

1 Metagenomics workflow for hybrid assembly, differential 2 coverage binning, transcriptomics and pathway analysis 3 (MUFFIN)

4 Renaud Van Damme^{1,2}, Martin Hölzer⁴, Adrian Viehweger^{3,4}, Bettina Müller¹, Erik Bongcam-
5 Rudloff², Christian Brandt^{2,5}

6
7 1. Department of Molecular Sciences, Swedish University of Agricultural Sciences, Uppsala, Sweden

8 2. Dept. Animal Breeding and Genetics, Bioinformatics section. Swedish University of Agricultural,
9 Sciences, Uppsala, Sweden

10 3. Department of Medical Microbiology, University Hospital Leipzig, Germany

11 4. RNA Bioinformatics and High-Throughput Analysis, Friedrich Schiller University Jena, Jena,
12 Germany

13 5. Institute for Infectious Diseases and Infection Control, Jena University Hospital, Jena, Germany
14

15 Abstract

16 Metagenomics has redefined many areas of microbiology. However, metagenome-
17 assembled genomes (MAGs) are often fragmented, primarily when sequencing was
18 performed with short reads. Recent long-read sequencing technologies promise to improve
19 genome reconstruction. However, the integration of two different sequencing modalities
20 makes downstream analyses complex. We, therefore, developed MUFFIN, a complete
21 metagenomic workflow that uses short and long reads to produce high-quality bins and their
22 annotations. The workflow is written by using Nextflow, a workflow orchestration software, to
23 achieve high reproducibility and fast and straightforward use. This workflow also produces
24 the taxonomic classification and KEGG pathways of the bins and can be further used by
25 providing RNA-Seq data (optionally) for quantification and annotation. We tested the
26 workflow using twenty biogas reactor samples and assessed the capacity of MUFFIN to
27 process and output relevant files needed to analyze the microbial community and their
28 function. MUFFIN produces functional pathway predictions and if provided *de novo* transcript
29 annotations across the metagenomic sample and for each bin.

30 Author Summary

31 RVD did the development and design of MUFFIN and wrote the first draft; BM and EBR did
32 the critical reading and correction of the manuscript; MH did the critical reading of the
33 manuscript and the general adjustments for the metagenomic workflow; AV did the critical
34 reading of the manuscript and adjustments for the taxonomic classifications. CB supervised
35 the project, did the workflow design, helped with the implementation, and revised the
36 manuscript.

37 Introduction

38 Metagenomics is widely used to analyze the composition, structure, and dynamics of
39 microbial communities, as it provides deep insights into uncultivable organisms and their
40 relationship to each other¹⁻⁵. In this context, whole metagenome sequencing is mainly
41 performed using short-read sequencing technologies, predominantly provided by Illumina.
42 Not surprisingly, the vast majority of tools and workflows for the analysis of metagenomic
43 samples are designed around short reads. However, long-read sequencing technologies
44 such as provided by PacBio or Oxford Nanopore Technologies (ONT) retrieve genomes from
45 metagenomic datasets with higher completeness and less contamination⁶. The long-read
46 information bridges gaps in a short-read-only assembly that often occur due to intra- and
47 interspecies repeats⁶. Complete viral genomes can be already identified from environmental
48 samples without any assembly step via nanopore-based sequencing⁷. Combined with a
49 reduction in cost per gigabase⁸ and an increase in data output, the technologies for
50 sequencing long reads quickly became suitable for metagenomic analysis⁹⁻¹². In particular,
51 with the MinION, ONT offers mobile and cost-effective sequencing device for long reads that
52 paves the way for the real-time analysis of metagenomic samples. Currently, the combination
53 of both worlds (long reads and high-precision short reads) allows the reconstruction of more
54 complete and more accurate metagenome-assembled genomes (MAGs)⁶.

55 One of the main challenges and bottlenecks of current metagenome sequencing studies is
56 the orchestration of various computational tools into stable and reproducible workflows to

57 analyze the data. A recent study from 2019 involving 24,490 bioinformatics software
58 resources showed that 26 % of all these resources are not currently online accessible ¹³.
59 Among 99 randomly selected tools, 49 % were deemed 'difficult to install,' and 28 %
60 ultimately failed the installation procedure. For a large-scale metagenomics study, various
61 tools are needed to analyze the data comprehensively. Thus, already during the installation
62 procedure, various issues arise related to missing system libraries, conflicting dependencies
63 and environments or operating system incompatibilities. Even more complicating,
64 metagenomic workflows are computing intense and need to be compatible with high-
65 performance compute clusters (HPCs), and thus different workload managers such as
66 SLURM or LSF. We combined the workflow manager Nextflow¹⁴ with virtualization software
67 (so-called 'containers') to generate reproducible results in various working environments and
68 allow full parallelization of the workload on a higher degree.

