

# 1 To assemble or not to resemble – A validated Comparative Metatranscriptomics Workflow

## 2 (CoMW)

3

## 4 Authors

5 Muhammad Zohaib Anwar<sup>1\*</sup>

6 Anders Lanzen<sup>2,3</sup>

7 Toke Bang-Andreasen<sup>1,4</sup>

8 Carsten Suhr Jacobsen<sup>1\*</sup>

9

## 10 Author Affiliations

11 1 Department of Environmental Science, Aarhus University RISØ Campus, Frederiksborgvej 399,  
12 4000 Roskilde, Denmark

13 2 AZTI, Herrera Kaia, Pasaia, Spain

14 3 IKERBASQUE, Basque Foundation for Science, 48011 Bilbao, Spain

15 4 Department of Biology, University of Copenhagen, Copenhagen, Denmark

16

## 17 \*Corresponding Authors:

18 Muhammad Zohaib Anwar: [mzanwar@envs.au.dk](mailto:mzanwar@envs.au.dk)

19 Carsten Suhr Jacobsen: [csj@envs.au.dk](mailto:csj@envs.au.dk)

## 20 **Abstract**

### 21 **Background**

22 Metatranscriptomics has been used widely for investigation and quantification of microbial  
 23 communities' activity in response to external stimuli. By assessing the genes expressed,  
 24 metatranscriptomics provide an understanding of the interactions between different major  
 25 functional guilds and the environment. Here, we present *de-novo* assembly-based Comparative  
 26 Metatranscriptomics Workflow (CoMW) implemented in a modular, reproducible structure,  
 27 significantly improving the annotation and quantification of metatranscriptomes.  
 28 Metatranscriptomics typically utilize short sequence reads, which can either be directly aligned  
 29 to external reference databases ("assembly-free approach") or first assembled into contigs  
 30 before alignment ("assembly-based approach"). We also compare CoMW (assembly-based  
 31 implementation) with assembly-free alternative workflow, using simulated and real-world  
 32 metatranscriptomes from Arctic and Temperate terrestrial environments. We evaluate their  
 33 accuracy in precision and recall using generic and specialized hierarchical protein databases.

### 34 **Results**

35 CoMW provided significantly fewer false positives resulting in more precise identification and  
 36 quantification of functional genes in metatranscriptomes. Using the comprehensive database  
 37 M5nr, the assembly-based approach identified genes with only 0.6% false positives at  
 38 thresholds ranging from inclusive to stringent compared to the assembly-free approach yielding  
 39 up to 15% false positives. Using specialized databases (Carbohydrate Active-enzyme and  
 40 Nitrogen Cycle), the assembly-based approach identified and quantified genes with 3-5x less  
 41 false positives. We also evaluated the impact of both approaches on real-world datasets.

## 42 **Conclusions**

43 We present an open source *de-novo* assembly-based Comparative Metatranscriptomics  
44 Workflow (CoMW). Our benchmarking findings support the argument of assembling short reads  
45 into contigs before alignment to a reference database, since this provides higher precision and  
46 minimizes false positives.

## 47 Key Words

48 Metatranscriptomics, Benchmarking, Assembly, Alignment, Precision, Recall, False positives

## 49 1 Introduction

50 Metatranscriptomics provides an unprecedented insight to complex functional dynamics of  
 51 microbial communities in various environments. The method has been applied to study the  
 52 microbial activity in thawing permafrost and the related biogeochemical mechanisms  
 53 contributing to greenhouse gas emissions [1], and Gonzalez *et al.* [2] applied  
 54 metatranscriptomics to evaluate root microbiome response to soil contamination.  
 55 Metatranscriptomics has also been used to study the functional human gut microbiota [3,4].  
 56 The method is typically used to identify, quantify and compare the functional response of  
 57 microbial communities in natural habitats or in relation to environmental or physio-chemical  
 58 impacts.

59 Using high-throughput sequencing techniques such as Illumina, metatranscriptomics offers a  
 60 non PCR biased method for looking at transcriptional activity occurring within a complex and  
 61 diverse microbial population at a specific point in time [5]. However, curation and annotation of  
 62 this complex data has emerged as a major challenge. To date, several studies have used various  
 63 analytic workflows. Typically, short sequence reads are utilized, which can either be individually  
 64 aligned directly to external reference databases (hereafter “assembly-free”) or assembled into  
 65 longer contiguous fragments (contigs) for alignment (hereafter “assembly-based”). Various  
 66 studies have used either of these two general approaches. For example, Poulsen *et al.* [6] used  
 67 an assembly-based approach. An open-source pipeline, IMP [7] also uses this approach in

integrated metagenomic and metatranscriptomic analyses. The assembly-free Approach has instead been used by e.g. Jung *et al.* [8], aligning short reads to reference genomes of lactic acid bacterial strains associated with the kimchi microbial community. Similarly, an open source pipeline developed by Martinez *et al.* [9] to analyse metatranscriptomics data-sets also aligns short reads directly to a protein database before annotation. The choice of either of these two alternatives for metatranscriptomics analyses may depend on lack of thorough comparisons. Since no independent and direct comparison between them has been performed presently, various metatranscriptomics analysis approaches may at times produce inconsistent observations, even if identical databases are used in the analysis. Thus, standardization of computational analysis is necessary to enable further propagation of metatranscriptomics approaches and their integration into microbial ecology research. Benchmarking provides a critical view of the efficiency and precision of different workflows and use of simulated communities for benchmarking enables the analysis to be independent of experimental variation and biases [10].

Here, we present Comparative Metatranscriptomic Workflow (CoMW) implemented using the *de-novo* assembly-based approach, standardized and validated for functional annotation and quantitative expression analysis. We validated the suitability of CoMW for functional analysis by comparing it to a typical assembly-free approach using simulated datasets and evaluated the accuracy of both approaches using precision, recall and False Discovery Rates (FDR). Three different protein databases were selected for this benchmarking in order to include a representative selection of three different degrees of specialization, on a range from a more inclusive database with wide coverage (universality) and low degree of expert curation, to a

90 smaller, highly curated database, with more narrow coverage: 1) M5nr [11] :-- an inclusive and  
 91 comprehensive non-redundant protein database in combination with eggNOG hierarchical  
 92 annotation 2) Carbohydrate-Active Enzymes (CAZymes) [12] :-- a database dedicated to  
 93 describing the families of structurally-related catalytic and carbohydrate-binding modules of  
 94 enzymes and 3) Nitrogen Cycling Database (NCycDB) [13] :-- a specialized and manually curated  
 95 database covering only N cycle genes. Finally, in order to estimate the consistency and variance  
 96 in the results caused by the choice of approach we then applied them to real world  
 97 metatranscriptomes from microbial communities in 1) active-layer permafrost soil from  
 98 Svalbard [14] and 2) Ash impacted Danish Forest soil [15].

