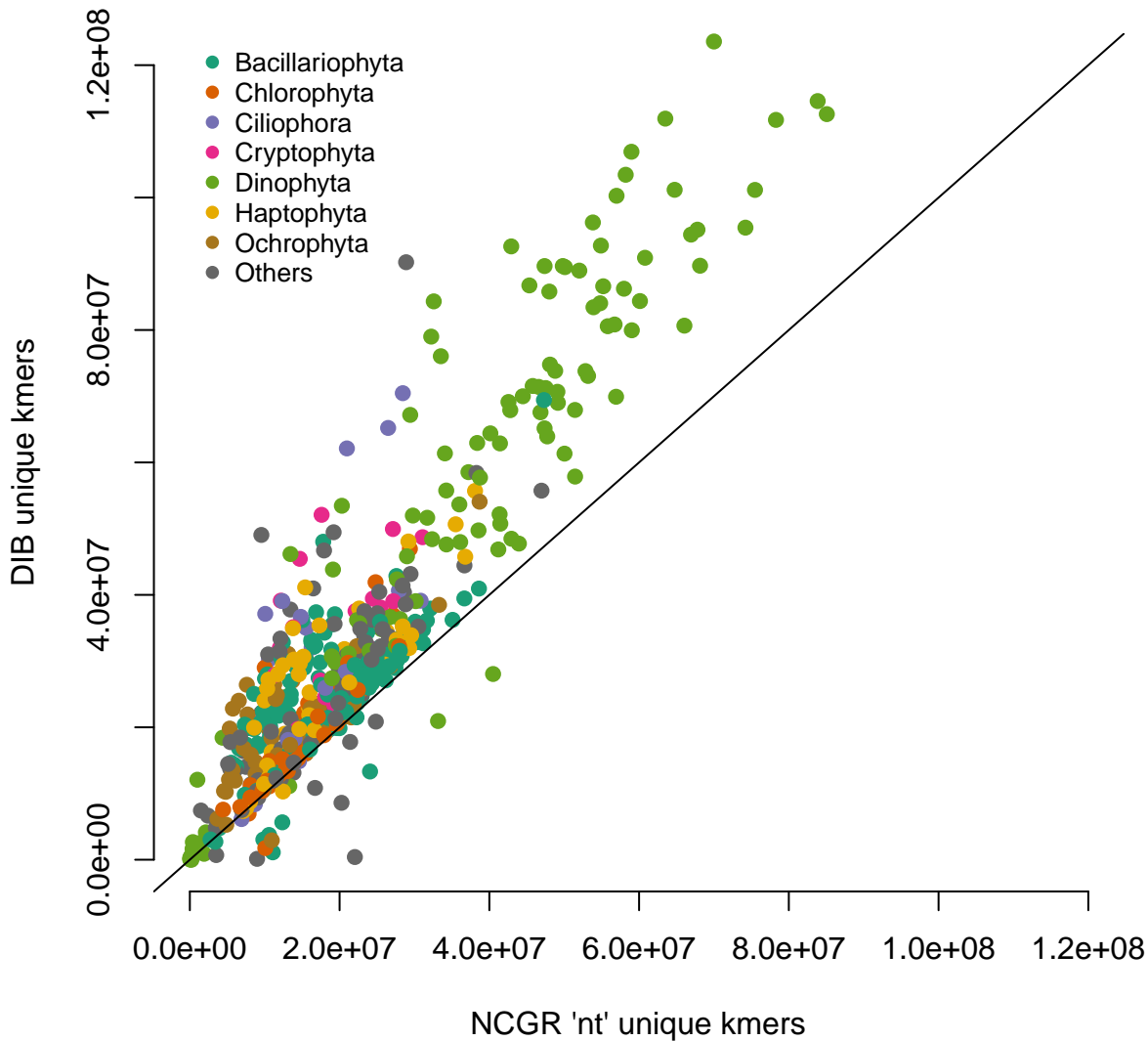
**C**[Click here to access/download;Figure;Figure2c\\_hist\\_tr](#) 