69 Several workflows for metagenomic analyses have been published, including
70 MetaWRAP(v1.2.1)¹⁵, Anvi'o¹⁶, SAMSA2¹⁷, Humann¹⁸, or MG-Rast¹⁹. Unlike those, MUFFIN
71 allows for a hybrid metagenomic approach combining the strengths of short and long reads.
72 It ensures reproducibility through the use of a workflow manager and reliance on either install
73 recipes (Conda ²⁰) or containers (Docker²¹).

74 Design and implementation

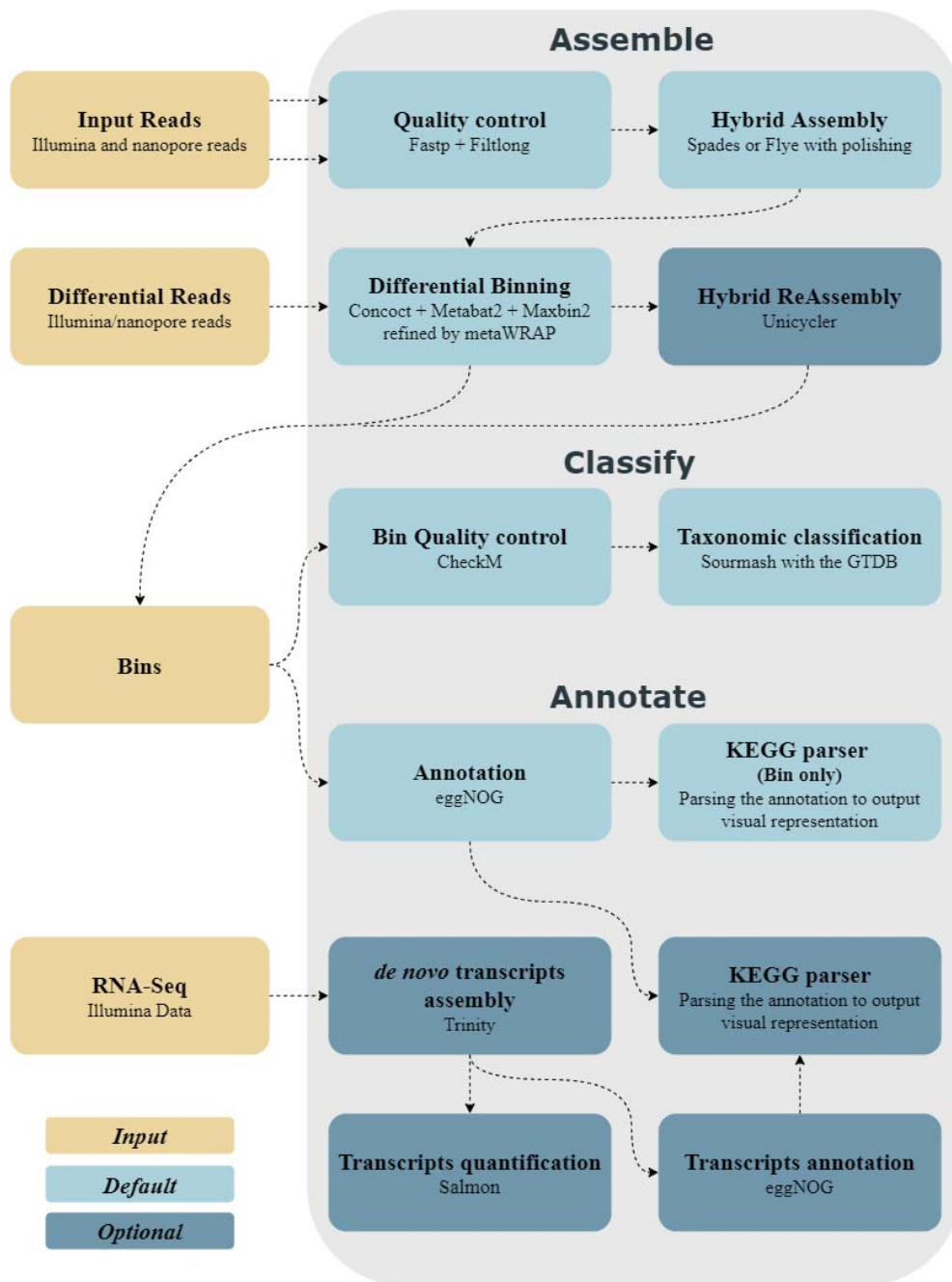
75 MUFFIN integrates state-of-the-art bioinformatic tools via Conda recipes or Docker
76 containers for the processing of metagenomic sequences in a Nextflow workflow
77 environment (Figure 1). MUFFIN executes three steps subsequently or separately if
78 intermediate results, such as MAGs, are available. As a result, a more flexible workflow
79 execution is possible. The three steps represent common metagenomic analysis tasks and
80 are summarized in Figure 1:

