

1 **ATLAS (Automatic Tool for Local Assembly Structures) - a**
2 **comprehensive infrastructure for assembly, annotation, and**
3 **genomic binning of metagenomic and metatranscriptomic data**

4 Richard Allen White III^{1*#}, Joseph Brown^{1#}, Sean Colby¹, Christopher C. Overall²,
5 Joon-Yong Lee¹, Jeremy Zucker¹, Kurt R. Glaesemann³, Christer Jansson⁴, Janet K.
6 Jansson^{1*}

7

8 ¹Biological Sciences Division, Pacific Northwest National Laboratory, Richland,
9 Washington 99352, USA

10 ²Department of Neuroscience, University of Virginia, Charlottesville, VA, United States

11 ³Information technology, High Performance Computing (HPC) and Cloud services,
12 Pacific Northwest National Laboratory, Richland, Washington 99352, USA

13 ⁴Environmental and Molecular Sciences Laboratory (EMSL), Pacific Northwest
14 National Laboratory, Richland, WA 99352, USA.

15 *To whom correspondence should be addressed

16 #These authors contributed equally to this work.

17

18

19

20

21

22 **Abstract**

23 **Summary:** ATLAS (Automatic Tool for Local Assembly Structures) is a comprehensive
24 multi-omics data analysis pipeline that is massively parallel and scalable. ATLAS contains a
25 modular analysis pipeline for assembly, annotation, quantification and genome binning of
26 metagenomics and metatranscriptomics data and a framework for reference metaproteomic
27 database construction. ATLAS transforms raw sequence data into functional and taxonomic
28 data at the microbial population level and provides genome-centric resolution through
29 genome binning. ATLAS provides robust taxonomy based on majority voting of
30 protein-coding open reading frames (ORFs) rolled-up at the contig level using modified lowest
31 common ancestor (LCA) analysis. ATLAS is user-friendly, easy install through bioconda
32 maintained as open-source on GitHub, and is implemented in Snakemake for modular
33 customizable workflows.

34 **Availability and implementation:** ATLAS is written in python and distributed under a BSD
35 license. ATLAS is compatible with python 3.5+ and anaconda 3+ versions. ATLAS functions
36 on both MacOS and Linux. The source code of ATLAS is freely available at
37 <https://github.com/pnnl/atlas>.

38 **Contact:** Richard Allen White III and Janet Jansson, Earth and Biological Sciences
39 Directorate, Pacific Northwest National Laboratory, Richland, Washington 99352, USA. Email:
40 Richard.white@pnnl.gov or raw937@gmail.com, Janet.jansson@pnnl.gov

41 **Keywords:** Next-generation sequencing, Metagenomics, Metatranscriptomics, Genome
42 binning, Lowest Common Ancestor (LCA) analysis, Snakemake, Modified Majority
43 Voting-Method (MMVM)

44

45 1 Introduction

46 Whole community sequencing of DNA (i.e., metagenomics) and RNA (i.e.,
47 metatranscriptomics) has provided a wealth of information about microbial communities in a
48 variety of habitats, including community compositions, predicted functions, and metabolic
49 potential and activities (Jansson, 2011; Mason *et al.*, 2014; Prosser, 2015; Hultman *et al.*,
50 2015; Butterfield *et al.*, 2016; White III *et al.*, 2016a). Recent improvements in metagenome
51 assembly have enabled direct assembly of large and complex metagenomes (Howe *et al.*,
52 2014; Li *et al.*, 2015; White III *et al.*, 2016b). In addition, new algorithms have been developed
53 and applied for binning genomes from metagenome data (Albertsen *et al.*, 2013; Imelfort *et*
54 *al.*, 2014; Wu *et al.*, 2016). These approaches provide valuable insight into the function of
55 microbial populations that are yet to be cultivated.

56 Current sequencing technologies can reach very high throughput >1 Terabytes (TB)
57 of data in a single run (White III *et al.*, 2016a). With increasing sequencing throughput, a
58 framework for rapid, modular, customizable workflows, and integrated data analysis is needed
59 to obtain meaning from microbial community derived sequencing data. While some
60 metagenomic data analysis pipelines and frameworks exist, such as IMG (Chen *et al.*, 2017),
61 Parallel-META (Su *et al.*, 2014), MG-RAST (Meyer *et al.*, 2008), MetaAMOS (Treangen *et al.*,
62 2013), and MetaPathways2 (Konwar *et al.*, 2013), none include every key element required
63 for metagenome and metatranscriptome analysis. These key elements include quality control
64 of raw data, assembly, genomic binning, coverage estimation, functional annotation,
65 taxonomic annotation using lowest common ancestor (LCA) and quantitative analysis of
66 reads. Here we introduce ATLAS (Automatic Tool for Local Assembly Structures) as an

67 integrated and customizable pipeline for metagenome/metatranscriptome data quality control,
68 assembly and annotation, metagenome binning, coverage estimation, and expression
69 analysis.