## 99 **2 Findings**

### 100 **2.1 Comparative Metatranscriptomics Workflow (CoMW)**

101 We have standardized, implemented, and validated a metatranscriptomic workflow (CoMW)  
 102 using de-novo assembly-based approach that can assist in analysing large metatranscriptomics  
 103 data. It makes each step of the metatranscriptomic workflow straightforward and help to make  
 104 these complex analyses more reproducible and the components re-useable in different  
 105 contexts. The core processes such as ORF detection and alignment against the functional  
 106 database are vital in any metatranscriptomic analyses and are, therefore, present uniformly in  
 107 all workflows. However, since most of the tools performing these core processes are ever  
 108 improving, the workflow is implemented in modular format in order to have the possibility of  
 109 using alternative tools and databases if preferred or use a newer version of these tools.  
 110 Modularity additionally also provides choice where optional steps can be skipped, changed or

even improved in a structural manner for example the scripts are designed to cater contigs from more than one assembler. In addition to core process CoMW has a couple of optional steps such as abundance based and non-coding RNA filtering which can be different in data sets from a different environment. CoMW is open source workflow written in python available at (<https://github.com/anwarMZ/CoMW>) and published as a computational capsule on codeocean [16]. An Anaconda cloud environment is created with the provided configuration file to install third-party tools and dependencies. Help regarding input, output and parameters is provided with each script and a comprehensive tutorial is presented in the GitHub repository.

## **2.2 Evaluation of CoMW (assembly-based Approach) and comparison to an assembly-free method**

In order to compare the performance of the assembly-based workflow CoMW and assembly-free approaches, we simulated community transcript data using 4943 full length genes provided by Martinez *et al.* [9]. We analysed both approaches separately and compared against direct annotation of full-length genes. The full-length genes were annotated using all three databases (M5nr, CAZy and NCycDB) independently to classify them into functional subsystems and gene families. Figure 1 shows detailed workflow of comparative analysis using both approaches.

*Figure 1: Flowchart illustrating the evaluation and benchmarking scheme used for the comparison of alternative approaches. Red path indicates the full-length genes workflow, Green indicates the steps in the assembly-based workflow CoMW and Blue indicates the steps in the assembly-free approach.*

131

## 132 **2.2.1 Functional assignment**

133 **M5nr Alignment** Full length genes of the simulated community dataset were aligned and  
 134 identified into 671 unique eggNOG orthologs, belonging to 19 distinct functional subsystems  
 135 (level II). At the default confidence threshold (bit score 50), the, assembly-free approach  
 136 produced alignments to 820 orthologs with a precision of 85% (14.9% FPs), whereas CoMW  
 137 identified 665 orthologs with a precision of 99.3% (0.6% FPs) at the default confidence threshold  
 138 of 1E-5. Repeating the alignments using a gradient of 15 varying confidence thresholds for each  
 139 approach (Low -  $T_L$ , Medium -  $T_M$  and High -  $T_H$ ; five thresholds / category) resulted in dissimilar  
 140 performance for both approaches. The precision and recall of CoMW did not decrease below  
 141 99.3% and 98.5% respectively throughout all categories whereas the assembly-free approach  
 142 had a maximum precision of 96.3% at  $T_M$  and decreases to 85% at  $T_L$  and  $T_H$ . CoMW also  
 143 produced fewer (only 0.6%) FPs consistently compared to the assembly-free Approach of FPs  
 144 ranging from 14.9% to minimum 3.6% at highest precision. Based on F-Score the most optimal  
 145 alignment for each approach is given in Table 1, whereas detailed values for precision, recall, F-  
 146 Score and FDR are listed in Supplementary Table S1. We then also evaluated both approaches  
 147 by selectively removing sequences belonging to a certain functional subsystem from the M5nr  
 148 database in a controlled manner (segmented cross validation) in order to replicate real world  
 149 metatranscriptomes where a certain functional subsystem can be completely or partially absent  
 150 from the reference database. We removed four (level II) subsystems (“[D] Cell cycle control, cell  
 151 division, chromosome partitioning”; “[L] Replication, recombination and repair”; “[E] Amino  
 152 acid transport and metabolism” and “[R] General function prediction only” and “[S] Function



unknown"). The level II subsystems were randomly removed (see data availability for the script used for the removal) one at a time realigning full-length genes and simulated reads using both CoMW and assembly-free approaches to the cropped database to compare identification consistency. In each validation round, both precision and recall of CoMW were significantly higher than assembly-free approach. Recalling ability of assembly-free approach dropped significantly in this validation as compared to full database comparison. CoMW also produced less FPs as compared to assembly-free approach. Table 2 provides details for each validation cycle.

**CAZY Alignment** From 2395 full length genes, 500 sequences were aligned to 395 unique functional genes in the CAZY database, which belonged to 130 gene families and were further classified as seven enzyme classes. Using default confidence thresholds (BTS 50, 1E-5), the assembly-free approach identified 765 functional genes belonging to 112 unique families and six enzyme classes with a precision of 28.5% (71.4% FPs). CoMW identified 488 functional genes from CAZY database that were classified into 147 gene families from seven enzyme classes with a precision of 66% (FDR 33.9%) at the default confidence threshold. However, when we repeated the process with 15 various confidence thresholds, precision improved consistently and FPs decreased, whereas for the assembly-free approach, precision dropped significantly with increasing confidence threshold (see Table 1 and Supplementary Table S2).