- 81 1. Assemble: Hybrid assembly and binning
- 82 2. Classify: Bin quality control and taxonomic assessment

83 3. Annotate: Bin annotation and KEGG pathway summary

84 The workflow takes paired-end Illumina reads (short reads) and nanopore-based reads (long
85 reads) as input for the assembly and binning and allows for additional user-provided read
86 sets for differential coverage binning. Differential coverage binning facilitates genome bins
87 with higher completeness than other currently used methods²². Step 2 will be executed
88 automatically after the assembly and binning procedure or can be executed independently by
89 providing MUFFIN a directory containing MAGs in FASTA format. In step 3, paired-end RNA-
90 Seq data can be optionally supplemented to improve the annotation of bins.

91 On completion, MUFFIN provides various outputs such as the MAGs, KEGG pathways, and
92 bin quality/annotations. Additionally, all mandatory databases are automatically downloaded
93 and stored in the working directory or can be alternatively provided via an input flag.



94

95 *Figure 1: Simplified overview of the MUFFIN workflow. All three steps (Assemble, Classify, Annotate) from top to*
 96 *bottom are shown. The RNA-Seq data for Step 3 (Annotate) is optional.*

97 **Step 1 - Assemble: Hybrid assembly and binning**

98 The first step (**Assembly and binning**), uses metagenomic nanopore-based long reads and

99 Illumina paired-end short reads to obtain high-quality and highly complete bins. The short-

100 read quality control is operated using fastp (v0.20.0)²³. Optionally, Filtlong (v0.2.0)²⁴ can be

101 used to discard long reads below a length of 1000 bp²⁴. The hybrid assembly can be
102 performed according to two principles, which differ substantially in the read set to begin with.
103 The default approach starts from a short-read assembly where contigs are bridged via the
104 long reads using metaSPAdes (v3.13.1)²⁵⁻²⁷. Alternatively, MUFFIN can be executed starting
105 from a long-read-only assembly using metaFlye (v2.6)^{28,29} followed by polishing the
106 assembly with the long reads using Racon (v1.4.7)³⁰ and medaka (v0.11.0)³¹ and finalizing
107 the error correction by incorporating the short reads using multiple rounds of Pilon (v1.23)³².
108 Binning is the most crucial step during metagenomic analysis. Therefore, MUFFIN combines
109 three different binning software tools, respectively CONCOCT (v1.0.0)³³, MaxBin2 (v2.2.4)
110³⁴, and MetaBAT2 (v2.14)³⁵ and refine these bins via MetaWRAP (v1.2.1)¹⁵. The user can
111 provide additional read data sets (short or long reads) to perform automatically differential
112 coverage binning to assign contigs to their bins better.

113 Moreover, an additional reassembly of bins has shown the capacity to increase the
114 completeness and N50 while decreasing the contamination of the bins¹⁵. Therefore, MUFFIN
115 allows for an optional reassembly to improve the continuity of the MAGs further. This re-
116 assembly is performed by retrieving the reads belonging to one bin and doing an assembly
117 with Unicycler (v0.4.8)³⁶.

118 To support a transparent and reproducible metagenomics workflow, all reads that cannot be
119 mapped back to the existing high-quality bins (after the refinement) are available as an
120 output for further analysis. These reads could be further analyzed by other tools or, e.g.,
121 used as a new input to run MUFFIN while providing other read sets for the differential
122 coverage binning to extract additional high-quality bins.