70 **2 DESCRIPTION OF TOOL**

71 ATLAS has five analysis steps: (1) quality control, (2) assembly (3) annotation, (4)
72 genome binning, and (5) taxonomic, functional and expression quantification analyses (Figure
73 1). Default input data are Illumina paired-end reads in FASTQ format; however, single-end
74 Illumina, Ion Proton, and SOLiD reads in FASTQ format are also supported.

75 The quality control module (Step 1) involves quality filtering of the
76 metagenome/metatranscriptome sequence read data using the decontamination tool BBduk2
77 within the BBDuk tool suite (<https://sourceforge.net/projects/bbmap/>). This approach uses
78 k-mers to find and trim adapter sequences, performs quality based read trimming, and filters
79 reads based on a minimum length threshold. The reads have the option of error correction
80 based on both k-mer overlaps and read pair overlaps using Tadpole within BBMap.
81 Decontamination can be performed across any reference read set and reads will be grouped
82 into reference bins or non-hits using BBSplit. ATLAS provides references for common Illumina
83 DNA spike-ins (i.e., bacteriophage phiX) and ribosomal RNA as default contaminant
84 databases. Any additional contamination references in FASTA format are supported and be
85 user supplied. Following decontamination, quality controlled read sets are used in read
86 quantification of Step 5. A future version of ATLAS will include MerCat (i.e., “Mer-Cat enate”),
87 a *de novo* assembly free direct read analysis module plug-in (Figure 1, White III *et al.*, 2017).

88 MerCat will provide alpha diversity and feature abundance calculations from quality controlled
89 reads supplied by ATLAS using k-mer counting of any length k, specified by end user, without
90 a reference sequence database dependency (i.e., database independent property analysis -
91 DIPA) (White III *et al.*, 2017).

92 The assembly module uses quality controlled sequence reads for *de novo*
93 assembly (Step 2). Pre-assembly sub-setting uses the quality controlled reads as input then
94 uses a read coverage normalization step based on k-mer frequency. The data is then subset
95 to a target coverage using BBNorm (in BBMap tool suite
96 (<https://sourceforge.net/projects/bbmap/>). This subset of high-quality reads is then used as
97 input to ATLAS default assemblers SPAdes (i.e., metagenomic mode) (Bankevich *et al.*, 2012)
98 for datasets <100 GB and MEGAHIT (Li *et al.*, 2015) for larger more complex datasets (e.g.,
99 soil). Assembled contigs are assessed for total length and percent read coverage. The final
100 contigs can optionally be trimmed prior to determining open reading frames. Assembly output
101 defaults include quality controlled contigs >1 kbp in length, with read coverage estimations
102 >2x per contig, and with at least 40% coverage of reads across the entire contig.

103 The annotation module (Step 3) performs functional and taxonomic annotation of
104 quality control contigs. Quality controlled contigs are translated to protein coding open reading
105 frames (ORFs) using Prodigal (Hyatt *et al.*, 2012) in metagenome mode and annotated using
106 DIAMOND (Buchfink *et al.*, 2015) blastp for protein-protein searching. DIAMOND blastp
107 high-scoring pairs are filtered to user specified bitscore and e-value cut-offs (defaults >200
108 and <1x10⁷, respectively). Functional annotation utilizes non-redundant RefSeq (O'Leary *et al.*
109 *et al.*, 2016), EggNOG (Huerta-Cepas *et al.*, 2016), dbCAN for CAZy families (Yin *et al.*, 2012),

110 ENZYME for enzyme commission number (EC) (Bairoch, 2000), and COG (Tatusov *et al.*,
111 2003) databases. ATLAS obtains KEGG (Kanehisa and Goto, 2000) (i.e., KO number)
112 annotations from EggNOG reference database. ATLAS provides pre-formatted databases via
113 FTP with subcommands to simplify the process of database downloading, formatting, and
114 version tracking. The database ontologies and hierarchies are included within the annotation
115 references for downstream analysis. A DNA-DNA database search module using the Lambda
116 search tool (Hauswedell *et al.*, 2014) will be added to a future version ATLAS (Figure 1). This
117 DNA-DNA database search module will annotate ribosomal internal transcribed spacers (ITS),
118 small subunit (SSU), and large subunit (LSU) genes using Unite ITS (Abarenkov *et al.*, 2016)
119 and the Silva (SSU/LSU) (Pruesse *et al.*, 2007) databases (Figure 1).