**NCycDB Alignment** 410 out of 2395 full-length genes were aligned to this database, identified as 29 unique Nitrogen cycle genes and further belonging to 15 functional gene families in five pathways. Using default confidence thresholds, the assembly-free approach identified 1541 functional genes belonging to 25 functional gene families classified into six pathways with a

precision of 0.9% (99% FPs). CoMW identified 42 Nitrogen cycle genes classified into 25 gene families from six pathways with a precision of 59.5% (40.4% FPs) at a default confidence threshold of 1E-5. Like comparisons against M5nr and CAZY we repeated the process with different confidence thresholds for each approach. Precision improved significantly for CoMW at stringent thresholds whereas for the assembly-free approach, the best precision achieved was 5.8%. (Table 1, Supplementary Table S3).

181

182 *Table 1 Comparison of Precision, Recall, F Score and FDR for the assembly-free and the CoMW (assembly-based) approaches*  
 183 *using all three databases based on best F-Score (Full table for both approaches and databases can be seen in Table S1, S2 and*  
 184 *S3). Bold emphasizes better precision, recall, F-Score and FDR in each database between both approaches*

Databases	Approach	Threshold	Threshold Category	Recall	Precision	F-Score	FDR (%)
eggNOG	assembly-free	<i>BTS 120</i>	<i>Strict [TH]</i>	<b>0.9880</b>	0.9540	0.9707	4.5977
	CoMW	<i>1.00E-15</i>	<i>Strict [TH]</i>	0.9851	<b>0.9939</b>	<b>0.9895</b>	<b>0.6006</b>
CAZY	assembly-free	<i>BTS 110</i>	<i>Strict [TH]</i>	0.3510	0.5325	0.4231	46.7433
	CoMW	<i>1.00E-08</i>	<i>Medium [TM]</i>	<b>0.8131</b>	<b>0.7759</b>	<b>0.7940</b>	<b>22.4096</b>
NCycDB	assembly-free	<i>BTS150</i>	<i>Strict [TH]</i>	0.1666	0.0581	0.0862	94.1860
	CoMW	<i>1.00E-14</i>	<i>Strict [TH]</i>	<b>0.6666</b>	<b>0.8333</b>	<b>0.7407</b>	<b>16.6666</b>

185

186 *Table 2 Comparison of Precision, Recall, F Score and FDR for the assembly-free and CoMW (assembly-based) approaches using*  
 187 *the selective removal of functional subsystems from eggNOG database (segmented cross-validation) to evaluate the consistency*  
 188 *of both approaches. Bold emphasizes better consistency compared to Full length genes*

Removed Subsystem	Approach	Recall	Precision	F-Score	FDR (%)
Cell wall/membrane/envelope biogenesis [M]	assembly-free	0.8726	0.9580	0.9133	4.1958
	CoMW	<b>0.9792</b>	<b>0.9855</b>	<b>0.9824</b>	<b>1.4423</b>
Replication, recombination and repair [L]	assembly-free	0.8734	0.9588	0.9141	4.1166
	CoMW	<b>0.9796</b>	<b>0.9858</b>	<b>0.9827</b>	<b>1.415</b>
Amino acid transport and metabolism [E]	assembly-free	0.8750	0.9589	0.9150	4.1095

	CoMW	<b>0.9812</b>	<b>0.9874</b>	<b>0.9843</b>	<b>1.2578</b>
General function prediction only and Function unknown [R], [S]	assembly- free	0.8933	0.9281	0.9104	7.1856
	CoMW	<b>0.9884</b>	<b>0.97443</b>	<b>0.9814</b>	<b>2.5568</b>

189

190 **2.2.2 Expression Quantification**

191 We also compared the ability of both approaches to quantify the expression of identified

192 transcripts by performing differential expression analysis of two groups in simulated

193 communities and compared against the full-length gene expression simulated. We selected

194 three best identification thresholds for both approaches based on highest F-Score and

195 performed differential expression analysis. This analysis for both approaches was carried out

196 against all three databases using the most specific level of hierarchy in the respective databases

197 in order to capture their ability to quantify expression levels of specific genes.

198 According to full-length gene alignments against eggNOG, 123 genes were significantly

199 upregulated and 270 were significantly downregulated. According to the assembly-free

200 Approach (with the best resulting F-Score), 73 genes were up-regulated (precision 94.5%, 5.4%

201 FPs) and 380 (precision 65.7%, 34.2% FPs) were down regulated. whereas using the assembly-

202 based Approach CoMW, 99 genes were identified as up-regulated (precision 94.9%, 5% FPs) and

203 249 down-regulated (precision 97.1%, 2.8% FPs). For the CAZy database full-length genes, 81

204 and 189 genes were identified as significantly up- and down regulated, respectively. Using the

205 assembly-free approach 31 up-regulated (precision 19.3%, 80.6% FPs) and 137 down-regulated

206 genes (precision 52.5%, 47.4% FPs) were identified, whereas the CoMW identified 83

207 (precision 71%, 28.9% FPs) and 191 (precision 73.8%, 26.1% FPs), respectively- In the NCyc

208 database expression analysis, three and 14 genes were seen as significantly up and down-

209 regulated respectively using full-length genes. According to the assembly-free approach, 26  
 210 (precision 0%, 100% FPs) and 107 (precision 4.6%, 95.3% FPs) genes were up and down  
 211 regulated respectively, whereas according to CoMW, three (precision 33.3%, 66.6% FPs) genes  
 212 were up-regulated and 18 (precision 55.5%, 44% FPs) were down-regulated. Precision, Recall  
 213 and FDR for both approaches against all three databases are available in Supplementary Table  
 214 S4. Additionally, we collapsed the functional genes into functional subsystems and gene  
 215 families to remove FPs produced due to identification of homologous proteins or proteins with  
 216 multiple inheritance. Fold change (log2 transformed) was then calculated for each  
 217 subsystem/gene family. (see Figure 2)

218

219 *Figure 2: Differential Expression comparison of the assembly-free and the CoMW assembly-based approaches using*  
 220 *A) M5nr database, B) NCycDB and C) CAZy database.*

221

### 222 **2.2.3 Real-World metatranscriptomes**

223 To evaluate the effect of the two approaches on real world data, two metatranscriptomes from  
 224 microbial communities were studied. In the first study we investigated the transcriptional  
 225 response during warming from -10 °C to 2 °C and subsequent cooling of 2 °C to -10 °C of an  
 226 Arctic tundra active layer soil from Svalbard, Norway . The aim of the study was to understand  
 227 taxonomic and functional shifts in microbial communities caused by climate change in the  
 228 Arctic. A pronounced shift during the incubation period was noticed by Schostag *et al.* [14]  
 229 which was not replicated by the assembly-free approach. However, using CoMW, we identified  
 230 an increase of genes in the subsystem “[P] Inorganic ion transport and metabolism”. During  
 231 cooling, CoMW also captured the upregulation and downregulation of genes related to “[J]

Translation, ribosomal structure and biogenesis” and “[C] Energy production and conversion” respectively (Figure 3) unlike the assembly-free approach. These findings may have implications for our understanding of carbon dioxide emission, Nitrogen cycling and plant nutrient availability in Arctic soils.