123 **Step 2 - Classify: Bin quality control and taxonomic assessment**

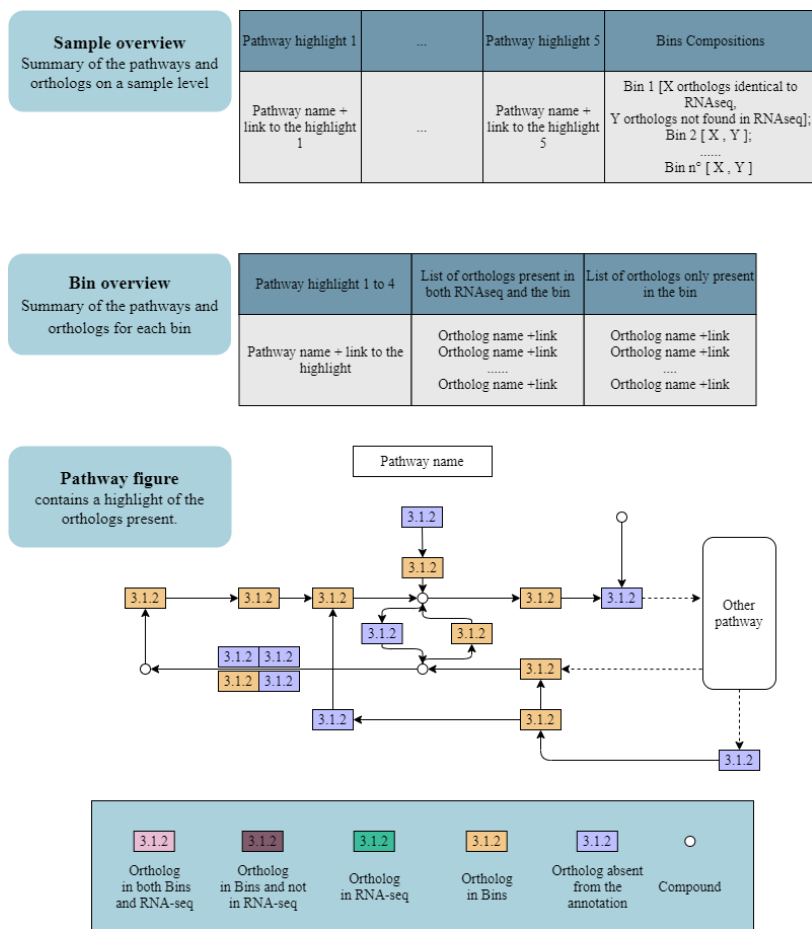
124 In the second step (**Bin quality control and taxonomic assessment**), the quality of the
125 bins is evaluated with CheckM (v1.0.18)³⁷ followed by assigning a taxonomic classification
126 to the bins using sourmash (v2.0.0a10)³⁸ and the Genome Taxonomy Database (GTDB

127 release r89)³⁹. The GTDB was chosen as it contains many unculturable bacteria and
128 archaea – this allows for monophyletic species assignments, which other databases do not
129 assure^{40,41}. GTDB substantially improved overall downstream results⁴⁰. The user can also
130 analyze other bin sets in this step regardless of their origin by providing a directory with
131 multiple FASTA files (bins).

132 Step 3 - Annotate: Bin annotation and KEGG pathway summary

133 The last step of MUFFIN (**Bin annotation and output summary**) comprises the annotation
134 of the bins using eggNOG-mapper (v2.0.1)⁴² and the eggNOG database (v5)⁴³. If RNA-Seq
135 data of the metagenome sample is provided (Illumina, paired-end), quality control using fastp
136 (v0.20.0)²³ and a *de novo* transcript assembly using Trinity (v2.8.5)⁴⁴ followed by a quasi-
137 mapping transcript quantification using Salmon (v0.15.0)⁴⁵ are performed. Lastly, the
138 transcripts are annotated using eggNOG-mapper (v2.0.1)⁴² again, followed by a parser to
139 output the activity of the pathway graphically in relation to the sample level. The expression
140 of low and high abundant genes present in the bins is shown. If only bin sets are provided
141 without any RNA-Seq data, the pathways of all the bins are created based on gene presence
142 alone. The KEGG pathway results are summarized in detail as interactive HTML files
143 (example snippet: Figure 2).

144 Like step 2, this step can be directly performed with a bin set created via another workflow.



145

146 *Figure 2: Example snippets of the sub-workflow results of step 3 (Annotate).*

147 Running MUFFIN and version control

148 MUFFIN requires only two dependencies, which allows an easy and user-friendly workflow
 149 execution. One of them is the workflow management system Nextflow ¹⁴ and the other can
 150 be either Conda ²⁰ as a package manager or Docker ²¹ to use containerized tools. A detailed
 151 Installation process is available on <https://github.com/RVanDamme/MUFFIN>. Each MUFFIN
 152 release specifies the Nextflow version it was tested on, to avoid any version conflicts
 153 between MUFFIN and Nextflow at any time. A Nextflow-specific version can always be
 154 directly downloaded as an executable file from [https://github.com/nextflow-](https://github.com/nextflow-io/nextflow/releases)
 155 [io/nextflow/releases](https://github.com/nextflow-io/nextflow/releases), which can then be paired with a compatible MUFFIN version via the -r
 156 flag.

