# GigaScience

## To assemble or not to resemble – benchmarking of metatranscriptomic practices and a validated Comparative Metatranscriptomics Workflow (CoMW)
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-19-00009 |
| Full Title: | To assemble or not to resemble – benchmarking of metatranscriptomic practices and a validated Comparative Metatranscriptomics Workflow (CoMW) |
| Article Type: | Research |
| Funding Information: | H2020 Marie Skłodowska-Curie Actions (675546)     Mr. Muhammad Zohaib Anwar |

| | |
|---|---|
| Abstract: | **Background**<br><br>Metatranscriptomics has been used widely for investigation and quantification of microbial communities' activity in response to external stimuli. By assessing the genes expressed metatranscriptomics provides an understanding of the interactions between different major functional guilds and the environment. Metatranscriptomics typically utilize short sequence reads, which can either be directly aligned to external reference databases ("assembly-free approach") or first assembled into contigs before alignment ("assembly-based approach"). Here we compared workflows representing both alternatives, using simulated and real-world metatranscriptomes from Arctic and Temperate terrestrial environments. We evaluate their accuracy in precision and recall using generic and specialized hierarchical protein databases.<br><br>**Results**<br><br>We show that the assembly-based approach provides significantly fewer false positives resulting in more precise identification and quantification of functional genes in metatranscriptomes. Using the comprehensive database M5nr, the assembly-based approach identifies genes with only 0.6% false positives at thresholds ranging from inclusive to stringent compared to assembly-free approach (3.6 to 15% false positives). Using specialized databases (Carbohydrate Active-enzyme and Nitrogen Cycle) the assembly-based approach identifies and quantifies genes with 3-5x less false positives. We also evaluated the impact of both approaches on real-world datasets. Based on this benchmarking we present a standardized and optimized workflow for identifying functional genes from metatranscriptomes.<br><br>**Conclusions**<br><br>Our findings support the argument of assembling short reads into contigs before alignment to a reference database, since this provides higher precision and minimizes false positives. By virtue of the extensive benchmarking we also present the open source metatranscriptomics analysis workflow Comparative Metatranscriptomics Workflow CoMW. |

| | |
|---|---|
| Corresponding Author: | Muhammad Zohaib Anwar<br>Aarhus University<br>Roskilde, Copenhagen DENMARK |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | Aarhus University |
| Corresponding Author's Secondary Institution: | |
| First Author: | Muhammad Zohaib Anwar |
| First Author Secondary Information: | |
| Order of Authors: | Muhammad Zohaib Anwar |

| | |
|---|---|
| | Anders Lanzen |
| | Toke Bang-Andreasen |
| | Carsten Suhr Jacobsen |
| **Order of Authors Secondary Information:** | |
| **Additional Information:** | |

| Question | Response |
|---|---|
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>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 Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically | Yes |

appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.

Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?

1  **To assemble or not to resemble – benchmarking of metatranscriptomic practices and a**

2  **validated Comparative Metatranscriptomics Workflow (CoMW)**

3

4  **Authors**

5  Muhammad Zohaib Anwar[1]*

6  Anders Lanzen[2,3]

7  Toke Bang-Andreasen[1,4]

8  Carsten Suhr Jacobsen[1]*

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

19  Carsten Suhr Jacobsen: 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 provides an understanding of the interactions between different major

25 functional guilds and the environment. Metatranscriptomics typically utilize short sequence

26 reads, which can either be directly aligned to external reference databases ("assembly-free

27 approach") or first assembled into contigs before alignment ("assembly-based approach"). Here

28 we compared workflows representing both alternatives, using simulated and real-world

29 metatranscriptomes from Arctic and Temperate terrestrial environments. We evaluate their

30 accuracy in precision and recall using generic and specialized hierarchical protein databases.

31 **Results**

32 We show that the assembly-based approach provides significantly fewer false positives resulting

33 in more precise identification and quantification of functional genes in metatranscriptomes.

34 Using the comprehensive database M5nr, the assembly-based approach identifies genes with

35 only 0.6% false positives at thresholds ranging from inclusive to stringent compared to assembly-

36 free approach (3.6 to 15% false positives). Using specialized databases (Carbohydrate Active-

37 enzyme and Nitrogen Cycle) the assembly-based approach identifies and quantifies genes with

38 3-5x less false positives. We also evaluated the impact of both approaches on real-world datasets.

39 Based on this benchmarking we present a standardized and optimized workflow for identifying

40 functional genes from metatranscriptomes.

41 **Conclusions**

2

42  Our findings support the argument of assembling short reads into contigs before alignment to a

43  reference database, since this provides higher precision and minimizes false positives. By virtue

44  of the extensive benchmarking we also present the open source metatranscriptomics analysis

45  workflow <u>Co</u>mparative <u>M</u>etatranscriptomics <u>W</u>orkflow *CoMW*.

46 **Key Words**

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

48 **1. Introduction**

49 Metatranscriptomics provides an unprecedented insight to complex functional dynamics of

50 microbial communities in various environments. The method has been applied to study the

51 microbial activity in thawing permafrost and the related biogeochemical mechanisms

52 contributing to greenhouse gas emissions [1], and Gonzalez *et al.* [2] applied metatranscriptomics

53 to evaluate root microbiome response to soil contamination. The method is typically used to

54 identify, quantify and compare the functional response of microbial communities in natural

55 habitats or in relation to environmental or physio-chemical impacts.

56 Using high-throughput sequencing techniques such as Illumina, metatranscriptomics offers a non

57 PCR biased method for looking at transcriptional activity occurring within a complex and diverse

58 microbial population at a specific point in time [3]. However, curation and annotation of this

59 complex data has emerged as a major challenge. To date, several studies have used various

60 analytic workflows. Typically, short sequence reads are utilized, which can either be individually

61 aligned directly to external reference databases (hereafter "assembly-free") or assembled into

62 longer contiguous fragments (contigs) for alignment (hereafter "assembly-based"). Various

63 studies have used either of these two general approaches. For example, Jung *et al.* [4] used an

64 assembly-free approach (with BWA [5] to map reads to reference genomes of lactic acid bacterial

65 strains associated with the kimchi microbial community) while Poulsen *et al.* [6] used an

66 assembly-based approach (using SHE-RA [7] assembly before aligning to protein database).

67 Similarly, an open source pipeline developed by Martinez *et al.* [8] to analyze

68  metatranscriptomics data-sets also aligns short reads directly to the M5nr database [9] and

69  provides eggNOG annotation [10].  Most of the studies have used an assembly-free approach

70  [11] due to less computational expense in addition to lack of thorough comparison available.