*Figure 3: Relative abundance of eggNOG functional subsystems in Arctic permafrost soil identified and quantified using both CoMW and the assembly-free approach compares the differences in observed functional dynamics. Blue dotted line represents trends using CoMW (assembly-based) whereas Red Solid line represents assembly-free approach*

In the second study, we investigated the effects of wood ash amendment on Danish forest soils [15]. Ash was added in three different quantities (0/control, 3, 12 and 90 tonnes ash per hectare (t ha<sup>-1</sup>)) and the effect over time was analysed in soil communities at 0, 3, 30 and 100 days after ash addition. This resulted in strong effects on functional expression as seen in Figure 4. Both approaches once again displayed varying results such as changes in genes related to eggNOG functional subsystem “[W] Extracellular structures”. assembly-free approach also identified 75% of genes as “[S] Function unknown” consistently unlike assembly-based.

*Figure 4: Relative abundance of eggNOG functional subsystems in Ash deposited Danish forest soil with time identified using both the CoMW and an assembly-free approach. Blue dotted line represents trends using CoMW (assembly-based) whereas Red Solid line represents assembly-free approach*

### 3 Discussion

The application of metatranscriptomics is less common than other DNA-based genomics techniques and thus most analysis pipelines are built *ad hoc* [17]. An assembly-free approach is

257 used in a few pipelines/workflows such as COMAN [18], Metatrans [9], and SAMSA2 [19] , while  
 258 an assembly-based approach is used in a few such as IMP [7]. The lack of thorough  
 259 benchmarking studies and standardized workflows in metatranscriptomics has made it a more  
 260 challenging task to analyse the typically big datasets produced. Previous studies e.g. Zhao *et al.*  
 261 & Celaj *et al.* [20,21] have compared *de-novo* sequence assemblers including Trinity  
 262 [22], MetaVelvet [23], Oases [24], AbySS [25] and SOAPden-ovo [26]. Similarly, for assembly-  
 263 free approach direct short read mappers have been compared thoroughly such as DIAMOND  
 264 [27], BLASTX [28] and RAPSearch2 [29] but an independent comparison of the two different  
 265 approaches based on including assembly or directly aligning reads (here “assembly-free”) has  
 266 been lacking. Critical Assessment of Metagenomic Interpreter (CAMI) [30] is so far the most  
 267 comprehensive benchmarking effort, however it lacks any similar metatranscriptomics  
 268 benchmarking. IMP [7] uses an integrated approach of metagenomics and metatranscriptomics  
 269 and has some overlapping areas to CoMW and can be used together due to modular approach  
 270 of CoMW.

271 Using simulated samples comprised of genes collected from abundant genomes provided by  
 272 Martinez *et al.*, we show that both approaches provide similarly high recall rates against the  
 273 general comprehensive database M5nr. However, CoMW provided a significantly better  
 274 precision and a lower false discovery rate for identification and quantification. For relatively  
 275 compact and specialized databases, recall and precision drop for both approaches (especially  
 276 for the most compact database NCyc). Whereas, CoMW still appeared to be more precise,  
 277 meaning that fewer genes were mis-assigned against these database and significantly lower FPs  
 278 were produced.

279 We have attempted to assist this decision-making for processing metatranscriptomic analysis  
 280 by independently assessing the performance of the two most common approaches and provide  
 281 a road map for functional annotation and expression quantification against databases ranging  
 282 from inclusive to specialized. The significantly higher precision in identification and  
 283 quantification for gene families and functional subsystems in simulated samples, against all  
 284 three databases, confirmed that while an assembly step is challenging computationally, it holds  
 285 the potential to reveal information regarding the gene expressions that is not attainable  
 286 without it. Selecting a single best workflow or pipeline for all types of metatranscriptomics  
 287 studies is not a straightforward affair, and we believe that choice of approach changes the  
 288 outcome of study significantly as observed with real-world datasets from active-layer  
 289 permafrost soil from Svalbard and Ash impacted Danish Forest soil. In addition to choosing the  
 290 right workflow, combining that with the appropriate reference database is equally important to  
 291 ensure the best annotation performance. With databases specialized for one or more specific  
 292 environments or functional categories, the assembly-free Approach under-performs due to its  
 293 inability to identify alignments to homologs in the reference database. We also show that the  
 294 assembly-free Approach can increase the FDR in annotation when a database is dominant in  
 295 specific functional subsystem, which can also lead to wrong estimation of fold change in  
 296 expression

297 While taxonomic annotation is beyond the scope of CoMW and thus our benchmarking  
 298 analyses, it is important to consider the limited value of most functional genes for and thus  
 299 functional metatranscriptomics alone for structural profiling of environmental communities,  
 300 due to the high rate of horizontal gene transfer (HGT) [31]. Approaches for this purpose include

the identification of a limited set of “phylogenetic marker genes” (eg.[32]) or “total RNA” metatranscriptomics whereby the rRNA content is retained and utilized for taxonomic analysis [33]. Though not shown here, we expect that the former approach would also benefit in accuracy from assembling mRNA to full length transcripts before classification, based on our results regarding functional diversity. The total RNA approach also benefits from custom rRNA targeted assembly [15], which may be incorporated into CoMW thanks to its modularity.