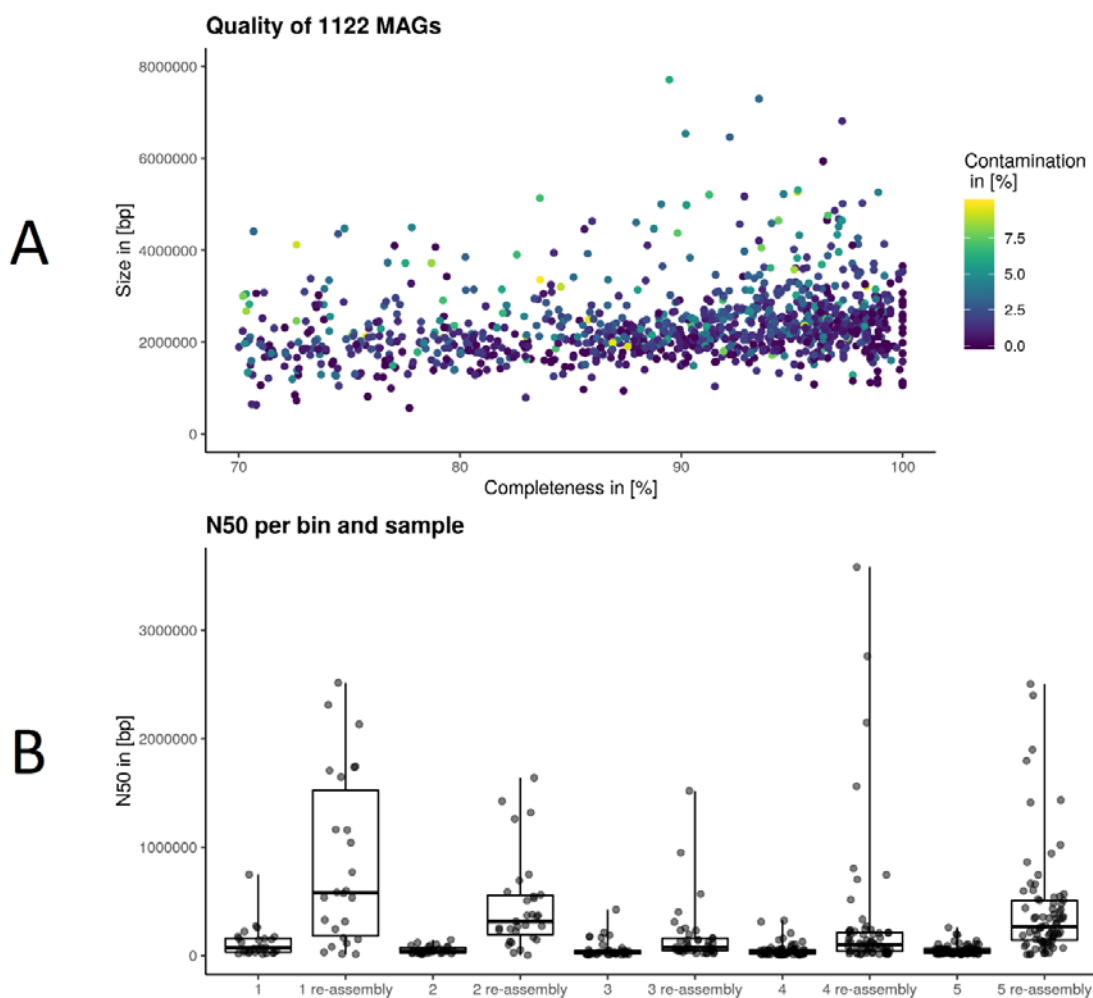
157 Results

158 We chose Nextflow for the development of our metagenomic workflow because of its direct
159 cloud computing support (Amazon AWS, Google Life Science, Kubernetes), various ready-
160 to-use batch schedulers (SGE, SLURM, LSF), state-of-the-art container support (Docker,
161 Singularity) and accessibility of a widely used software package manager (Conda).
162 Moreover, Nextflow ¹⁴ provides a practical and straightforward intermediary file handling with
163 process-specific work directories and the possibility to resume failed executions where the
164 work ceased. Additionally, the workflow code itself is separated from the 'profile' code (which
165 contains Docker, Conda, or cluster related code), which allows for a convenient and fast
166 workflow adaptation to different computing clusters without touching or changing the actual
167 workflow code.