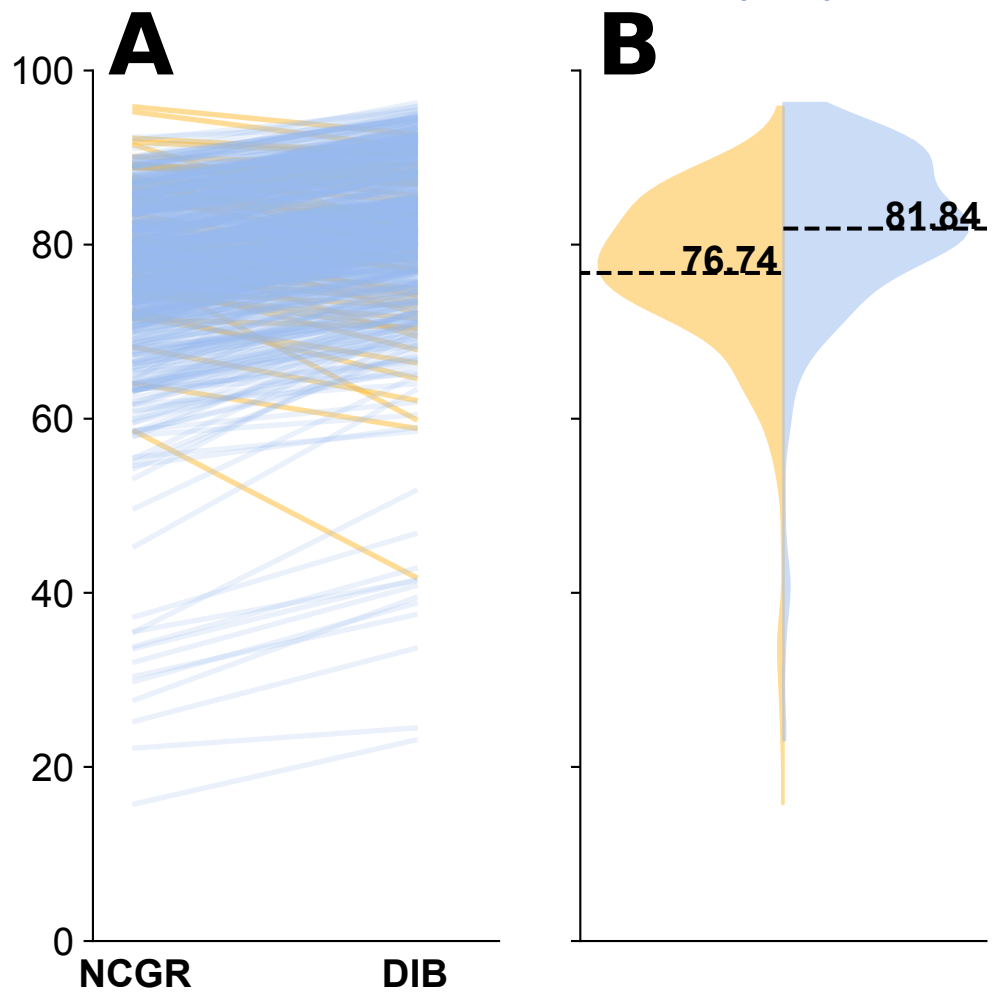
## A

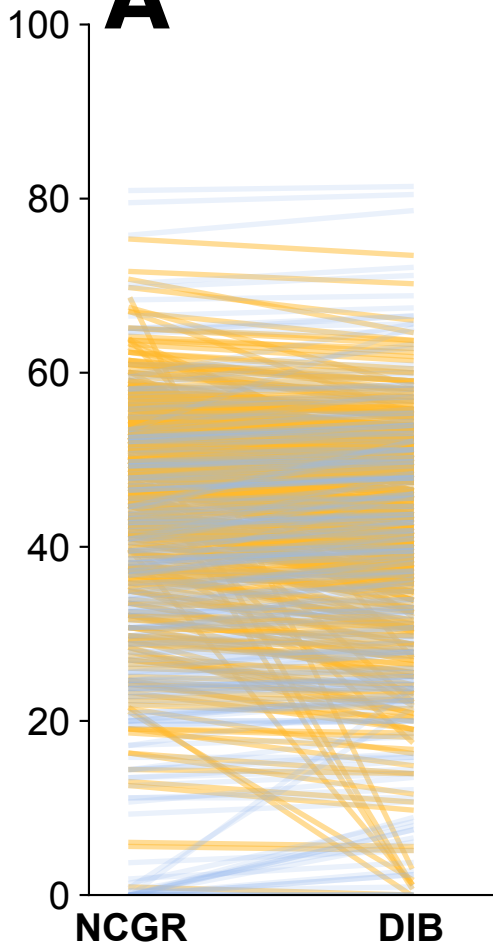


## B







**A****B**