71  Since no independent and direct comparison between these two alternative approaches has

72  been performed presently, various metatranscriptomics analysis approaches may at times

73  produce inconsistent observations, even if identical databases are used in the analysis. Thus,

74  standardization of computational analysis is necessary to enable further propagation of

75  metatranscriptomics approaches and their integration into microbial ecology research.

76  Benchmarking provides a critical view of the efficiency and precision of different workflows and

77  use of simulated communities for benchmarking enables the analysis to be independent of

78  experimental variation and biases [12].

79  Here, we compared the assembly-free vs. assembly-based approach using simulated datasets.

80  We evaluated the accuracy of both approaches using precision, recall and False Discovery Rates

81  (FDR) with three different databases ranging from a generic or inclusive to specialized database

82  dedicated to structurally or functionally related functional families: 1) M5nr: an inclusive and

83  comprehensive non-redundant protein database in combination with eggNOG hierarchical

84  annotation 2) Carbohydrate-Active Enzymes (CAZymes) [13]: a database dedicated to describing

85  the families of structurally-related catalytic and carbohydrate-binding modules of enzymes and

86  3) Nitrogen Cycling Database (NCycDB) [14] a specialized and manually curated database

87  covering only N cycle genes. In order to estimate the consistency and variance in the results

88  caused by the choice of approach we then applied them to real world metatranscriptomes from

89  microbial communities in 1) active-layer permafrost soil from Svalbard and 2) Ash impacted

90 Danish Forest soil. With the help of this comprehensive benchmarking and comparative analysis

91 we then standardized and developed an open source Comparative Metatranscriptomic Workflow

92 (*CoMW).*

93

## 94 2. Findings

### 95 *2.1 Evaluation*

96 In order to compare the performance of the assembly-based and assembly-free approaches, and

97 to standardize a workflow using either of these, we simulated community transcript data using

98 4943 full length genes provided by Martinez et al. [8]. We analyzed both approaches separately

99 and compared against direct annotation of full-length genes. The full-length genes were

100 annotated using all three databases (M5nr, CAZy and NCycDB) independently to classify them

101 into functional subsystems and gene families. Figure 1 shows detailed workflow of comparative

102 analysis using both approaches.

103

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

107

### 108 *2.1.1   Functional assignment*

**109 2.1.1.1   M5nr Alignment**

110 Full length genes of the simulated community dataset were aligned and identified into 671

111 unique eggNOG orthologs, belonging to 19 distinct functional subsystems (level II). At the default

112 confidence threshold (BTS score 50) of Diamond [15], assembly-free approach produced

113 alignments to 820 orthologs with a precision of 85% (14.9% FPs), whereas the assembly-based

6

114 approach identified 665 orthologs with a precision of 99.3% (0.6% FPs) at the default SWORD

115 [16] confidence threshold of 1E-5. Repeating the alignments using a gradient of 15 varying

116 confidence thresholds for each approach (Low - $T_L$, Medium - $T_M$ and High – $T_H$; 5 thresholds /

117 category) resulted in dissimilar performance for both approaches. The precision and recall of

118 assembly-based approach did not change from 99.3% and 98.5% respectively throughout all

119 categories whereas the assembly-free approach had a maximum precision of 96.3% at $T_M$ and

120 decreases to 85% at $T_L$ and $T_H$. The assembly-based approach also produced fewer (only 0.6%)

121 FPs consistently compared to assembly-free approach of FPs ranging from 14.9% to minimum

122 3.6% at highest precision. Based on F-Score the most optimal alignment for each approach is

123 given in Table 1, whereas detailed values for precision, recall, F-Score and FDR are listed in

124 Supplementary Table S1. We then also evaluated both approaches by selectively removing

125 sequences belonging to a certain functional subsystem from the M5nr database in a controlled

126 manner (segmented cross validation) in order to replicate real world metatranscriptomes where

127 a certain functional subsystem can be completely or partially absent from the reference

128 database. We removed four (level II) subsystems ("[D] Cell cycle control, cell division,

129 chromosome partitioning"; "[L] Replication, recombination and repair"; "[E] Amino acid

130 transport and metabolism" and "[R] General function prediction only" and "[S] Function

131 unknown"). The level II subsystems were removed one at a time realigning full-length genes and

132 simulated reads using both Assembly-based and assembly-free approaches to the cropped

133 database to compare identification consistency. In each validation round, the number of unique

134 (eggNOG) orthologs identified by the assembly-based approach were consistent to full length

135 gene alignment whereas from the assembly-free approach, the orthologs dropped significantly

136 along with its ability to recall TPs. Table 2 provides details for each validation cycle.

137

138 *Table 1 : Mean of Precision, Recall and F Score for both approaches against all three databases. Bold emphasizes better*
139 *precision, recall, F-Score and FDR in each category across approaches*

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

140

141 *Table 2 Selective removal of functional subsystems from eggnog database (segmented cross-validation) of approaches. Bold*
142 *emphasizes better consistency across approaches*

| | Full Length Genes Unique Orthologs | | Unique orthologs from the assembly-free approach | | Unique Orthologs from the Assembly-based approach | |
|---|---|---|---|---|---|---|
| Complete Database | 671 | | 784 | | 667 | |
| [D] removed | 628 | 93.59% | 572 | 72.95% | **624** | **93.55%** |
| [L] removed | 640 | 95.38% | 584 | 74.48% | **636** | **95.35%** |
| [E] removed | 640 | 95.38% | 583 | 74.36% | **636** | **95.35%** |
| [R], [S] removed | 347 | 51.7% | 334 | 42.60% | **352** | **52.77%** |

143

144 ### *2.1.1.2 CAZY Alignment*

145 Out of a total 2395 full length genes, 500 sequences were aligned to 395 unique functional genes

146 in the CAZY database, which belonged to 130 gene families and were further classified as 7

147 enzyme classes. Using default confidence thresholds (BTS 50 & 1E-5), the assembly-free approach

148 identified 765 functional genes belonging to 112 unique families and 6 enzyme classes with a

149 precision of 28.5% (71.4% FPs). The assembly-based approach identified 488 functional genes

150 from CAZY database that were classified into 147 gene families from 7 enzyme classes with a

151 precision of 66% (FDR 33.9%) at the default confidence threshold. However, when we repeated

152  the process with 15 various confidence thresholds, precision improved consistently for the

153  assembly-based approach and FPs decreased, whereas for assembly-free approach, precision

154  dropped significantly with increasing confidence threshold (see Table 1 and Supplementary Table

155  S2).

156  **2.1.1.3  *NCycDB Alignment***

157  410 of 2395 full length genes aligned to this database, identified as 29 unique Nitrogen cycle

158  genes and further belonging to 15 functional gene families in 5 pathways. Using default

159  confidence thresholds, the assembly-free approach identified 1541 functional genes belonging

160  to 25 functional gene families classified into 6 pathways with a precision of 0.9% (99% FPs). The

161  assembly-based approach identified 42 Nitrogen cycle genes classified into 25 gene families from

162  6 pathways with a precision of 59.5% (40.4% FPs) at a default confidence threshold of 1E-5. Like

163  comparisons against M5nr and CAZY we repeated the process with 15 different confidence

164  thresholds for each approach. Precision improved significantly for the assembly-based approach

165  at stringent thresholds whereas for the assembly-free approach, the best precision achieved was

166  5.8%. (Table 1, Supplementary Table S3).