168 The entire MUFFIN workflow was executed on 20 samples from the Bioproject PRJEB34573
169 (available at ENA or NCBI) using the Cloud Life Sciences API (google cloud) with docker
170 containers. This metagenomic bioreactor study provides paired-end Illumina and nanopore-
171 based data for each sample ⁴¹. We used five different Illumina read sets of the same project
172 for differential coverage binning, and the workflow runtime was less than two days for all
173 samples. MUFFIN was able to retrieve 1122 MAGs with genome completeness of at least 70
174 % and contamination of less than 10 % (Figure 3). In total, MUFFIN retrieved 654 MAGs with
175 genome completeness of over 90 %, of which 456 have less than 2% contamination out of
176 the 20 datasets. For comparison, a recent study was using 134 publicly available datasets
177 from different biogas reactors and retrieved 1,635 metagenome-assembled genomes with
178 genome completeness of over 50% ⁴⁶.

179 Exemplarily, we investigated the impact of additional re-assembly of each bin for five
180 samples (Figure 3). The N50 was increased by an average of 6-7 fold across all samples.
181 Twenty-six bins of the five samples had an N50 ranging between 1 to 3 Mbases. Some bins
182 benefit more of this step as the re-assembly performance depends on the number of reads
183 available for each bin.

184



186 *Figure 3: A: Quality overview of 1122 meta-assembled genomes (MAGs) by plotting size to completeness and*
187 *coloring based on contamination level. B: N50 comparison between each bin of five selected samples from the*
188 *Bioproject PRJEB34573 before and after individual bin reassembly.*

189 Discussion

190 The analysis of metagenomic sequencing data evolved as an emerging and promising
191 research field to retrieve, characterize, and analyze organisms that are difficult to cultivate.
192 There are numerous tools available for individual metagenomics analysis tasks, but they are
193 mainly developed independently and are often difficult to install and run. The MUFFIN
194 workflow gathers the different steps of a metagenomics analysis in an easy-to-install, highly
195 reproducible, and scalable workflow using Nextflow which makes them easily accessible to
196 researchers.

197

198 MUFFIN utilizes the advantages of two sequencing technologies, whereas short reads can
199 provide a better representation of low abundant species due to their higher coverage. This
200 aspect is further utilized via the final re-Assembly step after binning, which is an optional step
201 due to the additional computational burden which solely aims to improve genome continuity.
202 Another critical aspect is the full support of differential binning, for both long and short reads,
203 via a single input option. The additional coverage information from other read sets of similar
204 habitats allows for the generation of more concise bins with higher completeness and less
205 contamination because more coverage information is available for each binning tool to
206 decide which bin each contig belongs.

207 With supplied RNA-Seq data, MUFFIN is capable of enhancing the pathway results present
208 in the metagenomic sample by incorporating this data as well as the general expression level
209 of the genes. Such information is essential to further analyze a metagenomic data sets in-
210 depth, for example, to define the origin of a sample or to improve environmental parameters
211 for production reactors such as biogas reactors. Knowing whether an organism expresses a
212 gene is a crucial element in deciding whether a more detailed analysis of that organism in the
213 biotope where the sample was taken is necessary or not.

214 [Availability and future directions](#)

215 MUFFIN is an ongoing workflow project that gets further improved and adjusted. The
216 modular workflow setup of MUFFIN using Nextflow allows for fast adjustments as soon as
217 future developments in hybrid metagenomics arise, including the pre-configuration for other
218 workload managers. MUFFIN can directly benefit from the addition of new bioinformatics
219 software such as for differential expression analysis and short-read assembly that can be
220 easily plugged into the modular system of the workflow. Another improvement is the creation
221 of an advanced user and wizard user configuration file, allowing experienced users to tweak
222 all the different parameters of all the different software as desired.

223 MUFFIN will further benefit from different improvements, in particular by graphically
224 comparing the generated MAGs via a phylogenetic tree. Furthermore, a convenient approach
225 to include negative controls is under development to allow the reliable analysis of super-low
226 abundant organisms in metagenomic samples.

227 MUFFIN is publicly available at <https://github.com/RVanDamme/MUFFIN> under the GNU
228 general public license v3.0. Detailed information about the program versions used and
229 additional information can be found in the GitHub repository. All tools used by MUFFIN are
230 listed in the supplementary table S1. The Docker images used in MUFFIN are prebuilt and
231 publicly available at <https://hub.docker.com/u/nanozoo>, and the GTDB formatted for
232 sourmash(v2.0.0a10)³⁸ usage is publicly available at <https://osf.io/wxf9z/> and was created
233 by C. Titus Brown (associate professor at UC DAVIS, [http://ivory.idyll.org/blog/2019-](http://ivory.idyll.org/blog/2019-sourmash-lca-db-gtdb.html)
234 [sourmash-lca-db-gtdb.html](http://ivory.idyll.org/blog/2019-sourmash-lca-db-gtdb.html)).