120 For taxonomic annotation, ATLAS uses RefSeq high-scoring pairs along with
121 NCBI's taxonomy assignments reference tree via a modified majority voting-method (MMVM)
122 that utilizes lowest common ancestor (LCA) (Hanson *et al.*, 2016), to determine the lowest
123 common ancestor represented across all ORFs present within a single contig. Assembly and
124 annotation outputs from ATLAS can be directly used to create databases for proteome
125 searches or as inputs for quantitation analysis (step 5, below).

126 The binning module (Step 4) of ATLAS uses MaxBin2 (Wu *et al.*, 2016) to bin
127 genomes from metagenomes. There are two binning parameters for MaxBin2 in ATLAS; (1)
128 differential coverage estimation by user specified samples or (2) within a single sample
129 without multi-sample differential coverage mapping. For quality control of bins, we
130 recommend the CheckM package (Parks *et al.*, 2015). However, future versions of ATLAS will
131 include a bin quality control and annotation integrated into our MMVM taxonomic assignment

132 package.

133 The last module (Step 5) quantifies the coverage of the assembly by mapping
134 reads using annotations from metagenomes and metatranscriptomes. Functional and
135 taxonomic count data is obtained by mapping quality controlled reads to assembled contig
136 annotations using BMap, then parsed using featureCounts of the Subread package (Liao *et*
137 *al.*, 2014) to user specifications. This provides the final tabular output of functional
138 annotations, expressed functions (if RNA-Seq is available), taxonomy, and taxonomy based
139 functional annotations based on user specifications.

140 ATLAS is written in Python 3.5, implemented using Snakemake (Köster and
141 Rahmann, 2012) workflow infrastructure for flexible scalability, trivial parallelism of workflow
142 steps, and extensive data provenance for reproducibility. ATLAS is easily installed using
143 bioconda (<https://bioconda.github.io/>): `conda install --channel bioconda atlas`. The source
144 code of ATLAS is freely available at <https://github.com/pnnl/atlas>.

145 **3 SUMMARY**

146 ATLAS packages, databases, and workflows are easy to use, simple to install,
147 modular, and user customizable. ATLAS provides a robust bioinformatics framework for
148 metagenomic and metatranscriptomic data, where raw FASTQ files are fully processed into
149 annotated tabular files for downstream analysis and visualization. ATLAS fills a major
150 computational and analysis gap, namely the integration of quality control, assembly,
151 annotation, binning and expression analysis, and provides a framework for integrated 'omics
152 analysis.