167  **2.1.2  *Expression Quantification***

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

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

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

171  best identification thresholds for both approaches based on highest F-Score and performed

172  differential expression analysis using DESeq2 [16] algorithm in SARTools [17]. This analysis for

173  both approaches was carried out against all three databases using the most specific level of

174 hierarchy in the respective databases in order to capture their ability to quantify expression levels

175 of specific genes. According to full-length gene alignments against eggNOG, 123 genes were

176 significantly upregulated and 270 were significantly downregulated. Using assembly-free

177 approach with best F-Score 73 genes were up-regulated (precision 94.5%, 5.4% FPs) and 380

178 (precision 65.7%, 34.2% FPs) were down regulated. whereas in the assembly-based approach 99

179 (precision 94.9%, 5% FPs) genes were up-regulated and 249 (precision 97.1%, 2.8% FPs) down

180 regulated. For CAZy database full-length genes 81 and 189 genes were significantly up and down

181 regulated respectively. Using assembly-free approach 31 (precision 19.3%, 80.6% FPs) genes

182 significantly up regulated and 137 genes (precision 52.5%, 47.4% FPs) whereas the assembly-

183 based approach 83 (precision 71%, 28.9% FPs) genes were up-regulated and 191 (precision

184 73.8%, 26.1% FPs) genes were down regulated. In the NCyc database expression analysis, 3 and

185 14 genes were significantly up and down-regulated respectively.  Using assembly-free approach

186 26 (precision 0%, 100% FPs) and 107 (precision 4.6%, 95.3% FPs) genes up and down regulated

187 respectively, whereas using assembly-based approach 3 (precision 33.3%, 66.6% FPs) genes up

188 regulated and 18 (precision 55.5%, 44% FPs) were down regulated. Precision, Recall and FDR for

189 both approaches against all three databases are available in Supplementary Table S4.

190 Additionally, we then collapsed the functional genes into functional subsystems and gene

191 families to remove FPs produced due to identification of homologous proteins or proteins with

192 multiple inheritance. Fold change (log2 transformed) was then calculated for each

193 subsystem/gene family. (see Figure 3)

194

195 *Figure 2: Differential Expression comparison of Assembly-free and Assembly-based approaches using A) M5nr*
196 *database, B) NCycDB and C) CAZy database.*

197

### 2.1.3   Real-World metatranscriptomes

198

199  To evaluate the effect of the two approaches on real world data, two metatranscriptomes from

200  microbial communities were studied. The first study investigated the transcriptional response

201  during warming from -10 °C to 2 °C and subsequent cooling of 2 °C to -10 °C of an Arctic tundra

202  active layer soil from Svalbard, Norway . The aim of the study was to understand taxonomic and

203  functional shifts in microbial communities caused by climate change in the Arctic. A pronounced

204  shift during the incubation period was noticed by Schostag *et al.* [17] (under review Molecular

205  Ecology, SRA Bio-Project: PRJNA417839) which was not replicated by the assembly-free

206  approach. However, the assembly-based approach identified an increase of genes in "[P]

207  Inorganic ion transport and metabolism". For cooling, the assembly-based approach also

208  captured the upregulation and downregulation of genes related to "[J] Translation, ribosomal

209  structure and biogenesis" and "[C] Energy production and conversion" respectively (Figure 6)

210  unlike assembly-free approach. These findings may have implications for our understanding of

211  carbon dioxide emission, nitrogen cycling and plant nutrient availability in Arctic soils.

212

213  *Figure 3: Functional expression dynamics in Arctic permafrost soil identified against eggNOG functional subsystems*
214  *using Assembly-based and Assembly-free approach*

215

216  In the second study, Bang-Andreasen *et al.* [18] (under review ISME, SRA Bio-Project:

217  PRJNA512608) investigated the effects of wood ash amendment on Danish forest soils. Ash was

218  added in 3 different quantities (0/control, 3, 12 and 90 tonnes ash per hectare (t ha$^{-1}$)). In addition

219  to ash concentration, the effect over time was analysed in soil communities at 0, 3, 30 and 100

220 days after ash addition. This resulted in strong effects on functional expression as seen in Figure

221 7. Both approaches once again displayed varying results such as changes in genes related to

222 eggNOG functional subsystem "[W] Extracellular structures". Assembly-free approach also

223 identified 75% of genes as "[S] Function unknown" consistently unlike assembly-based.

224

225 *Figure 4: Functional expression dynamics in Danish forest soil due to Ash amendment and time elapsed, identified*
226 *against eggNOG functional subsystems using Assembly-based and Assembly-free approach*

227

## 2.2 Standardized Workflow (CoMW)

229 By the virtue of thorough benchmarking we standardized, implemented, and validated a

230 metatranscriptomic workflow (CoMW) using Assembly-based approach. The workflow was

231 implemented by keeping in mind the databases and tools for each step are ever improving thus

232 optional steps can be skipped, changed or even improved in a structural manner. CoMW is open

233 source workflow written in python available at (https://github.com/anwarMZ/CoMW). It is based

234 on four major steps: 1) Assembly and Mapping short reads to assembled contigs; 2) Filtering of

235 contigs; 3) Gene Prediction and Alignment and 4) Annotation. These scripts make each step of

236 the workflow straightforward and help to make these complex analyses more reproducible and

237 the components re-useable in different contexts. Help regarding input, output and parameters is

238 provided with each script and an overall tutorial is presented in the data repository at GitHub.

239 We here wanted to build an open source work flow for metatranscriptomics analysis that can

240 assist in analyzing large metatranscriptomics data. Processes like ORF detection, alignment

241 against the database and calculations of the gene expression are vital in any metatranscriptomic

242 analyses and are, therefore, present uniformly in all workflows. However, since we use the

12

243 assembly-based workflow where we assemble the reads into longer contigs we also propose a 2-

244 step filtering of the contigs to remove any chimeric or false contig made as a result of assembly

245 or sequencing error by removing contigs that have an expression level less than a specific

246 threshold and to remove any potential non-coding RNA contigs assembled.