235 Acknowledgment

236 We want to thank Hadrien Gourelé and Moritz Buck for the valuable insights into metagenomic
237 analysis and Annotation.

238 References

- 239 1. Handelsman, J., Rondon, M. R., Brady, S. F., Clardy, J. & Goodman, R. M. Molecular
240 biological access to the chemistry of unknown soil microbes: a new frontier for natural
241 products. *Chemistry & Biology* **5**, R245–R249 (1998).
- 242 2. De, R. Metagenomics: aid to combat antimicrobial resistance in diarrhea. *Gut Pathogens*
243 **11**, 47 (2019).
- 244 3. Mukherjee, A. & Reddy, M. S. Metatranscriptomics: an approach for retrieving novel
245 eukaryotic genes from polluted and related environments. *3 Biotech* **10**, 71 (2020).
- 246 4. Grossart, H.-P., Massana, R., McMahon, K. D. & Walsh, D. A. Linking metagenomics to
247 aquatic microbial ecology and biogeochemical cycles. *Limnology and Oceanography* **65**,
248 S2–S20 (2020).

- 249 5. Carabeo-Pérez, A., Guerra-Rivera, G., Ramos-Leal, M. & Jiménez-Hernández, J.
250 Metagenomic approaches: effective tools for monitoring the structure and functionality of
251 microbiomes in anaerobic digestion systems. *Appl Microbiol Biotechnol* **103**, 9379–9390
252 (2019).
- 253 6. Overholt, W. A. *et al.* Inclusion of Oxford Nanopore long reads improves all microbial and
254 phage metagenome-assembled genomes from a complex aquifer system. *bioRxiv*
255 2019.12.18.880807 (2019) doi:10.1101/2019.12.18.880807.
- 256 7. Beaulaurier, J. *et al.* Assembly-free single-molecule nanopore sequencing recovers
257 complete virus genomes from natural microbial communities. *bioRxiv* 619684 (2019)
258 doi:10.1101/619684.
- 259 8. Wetterstrand, K. A. DNA Sequencing Costs: Data.
260 www.genome.gov/sequencingcostsdata www.genome.gov/sequencingcostsdata.
- 261 9. Somerville, V. *et al.* Long-read based de novo assembly of low-complexity metagenome
262 samples results in finished genomes and reveals insights into strain diversity and an active
263 phage system. *BMC Microbiol* **19**, 143 (2019).
- 264 10. Warwick-Dugdale, J. *et al.* Long-read viral metagenomics captures abundant and
265 microdiverse viral populations and their niche-defining genomic islands. *PeerJ* **7**, (2019).
- 266 11. Driscoll, C. B., Otten, T. G., Brown, N. M. & Dreher, T. W. Towards long-read
267 metagenomics: complete assembly of three novel genomes from bacteria dependent on a
268 diazotrophic cyanobacterium in a freshwater lake co-culture. *Stand Genomic Sci* **12**,
269 (2017).
- 270 12. Suzuki, Y. *et al.* Long-read metagenomic exploration of extrachromosomal mobile
271 genetic elements in the human gut. *Microbiome* **7**, 119 (2019).
- 272 13. Mangul, S., Martin, L. S., Eskin, E. & Blekhman, R. Improving the usability and
273 archival stability of bioinformatics software. *Genome Biol.* **20**, 47 (2019).
- 274 14. Tommaso, P. D. *et al.* Nextflow enables reproducible computational workflows. *Nat*
275 *Biotechnol* **35**, 316–319 (2017).

- 276 15. Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP—a flexible pipeline for
277 genome-resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).
- 278 16. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics
279 data. *PeerJ* **3**, e1319 (2015).
- 280 17. Westreich, S. T., Treiber, M. L., Mills, D. A., Korf, I. & Lemay, D. G. SAMSA2: a
281 standalone metatranscriptome analysis pipeline. *BMC Bioinformatics* **19**, 175 (2018).
- 282 18. Abubucker, S. *et al.* Metabolic Reconstruction for Metagenomic Data and Its
283 Application to the Human Microbiome. *PLOS Computational Biology* **8**, e1002358 (2012).
- 284 19. Meyer, F. *et al.* The metagenomics RAST server – a public resource for the automatic
285 phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* **9**, 386 (2008).
- 286 20. Anaconda Software distribution. Anaconda | The World's Most Popular Data Science
287 Platform. <https://anaconda.com> <https://www.anaconda.com/>.
- 288 21. Boettiger, C. An introduction to Docker for reproducible research. *SIGOPS Oper.*
289 *Syst. Rev.* **49**, 71–79 (2015).
- 290 22. Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained by
291 differential coverage binning of multiple metagenomes. *Nat Biotechnol* **31**, 533–538
292 (2013).
- 293 23. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ
294 preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
- 295 24. Wick, R. *rrwick/Filtlong*. (2020).
- 296 25. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its
297 Applications to Single-Cell Sequencing. *Journal of Computational Biology* **19**, 455–477
298 (2012).
- 299 26. Antipov, D., Korobeynikov, A., McLean, J. S. & Pevzner, P. A. hybridSPAdes: an
300 algorithm for hybrid assembly of short and long reads. *Bioinformatics* **32**, 1009–1015
301 (2016).
- 302 27. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new
303 versatile metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).