In summary, we present the assembly-based workflow CoMW and show that this approach results in consistently better accuracy for functional analysis of metatranscriptomics data. Our benchmarking results show that the choice of approach (assembly-free *v* assembly-based) and database significantly affects the quality of the identification, annotation and expression results. Given the impact of each of these variables, it is inevitable that it significantly affects the results of an individual study and comparison of across studies. We believe that the work presented here will both provide a useful tool for and assist the microbial ecology research community to make more informed decisions about the most appropriate methodological approach to analyze large metatranscriptomic datasets with improved precision.

316

## 317 **4 Methods**

### 318 **4.1 CoMW Implementation**

CoMW (assembly-based) is based on four major steps: 1) *De-novo* Assembly and Mapping; 2) Filtering; 3) Gene Prediction and Alignment 4) Annotation.

*De-novo Assembly and Mapping* of short reads back to assembled contigs is done using Trinity [22] and BWA [34] respectively. Various tools have been developed for de-novo



323 metatranscriptome reconstruction that usually rely on graph-theory. Trinity however generates  
324 the most optimal assemblies for coding RNA reads [17,21,35]. Nevertheless, in CoMW, user can  
325 assemble short reads into contigs by any assembler preferred but it can reduce the quality of  
326 the following steps such as alignment of contigs.

327 *Filtering of Contigs* is done to remove variance in sequences/samples. Since CoMW is assembly-  
328 based, after we assemble the reads into longer contigs we also propose a 2-step filtering of the  
329 contigs to remove any chimeric or false contig made as a result of assembly or sequencing error  
330 by removing contigs that have an expression level less than a specific threshold and to remove  
331 any potential non-coding RNA contigs assembled. We can filter contig abundance data by  
332 removing all contigs with relative expression lower than a specific cut-off, e.g. 1% (selected  
333 based on dataset variance) of the number of sequences in the dataset with least number of  
334 sequences. This threshold is also flexible for different datasets and in some cases not required  
335 at all so CoMW allows user to bypass this step or change the threshold up and down based on  
336 data variation. The filtered contigs are subject to potential non-coding RNA filtration by aligning  
337 them against the RFam database [36] using infernal [37] which is a secondary-structure-aware  
338 aligner that predicts the secondary structure of RNA sequences and similarities based on the  
339 consensus structure models. Once again, the ncRNA filtering is an optional step in CoMW,  
340 though highly recommended in order to reduce FPs.

341 *Gene Prediction and Alignment* is done using Transeq from EMBOSS [38] to predict probable  
342 open reading frames (ORFs) of the contigs (customizable, by default six per contig). We used  
343 SWORD [39] as alignment tool against reference databases. SWORD can be used in parallel  
344 based on computational resources available and the aligned results are parsed and cut-off at a

specific confidence threshold of combination of e-value and alignment length (usually 1e-5, can be changed given the assembly distribution in datasets).

*Annotation* of aligned transcripts from the previous step can be done using the databases such as eggNOG which is a hierarchically structured annotation using a graph-based unsupervised clustering available algorithm to produce genome wide orthology inferences. Aligned proteins are then placed into functional subsystems based on their best hits.), CAZy which is a knowledge-based resource specialized in the Glycogenomics, and NCycDB; a Nitrogen cycle database. This results in a count table with a contig and eggNOG ortholog or CAZy gene or NCyc gene having a certain count from each sample depending upon database used. This count table can be then used for differential expression using state-of-the-art expression analysis suit such as DESeq2 [40] or its wrapper SARTools [41]. For evaluation of CoMW we used the template script provided by the SARTools for DeSeq2 analysis where we specified first group of samples as the reference samples and second group as condition with a parametric mean-variance and Benjamini & Hochberg method for P adjustment [42].

#### 4.2 Assembly-free Workflow

For the assembly-free approach we used the Metatrans pipeline [9], which uses FragGeneScan [43] for ORF predictions in short reads, CD-Hit [44] for gene clustering and Diamond [27] for alignment against the M5nr, CAZy and NCyc [11–13] database. We then used the same annotation script which is included in CoMW. For expression analysis gene counts were normalized between samples using the DESeq2 [40] algorithm. Significantly differentially expressed genes were analysed in SARTools [41] using parametric relationship and p-value 0.05 as significance threshold. The Benjamini and Hochberg correction procedure [42] was used to

adjust p-value. For parameters and versions of tools used in Metatrans see supplementary  
GitHub repository in data availability

### 4.3 Composition of Simulated Communities

In this study we utilised a set of simulated communities from Martinez *et al.* [9] where they collected 4943 genes (coding regions) from five abundant microbial genomes: *Bacteroides vulgatus* ATCC 8482, *Ruminococcus torques* L2-14, *Faecalibacterium prausnitzii* SL3/3, *Bacteroides thetaiotaomicron* VPI-5482 and *Parabacteroides distasonis* ATCC 8503. We simulated short reads into 100 samples using Polyester [45] embedded in a script provided by Martinez *et al.* [9] at coverage of 20x which resulted in a count table and short reads with 2395 genes to add the impact of sequencing coverage that the simulator mimics. The process of regulation of abundance was done by first dividing the 100 samples into two groups (“A” and “B”) and then abundance of randomly selected 10% genes was regulated up- and down up to 4-folds, in addition to this we also knocked out (0 abundance) 5% genes completely from both simulated reads and count tables. The process of selection of samples and genes was random but tracked. To include quality and coverage bias, we used the ART simulator [46] that mimics the coverage bias and thus some genes were removed to produce an equal number of reads in FASTQ format to those produced by Polyester. ART was initially trained with Hi-Seq 2500 Illumina quality error model from dataset discussed above to have a consistent error bias. After simulating FASTQ files we then extracted the quality data and bound it to the FASTA files generating new FASTQ files. With the coverage bias and quality training included we had a total of 62,035,912 reads ( $310,179 \pm 3,454$  reads/sample).