247 *Assembly and Mapping of short reads back to assembled contigs* is done using Trinity [19] and

248 BWA [5] respectively. Various tools have been developed for metatranscriptome reconstruction

249 that usually rely on graph-theory. Trinity however generates the most optimal assemblies for

250 coding RNA reads [11,20,21]. However, the user can generate contigs by any assembler preferred

251 but it can reduce the quality of the following steps such as alignment of contigs.

252 *Filtering of Contigs* is done to remove variance in sequences/samples. We can filter contig

253 abundance data by removing all contigs with relative expression lower than a specific cutoff, e.g.

254 1% (selected based on dataset variance) of the number of sequences in the dataset with least

255 number of sequences. This threshold is also flexible for different datasets and in some cases not

256 required at all so CoMW allows user to bypass this step or change the threshold up and down

257 based on data variation. The filtered contigs are subject to potential non-coding RNA filtration by

258 aligning them against the RFam database [22] using infernal [23] which is a secondary-structure-aware

259 aligner that predicts the secondary structure of RNA sequences and similarities based on the consensus

260 structure models. Once again, the ncRNA filtering is an optional step in CoMW, though highly

261 recommended in order to reduce FPs.

262 *Gene Prediction and Alignment* Transeq from EMBOSS [24] is used to predict probable open

263 reading frames (ORFs) of the contigs (customizable, by default 6 per contig). We used SWORD

264 [16] as alignment tool against reference databases. SWORD can be used in parallel based on

13

265 computational resources available and the aligned results are parsed and cut-off at a specific

266 confidence threshold of combination of e-value and alignment length (usually 1e-5, can be

267 changed given the assembly distribution in datasets).

268 *Annotation* of aligned transcripts from the previous step can be done using the databases such

269 as eggNOG which is a hierarchically structured annotation using a graph-based unsupervised

270 clustering available algorithm to produce genome wide orthology inferences. Aligned proteins

271 are then placed into functional subsystems based on their best hits.), CAZy which is a knowledge-

272 based resource specialized in the Glycogenomics, and NCycDB; a Nitrogen cycle database. This results

273 in a count table with a contig and eggNOG ortholog or CAZy gene or NCyc gene having a certain

274 count from each sample depending upon database used. This count table can be then used for

275 differential expression using state-of-the-art expression analysis.

276 CoMW is based on the results and findings from comparison of approaches. However, it has

277 multiple optional steps such as abundance based and non-coding RNA filtering which can be

278 different in data sets from a different environment. Similarly, the scripts are designed to cater

279 more than one assembler output to enable diverse range of environments to be studied.

280

## 3    Discussion

282 The application of metatranscriptomics is less common than other DNA-based genomics

283 techniques and thus most analysis pipelines are built ad hoc. The majority of these pipelines

284 follow the assembly-free approach [11] such as COMAN [25], Metatrans [8], and SAMSA2 [26].

285 The lack of thorough benchmarking studies and standardized workflows in metatranscriptomics

286 has made it a more challenging task to analyze the typically big datasets produced. Previous

14

287 studies have compared *de novo* sequence assemblers including Trinity, MetaVelvet [27],

288 Oases[28], AbySS [29] and SOAPden-ovo [30] but an independent comparison of the two

289 different approaches based on including assembly or directly aligning reads (here "assembly-

290 free") has been lacking. We have attempted to assist this decision-making for processing

291 metatranscriptomic analysis by independently assessing the performance of the two for

292 functional annotation and expression quantification against three databases ranging from

293 inclusive to specialized.

294 With simulated samples comprised of genes collected from abundant genomes provided by

295 Martinez *et al.* [8] we show that both approaches provide high recall rates against the general

296 comprehensive database M5nr. However, the assembly-based approach provided a significantly

297 better precision for identification and quantification. For relatively compact and specialized

298 databases, recall and precision drop for both approaches (especially for the most compact

299 database NCyc). However, the assembly-based approach still appeared to be more precise,

300 meaning that fewer genes were mis-assigned against these database and significantly lower FPs

301 were produced. The precision in identification and expressional fold change comparison of gene

302 families and functional subsystems for simulated samples against all three databases confirmed

303 that while an assembly step is challenging computationally it holds the potential to reveal

304 information regarding the gene expressions that is not attainable without it.

305 Selecting a single best workflow or pipeline for all types of metatranscriptomics studies is not a

306 straightforward affair, and we believe that choice of approach changes the outcome of study

307 significantly as observed with real-world datasets from active-layer permafrost soil from Svalbard

308 and Ash impacted Danish Forest soil.  In addition to choosing the right workflow, combining that

309 with the appropriate reference database is equally important to ensure the best annotation

310 performance. With databases specialized for one or more specific environments or functional

311 categories assembly-free approach underperforms due to its inability to identify conserved

312 sequences in reference database. We also show that assembly-free approach can increase the

313 rate of FPs in annotation when a database is dominant in specific functional subsystem or does

314 not possess certain category which can also lead to wrong estimation of fold change in expression

315 In summary, we show that the choice of approach (assembly-free or assembly-based) and

316 database significantly affects the quality of the identification, annotation and expression results.

317 Given the impact of each of these variables, it is inevitable that it significantly affects the results

318 of an individual study and comparison of across studies. By standardizing and developing CoMW

319 we believe our work presented here further assists the microbial ecology research community to

320 make more informed decisions about the most appropriate methodological approach to analyze

321 large metatranscriptomic datasets with improved precision.

322

## 4   Methods

324 For the assembly-free approach we used the Metatrans pipeline [8], which uses FragGeneScan

325 [31] for ORF predictions in short reads, CD-Hit [32] for gene clustering and Diamond [15] for

326 alignment to the M5nr database. For assembly-based approach we assembled the simulated

327 short reads using Trinity [19] which has been studied to outperform other de novo RNA-seq

328 assemblers  and aligned the resulting contigs using SWORD (an efficient protein database search

329 implementation especially optimized for large databases) [16] against the M5nr [9], CAZy [13]

330 and NCyc [14] database. We then wrote an annotation script which Is included in CoMW. For

16

331  expression analysis gene counts were normalized between samples using the DESeq2 [33]

332  algorithm. Significantly differentially expressed genes were analyzed in SARTools [34] using

333  parametric relationship and p-value 0.05 as significance threshold. The Benjamini and Hochberg

334  correction procedure [35] was used to adjust p-value.