- 304 28. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone
305 reads using repeat graphs. *Nat Biotechnol* **37**, 540–546 (2019).
- 306 29. Kolmogorov, M., Rayko, M., Yuan, J., Pevnikov, E. & Pevzner, P. metaFlye: scalable
307 long-read metagenome assembly using repeat graphs. *bioRxiv* 637637 (2019)
308 doi:10.1101/637637.
- 309 30. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome
310 assembly from long uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
- 311 31. *nanoporetech/medaka*. (Oxford Nanopore Technologies, 2020).
- 312 32. Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and
313 Genome Assembly Improvement.
314 <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0112963>.
- 315 33. Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat*
316 *Methods* **11**, 1144–1146 (2014).
- 317 34. Wu, Y.-W., Tang, Y.-H., Tringe, S. G., Simmons, B. A. & Singer, S. W. MaxBin: an
318 automated binning method to recover individual genomes from metagenomes using an
319 expectation-maximization algorithm. *Microbiome* **2**, 26 (2014).
- 320 35. MetaBAT, an efficient tool for accurately reconstructing single genomes from complex
321 microbial communities [PeerJ]. <https://peerj.com/articles/1165/>.
- 322 36. Unicycler: Resolving bacterial genome assemblies from short and long sequencing
323 reads. <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005595>.
- 324 37. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
325 assessing the quality of microbial genomes recovered from isolates, single cells, and
326 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 327 38. Brown, C. & Irber, L. sourmash: a library for MinHash sketching of DNA. *Journal of*
328 *Open Source Software* <https://joss.theoj.org> (2016) doi:10.21105/joss.00027.
- 329 39. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
330 substantially revises the tree of life. *Nat Biotechnol* **36**, 996–1004 (2018).

- 331 40. Méric, G., Wick, R. R., Watts, S. C., Holt, K. E. & Inouye, M. Correcting index
332 databases improves metagenomic studies. *bioRxiv* 712166 (2019) doi:10.1101/712166.
- 333 41. Brandt, C., Bongcam-Rudloff, E. & Müller, B. *Abundance tracking by long-read*
334 *nanopore sequencing of complex microbial communities in samples from 20 different*
335 *biogas/wastewater plants*. <https://www.mdpi.com/2076-3417/10/21/7518> (2020) doi:
336 10.3390/app10217518.
- 337 42. Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through Orthology
338 Assignment by eggNOG-Mapper. *Mol Biol Evol* **34**, 2115–2122 (2017).
- 339 43. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically
340 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*
341 *Res* **47**, D309–D314 (2019).
- 342 44. Haas, B. J. *et al.* De novo transcript sequence reconstruction from RNA-seq using the
343 Trinity platform for reference generation and analysis. *Nat Protoc* **8**, 1494–1512 (2013).
- 344 45. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast
345 and bias-aware quantification of transcript expression. *Nat Methods* **14**, 417–419 (2017).
- 346 46. Campanaro, S. *et al.* The anaerobic digestion microbiome: a collection of 1600
347 metagenome-assembled genomes shows high species diversity related to methane
348 production. *bioRxiv* 680553 (2019) doi:10.1101/680553.

349

350 Funding Disclosure

351 This study was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research
352 Foundation) – BR 5692/1-1 and BR 5692/1-2. This material is based upon work supported by
353 Google Cloud.

354 BM was funded by FORMAS, grant number 942-2015-1008. The funders had no role in
355 study design, data collection and analysis, decision to publish, or preparation of the
356 manuscript.

357 MH is supported by the Collaborative Research Centre AquaDiva (CRC 1076 AquaDiva) of
358 the Friedrich Schiller University Jena, funded by the DFG. MH appreciates the support of the
359 Joachim Herz Foundation by the add-on fellowship for interdisciplinary life science.

360