#### 388 4.4 Evaluation Measures

389 We used the standard measures of precision (also named positive predictive value, PPV),  
 390 accounting for how many annotations and identifications of significantly differentially  
 391 expressed gene families and subsystems are correct and defined as  $\frac{TP}{TP+FP}$  and recall (also  
 392 named sensitivity or true positive rate, TPR), accounting for how many correct annotations are  
 393 selected, defined as  $\frac{TP}{TP+FN}$  where TP indicates the number of orthologs that have been correctly  
 394 annotated, FN indicates the number of orthologs/genes/functional subsystem which are in the  
 395 simulated communities but were not found by a certain approach and FP indicates the number  
 396 of orthologs/genes/functional subsystem that have been wrongly annotated (because they do  
 397 not appear in the simulated communities). The F-score is the harmonic mean of precision and  
 398 recall, defined as  $\frac{2*Precision*Recall}{Precision+Recall}$ .

### 399 **Availability of source code and requirements**

- 400 • Project name: Comparative Metatranscriptomics Workflow (*CoMW*)
- 401 • Project home page: <https://github.com/anwarMZ/CoMW>
- 402 • Operating system(s): Platform independent
- 403 • Programming language: Python, R, and bash
- 404 • Other requirements: Requirements mentioned in detailed manual at GitHub
- 405 • License: GNU General Public License v3.0

### 406 **Availability of supporting data and materials**

- 407 • Raw sequence data generated using simulation of full-length genes were deposited in
- 408 the NCBI Sequence Read Archive and are accessible through BioProject accession
- 409 number PRJNA509064
- 410 • Project supplementary scripts: [https://github.com/anwarMZ/CoMW\\_supp](https://github.com/anwarMZ/CoMW_supp)
- 411 • Supplementary File 1 – Precision Recall Analysis of both approaches
- 412 • Supplementary File 2 – Differential Expression Analysis of all approaches using eggNOG
- 413 database
- 414 • Supplementary File 3 – Differential Expression Analysis of all approaches using CAZy
- 415 database
- 416 • Supplementary File 4 – Differential Expression Analysis of all approaches using NCyc
- 417 database

### 418 **Tracking and Reproducibility**

- 419 • CoMW is published as computational capsule on codeocean and can be accessed
- 420 through <https://doi.org/10.24433/CO.1793842.v1>

- 421 • CoMW is registered at SciCrunch.org with RRID – SCR\_017109.

## 422 ***List of abbreviations***

423 FDR: False Discovery Rate, FP: False Positives, TP: True Positives, FN: False Negatives, mRNA:  
424 messenger RNA

## 425 ***Ethical Approval***

426 Not applicable

## 427 ***Consent for publication***

428 Not applicable

## 429 ***Competing Interests***

430 The authors declare that they have no competing interests.

## 431 ***Funding***

432 This work was supported by a grant from the European Commission’s Marie Skłodowska Curie  
433 Actions program under project number 675546 (*MicroArctic*).

## 434 ***Author's Contributions***

435 MZA & CSJ conceived and designed the study. MZA, TBA and AL carried out the data  
436 production. MZA and AL carried out analysis. MZA drafted the manuscript and AL, TBA and CSJ  
437 revised and approved the final version.

## 438 ***Acknowledgements***

439 Authors would like to acknowledge European Commission’s MicroArctic project for the funding.  
440 We would also like to thank authors of Metatrans for providing the data used for simulation.  
441 Additionally, we would like to thank Robert Vaser author of Sword to make it available on  
442 anaconda cloud and helping in integration with CoMW.