335  **4.1    Composition of Simulated Communities**

336  In this study we utilised a set of simulated communities from Martinez *et al*. [8] where they

337  collected 4943 genes from five abundant microbial genomes: *Bacteroides vulgatus* ATCC 8482,

338  *Ruminococcus torques* L2-14, *Faecalibacterium prausnitzii* SL3/3, *Bacteroides thetaiotaomicron*

339  VPI-5482 and *Parabacteroides distasonis* ATCC 8503. We simulated short reads into 100 samples

340  using Polyester [36]  embedded in a script provided by Martinez *et al.* [8] at coverage of 20x

341  which resulted in a count table and short reads with 2395 genes to add the impact of sequencing

342  coverage. Their abundance was then regulated up and down and by knocking out few genes in a

343  controlled manner in order to make the composition similar to real world metatranscriptomes.

344  The process of regulation of abundance was done by first dividing the 100 samples into 2 groups

345  ("A" and "B") and then increasing the abundance of 10% genes up to 4-fold, decreasing the

346  abundance of another 10% of the genes 4-fold and completely removing another 5% of the genes

347  from both simulated reads and count tables. The process of selection of samples and genes was

348  random but tracked. To include quality bias, we used the ART simulator [37] to produce an equal

349  number of reads in FASTQ format to those produced by Polyester. ART was initially trained with

350  Hi-Seq 2500 Illumina quality error model from dataset discussed above to have a consistent error

351  bias. After simulating FASTQ files we then extracted the quality data and bound it to the FASTA

352 files generating new FASTQ files. With the coverage bias and quality training included we had a

353 total of 62,035,912 reads (310,179 ± 3,454 reads/sample).

354 **4.2   Evaluation Measures**

355 We used the standard measures of precision (also named positive predictive value, PPV),

356 accounting for how many annotations and identifications of significantly differentially expressed

357 gene families and subsystems are correct and defined as $\frac{TP}{TP+FP}$ and recall (also named sensitivity

358 or true positive rate, TPR), accounting for how many correct annotations are selected, defined as

359 $\frac{TP}{TP+FN}$ where TP indicates the number of orthologs that have been correctly annotated, FN

360 indicates the number of orthologs/genes/functional subsystem which are in the simulated

361 communities but were not found by a certain approach and FP indicates the number of

362 orthologs/genes/functional subsystem that have been wrongly annotated (because they do not

363 appear in the simulated communities). The F-score is the harmonic mean of precision and recall,

364 defined as $2 * \frac{Precision*Recall}{Precision+Recall}$

365 **Availability of source code and requirements**

366 - Project name: Comparative Metatranscriptomics Workflow [*CoMW*]

367 - Project home page:  https://github.com/anwarMZ/CoMW

368 - Operating system(s): Platform independent

369 - Programming language: Python, R, and bash

370 - Other requirements: Requirements mentioned in detailed manual at GitHub

371 - License: GNU General Public License v3.0

372 **Availability of supporting data and materials**

373 Raw sequence data generated using simulation of full-length genes were deposited in the NCBI

374 Sequence Read Archive and are accessible through BioProject accession number PRJNA509064

375 All databases can be accessed in one place at http://tiny.cc/CoMW_DBs

376 Supplementary File 1 – Precision Recall Analysis of both approaches

377 Supplementary File 2 – Differential Expression Analysis of all approaches using eggNOG

378 database

379 Supplementary File 3 – Differential Expression Analysis of all approaches using CAZy database

380 Supplementary File 4 – Differential Expression Analysis of all approaches using NCyc database

381 **Declarations**

382 *--*

383 *List of abbreviations*

384 FDR: False Discovery Rate, FP: False Positives, TP: True Positives, FN: False Negatives, mRNA:

385 messenger RNA

386 *Ethical Approval (optional)*

19

387  Not applicable

388  ***Consent for publication***

389  Not applicable

390  ***Competing Interests***

391  The authors declare that they have no competing interests.

395  ***Author's Contributions***

396  MZA & CSJ conceived and designed the study. MZA, TBA and AL carried out the data

397  production. MZA and AL carried out analysis. MZA drafted the manuscript and AL, TBA and CSJ

398  revised and approved the final version.

402  ***Authors' information (optional)***

403

## References

404  **References**

405  1. Coolen MJL, Orsi WD. The transcriptional response of microbial communities in thawing Alaskan
406  permafrost soils. Front Microbiol [Internet]. 2015 [cited 2018 Oct 25];6. Available from:
407  https://www.frontiersin.org/articles/10.3389/fmicb.2015.00197/full

408  2. Gonzalez E, Pitre FE, Pagé AP, Marleau J, Guidi Nissim W, St-Arnaud M, et al. Trees, fungi and bacteria:
409  tripartite metatranscriptomics of a root microbiome responding to soil contamination. Microbiome.
410  2018;6:53.

411  3. Leimena MM, Ramiro-Garcia J, Davids M, van den Bogert B, Smidt H, Smid EJ, et al. A comprehensive
412  metatranscriptome analysis pipeline and its validation using human small intestine microbiota datasets.
413  BMC Genomics. 2013;14:530.

414  4. Jung JY, Lee SH, Jin HM, Hahn Y, Madsen EL, Jeon CO. Metatranscriptomic analysis of lactic acid
415  bacterial gene expression during kimchi fermentation. Int J Food Microbiol. 2013;163:171–9.

416  5. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA
417  sequences to the human genome. Genome Biol. 2009;10:R25.

418  6. Poulsen M, Schwab C, Jensen BB, Engberg RM, Spang A, Canibe N, et al. Methylotrophic
419  methanogenic Thermoplasmata implicated in reduced methane emissions from bovine rumen. Nat
420  Commun. 2013;4:1428.

421  7. Rodrigue S, Materna AC, Timberlake SC, Blackburn MC, Malmstrom RR, Alm EJ, et al. Unlocking Short
422  Read Sequencing for Metagenomics. PLOS ONE. 2010;5:e11840.

423  8. Martinez X, Pozuelo M, Pascal V, Campos D, Gut I, Gut M, et al. MetaTrans: an open-source pipeline
424  for metatranscriptomics. Sci Rep. 2016;6:26447.

425  9. Wilke A, Harrison T, Wilkening J, Field D, Glass EM, Kyrpides N, et al. The M5nr: a novel non-
426  redundant database containing protein sequences and annotations from multiple sources and
427  associated tools. BMC Bioinformatics. 2012;13:141.

428  10. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, et al. eggNOG 4.5: a
429  hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and
430  viral sequences. Nucleic Acids Res. 2016;44:D286–93.

431  11. Aguiar-Pulido V, Huang W, Suarez-Ulloa V, Cickovski T, Mathee K, Narasimhan G. Metagenomics,
432  Metatranscriptomics, and Metabolomics Approaches for Microbiome Analysis. Evol Bioinforma Online.
433  2016;12:5–16.