153 Acknowledgments

154 We thank Nathan Johnson for his assistance in preparing an excellent figure.

155 Funding

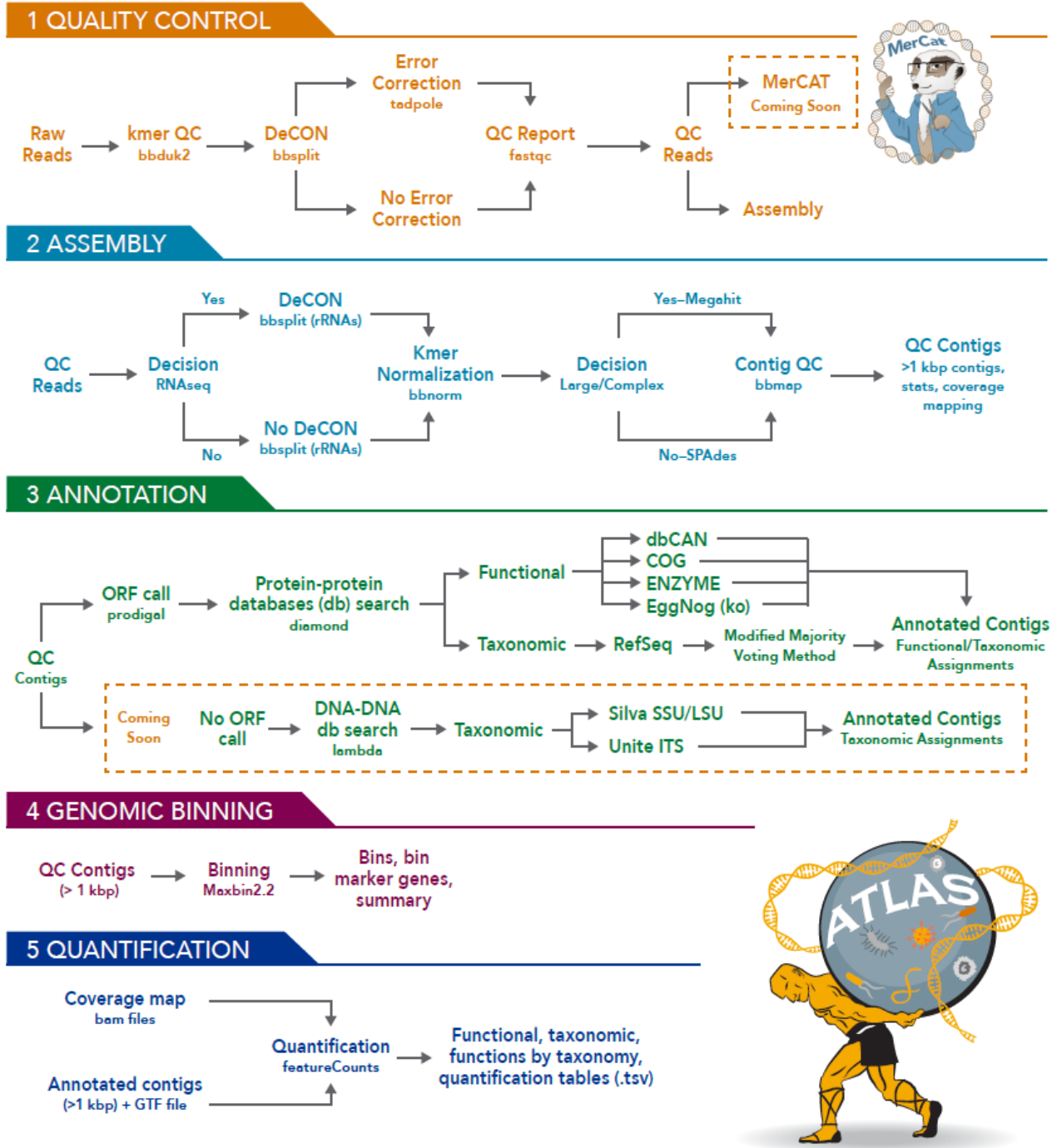
156 This research was provided by the PNNL Laboratory-Directed Research and
157 Development (LDRD) Program at PNNL through the Microbiomes in Transition (MinT)
158 Initiative and the Initiative integrated Plant-Atmosphere-Soil System (iPASS) Initiative. PNNL
159 is a national laboratory operated by Battelle for the Department of Energy (DOE) under
160 contract DE-AC06-76RL01830. A portion of the research was conducted using PNNL
161 Institutional Computing (PIC) at PNNL and at EMSL, a national scientific user facility
162 sponsored by the DOE Office of Biological and Environmental Research and located at
163 PNNL.

164 References

- 165 Abarenkov,K. et al. (2016) Annotating public fungal ITS sequences from the built environment
166 according to the MlxS-Built Environment standard – a report from a May 23-24, 2016
167 workshop (Gothenburg, Sweden). *MycoKeys*, **16**, 1-15.
- 168 Albertsen,M. *et al.* (2013) Genome sequences of rare, uncultured bacteria obtained by
169 differential coverage binning of multiple metagenomes. *Nat. Biotechnol.*, **6**, 533-538.
- 170 Bankevich, A. *et al.* (2012) SPAdes: a new genome assembly algorithm and its applications to
171 single-cell sequencing. *J. Comput. Biol.*, **19**, 455-477.
- 172 Bairoch,A. (2000) The ENZYME database in 2000. *Nucleic Acids Res.*, **28**, 304-305.
- 173 Buchfink,B. *et al.* (2015) Fast and sensitive protein alignment using DIAMOND. *Nat.*
174 *Methods.*, **12**, 59-60.
- 175 Butterfield,C.N. *et al.* (2016) Proteogenomic analyses indicate bacterial methylotrophy and
176 archaeal heterotrophy are prevalent below the grass root zone. *PeerJ*, **4**, e2687.
- 177 Chen,I.A. *et al.* (2017) IMG/M: integrated genome and metagenome comparative data
178 analysis system. *Nucleic Acids Res.*, **45**, D507-D516.