## 443 **References**

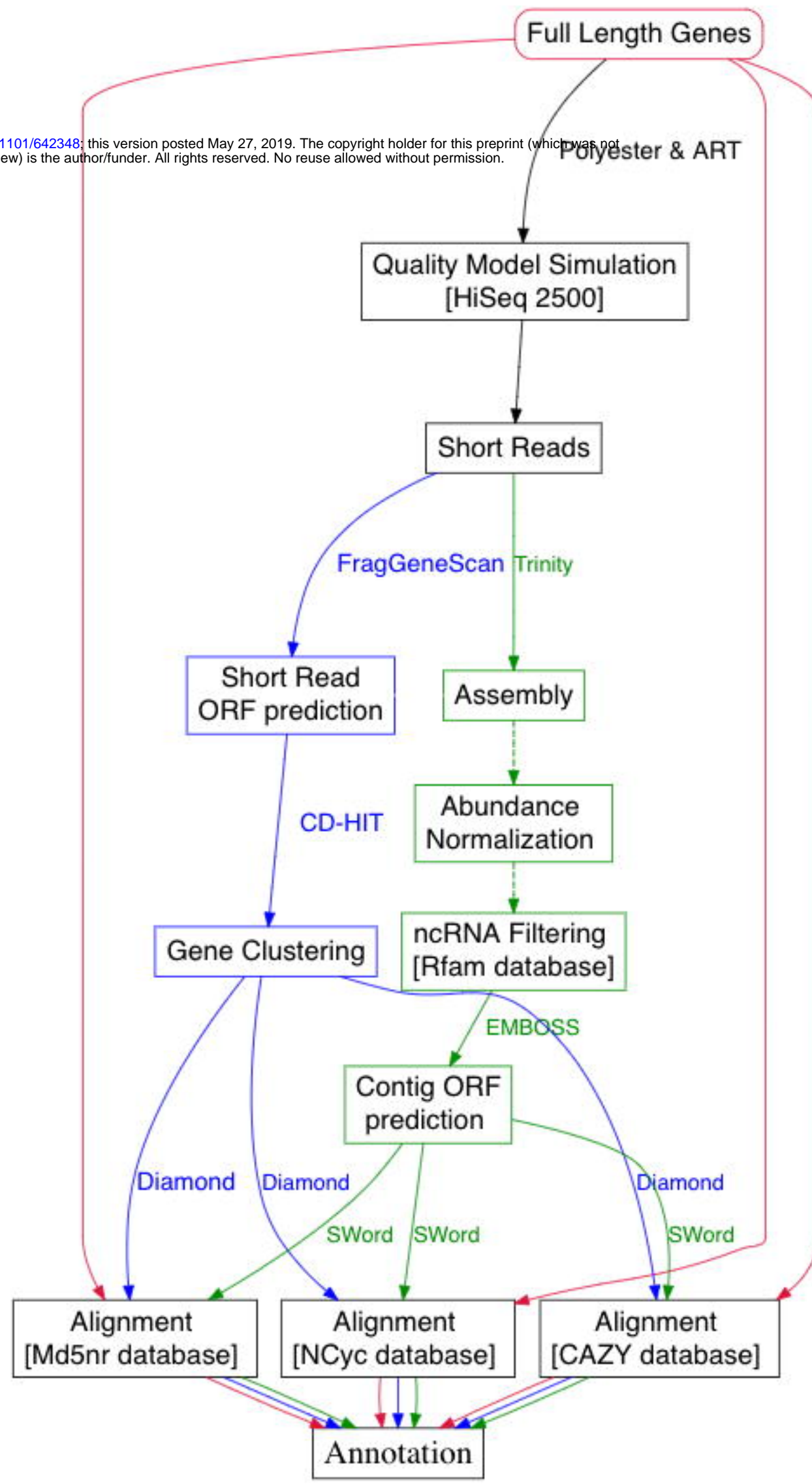
- 444 1. Coolen MJL, Orsi WD. The transcriptional response of microbial communities in thawing Alaskan  
445 permafrost soils. *Front Microbiol.* 2015;6.
- 446 2. Gonzalez E, Pitre FE, Pagé AP, Marleau J, Guidi Nissim W, St-Arnaud M, et al. Trees, fungi and bacteria:  
447 tripartite metatranscriptomics of a root microbiome responding to soil contamination. *Microbiome.*  
448 2018;6:53.
- 449 3. Gosalbes MJ, Durbán A, Pignatelli M, Abellan JJ, Jiménez-Hernández N, Pérez-Cobas AE, et al.  
450 Metatranscriptomic Approach to Analyze the Functional Human Gut Microbiota. *PLOS ONE.*  
451 2011;6:e17447.
- 452 4. Abu-Ali GS, Mehta RS, Lloyd-Price J, Mallick H, Branck T, Ivey KL, et al. Metatranscriptome of human  
453 faecal microbial communities in a cohort of adult men. *Nat Microbiol.* 2018;3:356.
- 454 5. Leimena MM, Ramiro-Garcia J, Davids M, van den Bogert B, Smidt H, Smid EJ, et al. A comprehensive  
455 metatranscriptome analysis pipeline and its validation using human small intestine microbiota datasets.  
456 *BMC Genomics.* 2013;14:530.
- 457 6. Poulsen M, Schwab C, Jensen BB, Engberg RM, Spang A, Canibe N, et al. Methylophilic  
458 methanogenic Thermoplasmata implicated in reduced methane emissions from bovine rumen. *Nat*  
459 *Commun.* 2013;4:1428.
- 460 7. Narayanasamy S, Jarosz Y, Muller EEL, Heintz-Buschart A, Herold M, Kaysen A, et al. IMP: a pipeline  
461 for reproducible reference-independent integrated metagenomic and metatranscriptomic analyses.  
462 *Genome Biol.* 2016;17:260.
- 463 8. Jung JY, Lee SH, Jin HM, Hahn Y, Madsen EL, Jeon CO. Metatranscriptomic analysis of lactic acid  
464 bacterial gene expression during kimchi fermentation. *Int J Food Microbiol.* 2013;163:171–9.
- 465 9. Martinez X, Pozuelo M, Pascal V, Campos D, Gut I, Gut M, et al. MetaTrans: an open-source pipeline  
466 for metatranscriptomics. *Sci Rep.* 2016;6:26447.
- 467 10. Almeida A, Mitchell AL, Tarkowska A, Finn RD. Benchmarking taxonomic assignments based on 16S  
468 rRNA gene profiling of the microbiota from commonly sampled environments. *GigaScience.* 2018;7.
- 469 11. Wilke A, Harrison T, Wilkening J, Field D, Glass EM, Kyrpides N, et al. The M5nr: a novel non-  
470 redundant database containing protein sequences and annotations from multiple sources and  
471 associated tools. *BMC Bioinformatics.* 2012;13:141.
- 472 12. Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. The Carbohydrate-Active  
473 EnZymes database (CAZy): an expert resource for Glycogenomics. *Nucleic Acids Res.* 2009;37:D233-238.
- 474 13. Tu Q, Lin L, Cheng L, Deng Y, He Z. NCycDB: a curated integrative database for fast and accurate  
475 metagenomic profiling of nitrogen cycling genes. *Bioinforma Oxf Engl.* 2018;
- 476 14. Schostag MD, Anwar MZ, Jacobsen CS, Larose C, Vogel TM, Maccario L, et al. Transcriptomic  
477 responses to warming and cooling of an Arctic tundra soil microbiome. *bioRxiv.* 2019;599233.