434  12. Almeida A, Mitchell AL, Tarkowska A, Finn RD. Benchmarking taxonomic assignments based on 16S
435  rRNA gene profiling of the microbiota from commonly sampled environments. GigaScience. 2018;7.

436  13. Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. The Carbohydrate-Active
437  EnZymes database (CAZy): an expert resource for Glycogenomics. Nucleic Acids Res. 2009;37:D233-238.

438 14. Tu Q, Lin L, Cheng L, Deng Y, He Z. NCycDB: a curated integrative database for fast and accurate
439 metagenomic profiling of nitrogen cycling genes. Bioinforma Oxf Engl. 2018;

440 15. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. Nat Methods.
441 2015;12:59–60.

442 16. Vaser R, Pavlović D, Šikić M. SWORD—a highly efficient protein database search. Bioinformatics.
443 2016;32:i680–4.

444 17. Schostag MD, Anwar MZ, Jacobsen CS, Larose C, Vogel TM, Maccario L, et al. Transcriptomic
445 responses to warming and cooling of microbial communities in an Arctic tundra soil. Moleuclar Ecol.
446 Under Review;

447 18. Bang-Andreasen T, Anwar MZ, Lanzen A, Kjoller R, Ronn R, Ekelund F, et al. Total RNA-sequencing
448 reveals complex taxonomic and functional responses to wood ash application in agricultural and forest
449 soil multi-level microbial communities. ISME. Under Review;

450 19. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Trinity: reconstructing a full-
451 length transcriptome without a genome from RNA-Seq data. Nat Biotechnol. 2011;29:644–52.

452 20. Celaj A, Markle J, Danska J, Parkinson J. Comparison of assembly algorithms for improving rate of
453 metatranscriptomic functional annotation. Microbiome. 2014;2:39.

454 21. Lau MCY, Harris RL, Oh Y, Yi MJ, Behmard A, Onstott TC. Taxonomic and Functional Compositions
455 Impacted by the Quality of Metatranscriptomic Assemblies. Front Microbiol [Internet]. 2018 [cited 2018
456 Oct 25];9. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6019464/

457 22. Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR. Rfam: an RNA family database. Nucleic
458 Acids Res. 2003;31:439–41.

459 23. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. Bioinformatics.
460 2013;29:2933–5.

461 24. Rice P, Longden I, Bleasby A. EMBOSS: The European Molecular Biology Open Software Suite. Trends
462 Genet. 2000;16:276–7.

463 25. Ni Y, Li J, Panagiotou G. COMAN: a web server for comprehensive metatranscriptomics analysis. BMC
464 Genomics. 2016;17:622.

465 26. Westreich ST, Treiber ML, Mills DA, Korf I, Lemay DG. SAMSA2: a standalone metatranscriptome
466 analysis pipeline. BMC Bioinformatics. 2018;19:175.

467 27. Namiki T, Hachiya T, Tanaka H, Sakakibara Y. MetaVelvet: an extension of Velvet assembler to de
468 novo metagenome assembly from short sequence reads. Nucleic Acids Res. 2012;40:e155.

469 28. Schulz MH, Zerbino DR, Vingron M, Birney E. Oases: robust de novo RNA-seq assembly across the
470 dynamic range of expression levels. Bioinformatics. 2012;28:1086–92.

471    29. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol İ. ABySS: A parallel assembler for short
472    read sequence data. Genome Res. 2009;19:1117–23.

473    30. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-
474    efficient short-read de novo assembler. GigaScience. 2012;1:18.

475    31. Rho M, Tang H, Ye Y. FragGeneScan: predicting genes in short and error-prone reads. Nucleic Acids
476    Res. 2010;38:e191.

477    32. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or
478    nucleotide sequences. Bioinforma Oxf Engl. 2006;22:1658–9.

479    33. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data
480    with DESeq2. Genome Biol [Internet]. 2014 [cited 2018 Oct 25];15. Available from:
481    https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4302049/

482    34. Varet H, Brillet-Guéguen L, Coppée J-Y, Dillies M-A. SARTools: A DESeq2- and EdgeR-Based R Pipeline
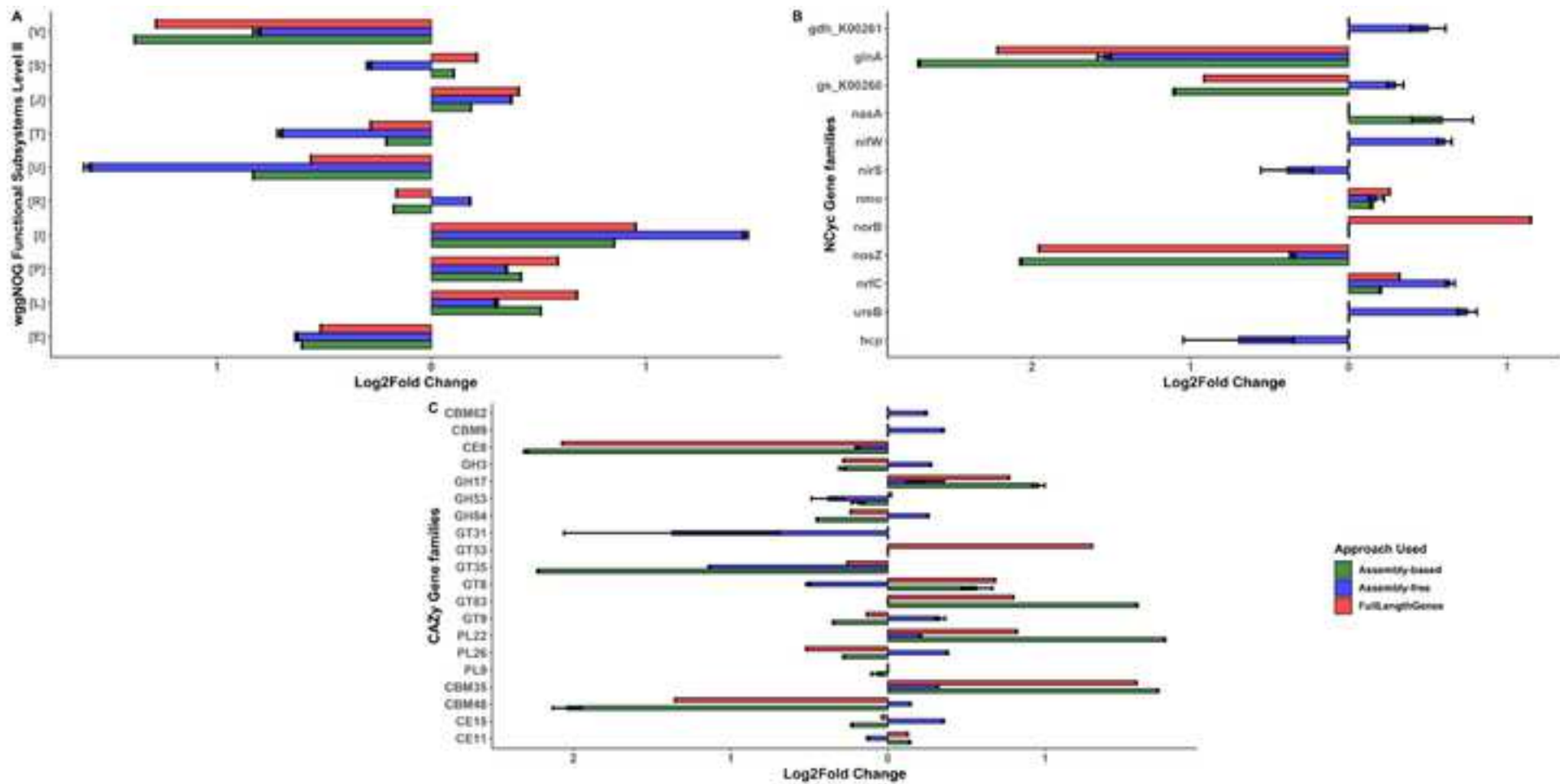483    for Comprehensive Differential Analysis of RNA-Seq Data. PLOS ONE. 2016;11:e0157022.

484    35. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to
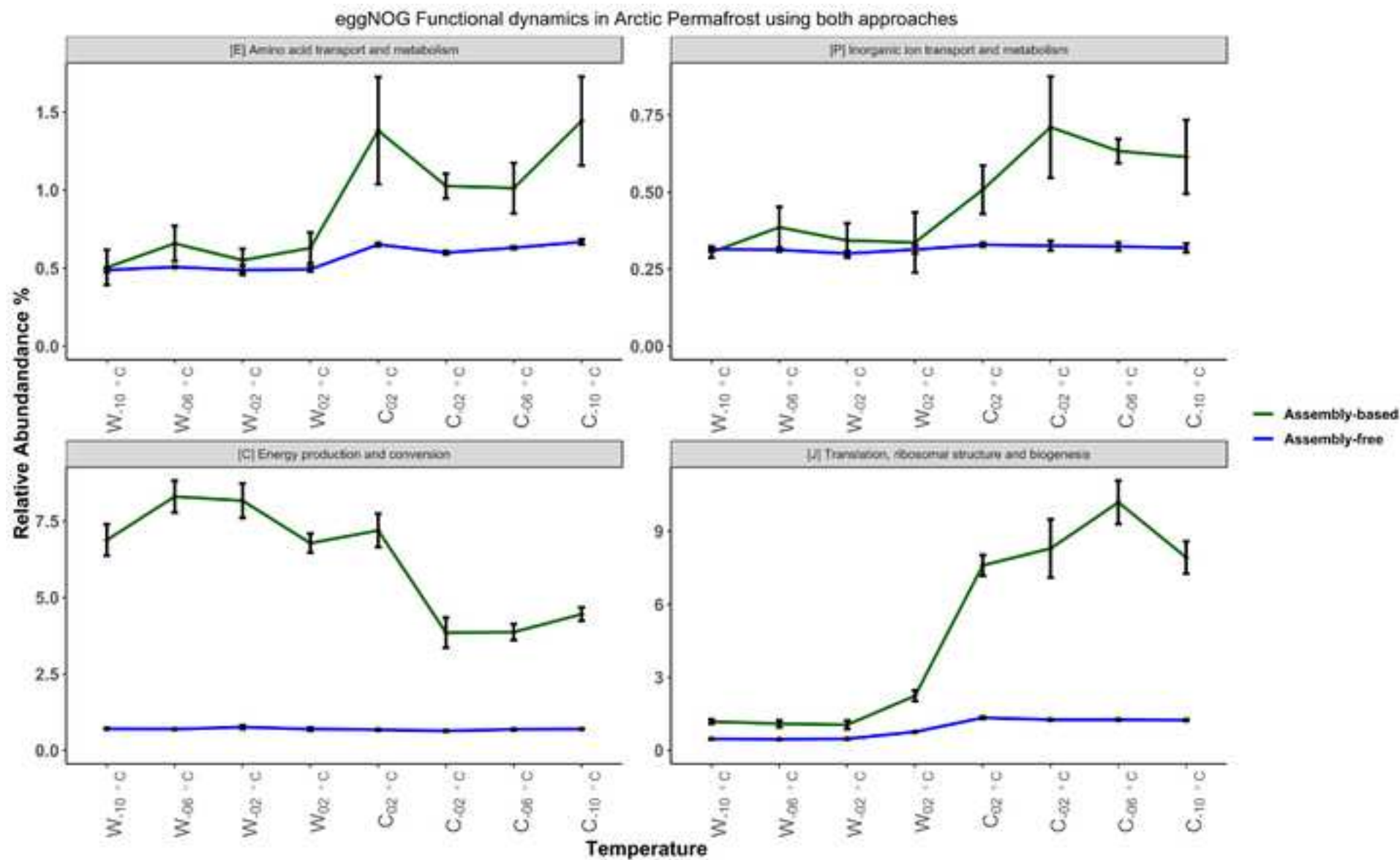485    Multiple Testing. J R Stat Soc Ser B Methodol. 1995;57:289–300.