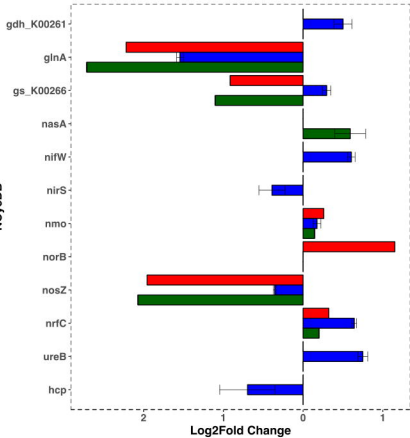
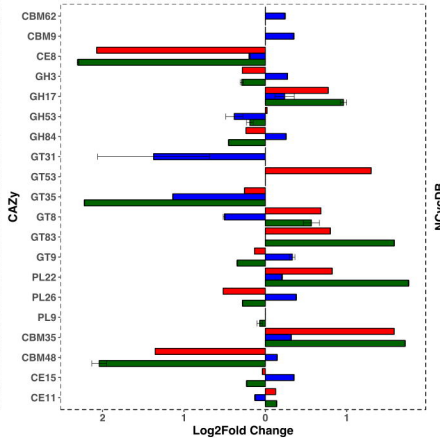
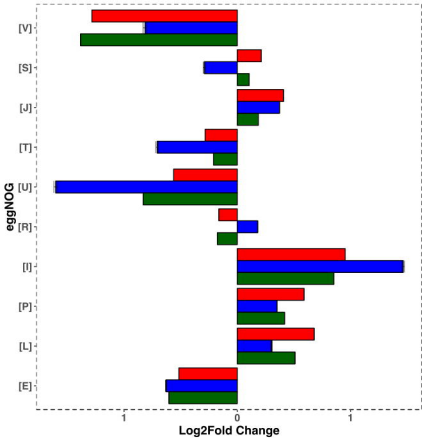
- 478 15. Bang-Andreasen T, Anwar MZ, Lanzen A, Kjølner R, Rønn R, Ekelund F, et al. Total RNA-sequencing  
479 reveals multi-level microbial community changes and functional responses to wood ash application in  
480 agricultural and forest soil. *bioRxiv*. 2019;621557.
- 481 16. Anwar MZ, Lanzen A, Bang-Andreasen T, Jacobsen CS. Comparative Metatranscriptomic Workflow  
482 (CoMW) [source code]. *codeocean*; 2019. Available from: <https://doi.org/10.24433/CO.1793842.v1>
- 483 17. Aguiar-Pulido V, Huang W, Suarez-Ulloa V, Cickovski T, Mathee K, Narasimhan G. Metagenomics,  
484 Metatranscriptomics, and Metabolomics Approaches for Microbiome Analysis: Supplementary Issue:  
485 Bioinformatics Methods and Applications for Big Metagenomics Data. *Evol Bioinforma*.  
486 2016;12s1:EBO.S36436.
- 487 18. Ni Y, Li J, Panagiotou G. COMAN: a web server for comprehensive metatranscriptomics analysis. *BMC*  
488 *Genomics*. 2016;17:622.
- 489 19. Westreich ST, Treiber ML, Mills DA, Korf I, Lemay DG. SAMSA2: a standalone metatranscriptome  
490 analysis pipeline. *BMC Bioinformatics*. 2018;19:175.
- 491 20. Zhao Q-Y, Wang Y, Kong Y-M, Luo D, Li X, Hao P. Optimizing de novo transcriptome assembly from  
492 short-read RNA-Seq data: a comparative study. *BMC Bioinformatics*. 2011;12:S2.
- 493 21. Celaj A, Markle J, Danska J, Parkinson J. Comparison of assembly algorithms for improving rate of  
494 metatranscriptomic functional annotation. *Microbiome*. 2014;2:39.
- 495 22. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Trinity: reconstructing a full-  
496 length transcriptome without a genome from RNA-Seq data. *Nat Biotechnol*. 2011;29:644–52.
- 497 23. Namiki T, Hachiya T, Tanaka H, Sakakibara Y. MetaVelvet: an extension of Velvet assembler to de  
498 novo metagenome assembly from short sequence reads. *Nucleic Acids Res*. 2012;40:e155.
- 499 24. Schulz MH, Zerbino DR, Vingron M, Birney E. Oases: robust de novo RNA-seq assembly across the  
500 dynamic range of expression levels. *Bioinformatics*. 2012;28:1086–92.
- 501 25. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol I. ABySS: A parallel assembler for short  
502 read sequence data. *Genome Res*. 2009;19:1117–23.
- 503 26. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-  
504 efficient short-read de novo assembler. *GigaScience*. 2012;1:18.
- 505 27. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Methods*.  
506 2015;12:59–60.
- 507 28. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*.  
508 1990;215:403–10.
- 509 29. Zhao Y, Tang H, Ye Y. RAPSearch2: a fast and memory-efficient protein similarity search tool for next-  
510 generation sequencing data. *Bioinformatics*. 2012;28:125–6.



- 511 30. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical Assessment of  
512 Metagenome Interpretation – a benchmark of computational metagenomics software. *Nat Methods*.  
513 2017;14:1063–71.
- 514 31. Simonson AB, Servin JA, Skophammer RG, Herbold CW, Rivera MC, Lake JA. Decoding the genomic  
515 tree of life. *Proc Natl Acad Sci U S A*. 2005;102:6608–13.
- 516 32. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, et al. Metagenomic biomarker  
517 discovery and explanation. *Genome Biol*. 2011;12:R60.
- 518 33. Urich T, Lanzén A, Qi J, Huson DH, Schleper C, Schuster SC. Simultaneous Assessment of Soil  
519 Microbial Community Structure and Function through Analysis of the Meta-Transcriptome. *PLoS ONE*.  
520 2008;3.
- 521 34. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.  
522 *Bioinformatics*. 2009;25:1754–60.
- 523 35. Lau MCY, Harris RL, Oh Y, Yi MJ, Behmard A, Onstott TC. Taxonomic and Functional Compositions  
524 Impacted by the Quality of Metatranscriptomic Assemblies. *Front Microbiol*. 2018;9.
- 525 36. Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR. Rfam: an RNA family database. *Nucleic*  
526 *Acids Res*. 2003;31:439–41.
- 527 37. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*.  
528 2013;29:2933–5.
- 529 38. Rice P, Longden I, Bleasby A. EMBOSS: The European Molecular Biology Open Software Suite. *Trends*  
530 *Genet*. 2000;16:276–7.
- 531 39. Vaser R, Pavlović D, Šikić M. SWORD—a highly efficient protein database search. *Bioinformatics*.  
532 2016;32:i680–4.
- 533 40. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data  
534 with DESeq2. *Genome Biol*. 2014;15.
- 535 41. Varet H, Brillet-Guéguen L, Coppée J-Y, Dillies M-A. SARTools: A DESeq2- and EdgeR-Based R Pipeline  
536 for Comprehensive Differential Analysis of RNA-Seq Data. *PLOS ONE*. 2016;11:e0157022.
- 537 42. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to  
538 Multiple Testing. *J R Stat Soc Ser B Methodol*. 1995;57:289–300.
- 539 43. Rho M, Tang H, Ye Y. FragGeneScan: predicting genes in short and error-prone reads. *Nucleic Acids*  
540 *Res*. 2010;38:e191.
- 541 44. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or  
542 nucleotide sequences. *Bioinforma Oxf Engl*. 2006;22:1658–9.
- 543 45. Frazee AC, Jaffe AE, Langmead B, Leek JT. Polyester: simulating RNA-seq datasets with differential  
544 transcript expression. *Bioinformatics*. 2015;31:2778–84.

545 46. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator.  
546 Bioinformatics. 2012;28:593–4.





Assembly-based Assembly-free FullLengthGenes  
(CoMW)