486    36. Frazee AC, Jaffe AE, Langmead B, Leek JT. Polyester: simulating RNA-seq datasets with differential
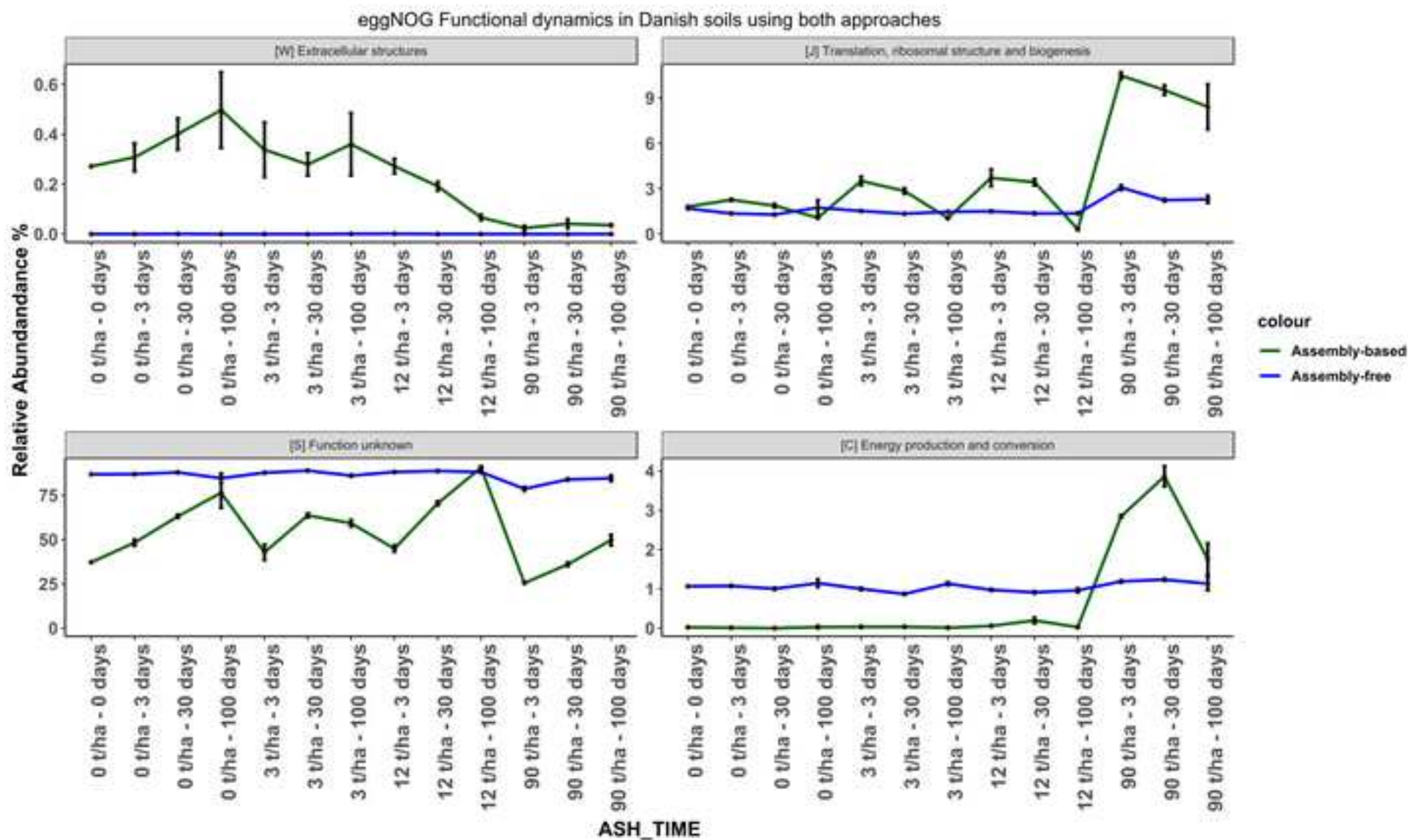487    transcript expression. Bioinformatics. 2015;31:2778–84.

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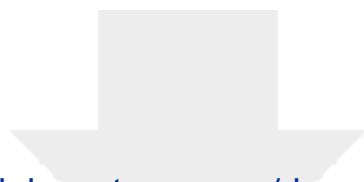489    Bioinformatics. 2012;28:593–4.

23

Figure1

Figure2

Click here to download Figure Figure 2.tiff

Figure3

Click here to download Figure Figure 3.tiff ⬇



eggNOG Functional dynamics in Arctic Permafrost using both approaches

Figure4

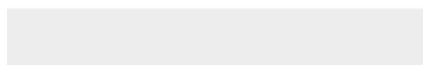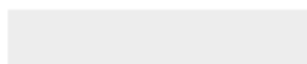eggNOG Functional dynamics in Danish soils using both approaches

Click here to access/download
**Supplementary Material**
SupplementaryFile1_PrecisionRecall.docx

Click here to access/download

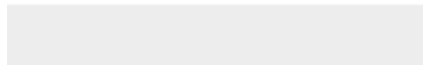**Supplementary Material**

SupplementaryFile2_eggNOG_DEAnalysis.xlsx

SupplementaryFile3

SupplementaryFile4

Click here to access/download
**Supplementary Material**
SupplementaryFile4_NCyC_DEAnalysis.xlsx

**AARHUS**
**UNIVERSITY**
DEPARTMENT OF ENVIRONMENTAL SCIENCE

To: Editor, *GigaScience*

Re: Submission of manuscript entitled "To assemble or not to resemble – benchmarking of metatranscriptomic practices and a validated Comparative Metatranscriptomics Workflow (CoMW)" to *GigaScience*
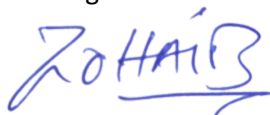
Metatranscriptomics has recently gained popularity, thanks to its ability to uncover the active functional profiles of microbial communities. Being a relatively recent approach, there are still several analytical obstacles that limit its large-scale application. Sequence reference databases are also limited in their coverage and thus the use of different workflows and databases can lead to different outcomes, rendering it difficult to compare results between independent studies.

We have conducted a comprehensive comparison of workflows representing the two main alternatives for metatranscriptome analysis, namely assembly-based or assembly-free. This comparison was done using both simulated datasets and real world metatranscriptomes using three different hierarchical databases. To the best of our knowledge this is first independent comparison of these alternatives that will assist decision making and analysis of metatranscriptomics. Subsequently we also present a validated workflow using assembly-based analysis, which provided the best results according to simulated datasets.

We believe that GigaScience would be an outstanding forum for this manuscript, due to its intention of featuring interdisciplinary research; it would be of interest both to microbial ecologists, clinical microbiologists and bioinformaticians. To maintain open data and transparency in our benchmarking we have made all code, test data, results and supporting documents for CoMW available at different links provided within manuscript as per GigaScience policies. This manuscript presents material that has not previously been published and is not under consideration for publication elsewhere and all authors have seen and approved the final version submitted.

Thank you again for considering our manuscript.

Kind regards

Muhammad Zohaib Anwar
On behalf of Authors

**Environmental microbiology & biotechnology**

**Muhammad Zohaib Anwar**
PhD Student

Date: 10 January 2019

E-mail:
mzanwar@envs.au.dk
Web:
au.dk/en/mzanwar@envs

Sender's CVR no.:
31119103

Page 1/1

**Environmental microbiology & biotechnology**
Aarhus University
Frederiksborgvej 399
PO box 358
DK-4000 Roskilde
Denmark

Tel.: +45 8715 0000
Fax: +45 8715 5010
E-mail: envs@au.dk
Web: envs.au.dk/en