- 179 Kanehisa,M. and Goto,S. (2000) KEGG: Kyoto Encyclopedia of Genes and Genomes.
180 *Nucleic Acids Res.*, **28**, 27-30.
- 181 Hanson,N.W. *et al.* (2016) LCA*: an entropy-based measure for taxonomic assignment within
182 assembled metagenomes. *Bioinformatics*, **32**, 3535-3542.
- 183 Howe,A.C. *et al.* (2014) Tackling soil diversity with the assembly of large, complex
184 metagenomes. *Proc. Natl. Acad. Sci.*, **111**, 4904-4909.
- 185 Hauswedell,H. *et al.* (2014) Lambda: the local aligner for massive biological data.
186 *Bioinformatics*, **30**, 49-355.
- 187 Huerta-Cepas,J. *et al.* (2016) eggNOG 4.5: a hierarchical orthology framework with
188 improved functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic
189 Acids Res.*, **44**, D286-D293.
- 190 Hyatt,D. *et al.* (2012) Gene and translation initiation site prediction in metagenomic
191 sequences. *Bioinformatics*, **28**, 2223-2230.
- 192 Imelfort,M. *et al.* (2014) GroopM: an automated tool for the recovery of population genomes
193 from related metagenomes. *PeerJ*, **2**, e603.
- 194 Jansson,J.K. (2011) Towards "Tera-Terra": terabase sequencing of terrestrial metagenomes.
195 *ASM Microbe*, **6**, 309-315.
- 196 Konwar,K.M. *et al.* (2013) MetaPathways: a modular pipeline for constructing
197 pathway/genome databases from environmental sequence information. *BMC Bioinformatics*,
198 **14**, 202.
- 199 Köster,J. and Rahmann, S. (2012) Snakemake--a scalable bioinformatics workflow engine.
200 *Bioinformatics*, **28**, 2520-2522.
- 201 Li,D. *et al.* (2015) MEGAHIT: an ultra-fast single-node solution for large and complex
202 metagenomics assembly via succinct *de Bruijn* graph. *Bioinformatics*, **31**, 1674-1676.
203
- 204 Liao,Y. *et al.* (2014) featureCounts: an efficient general purpose program for assigning
205 sequence reads to genomic features. *Bioinformatics*, **30**, 923-930.
206
- 207 Mason,O.U. *et al.* (2012) Metagenome, metatranscriptome and single-cell sequencing reveal
208 microbial response to Deepwater Horizon oil spill. *ISMEJ*. **6**, 1715-1727.
209
- 210 Meyer,F. *et al.* (2008) The metagenomics RAST server - a public resource for the automatic
211 phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics*, **9**, 386.
212
- 213 O'Leary,N.A. *et al.* (2016) Reference sequence (RefSeq) database at NCBI: current status,

- 214 taxonomic expansion, and functional annotation. *Nucleic Acids Res.*, **44**, D733-D745.
- 215 Parks,D.H. *et al.* (2015) CheckM: assessing the quality of microbial genomes recovered from
216 isolates, single cells, and metagenomes. *Genome Res.*, **25**, 1043-1055.
- 217 Prosser,J.I. (2015) Dispersing misconceptions and identifying opportunities for the use of
218 “omics” in soil microbial ecology. *Nat. Rev. Microbiol.*, **13**, 439-446.
- 219 Pruesse,E. (2007) SILVA: a comprehensive online resource for quality checked and aligned
220 ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res.*, **35**, 7188-7196.
- 221 Su,X. *et al.* (2014) Parallel-META 2.0: enhanced metagenomic data analysis with functional
222 annotation, high performance computing and advanced visualization. *PLoS One*, **9**, e89323
- 223 Tatusov,R.L. *et al.* (2003) The COG database: an updated version includes eukaryotes. *BMC*
224 *Bioinformatics*, **4**, 41.
- 225 Treangen,T.J. *et al.* (2013) MetAMOS: a modular and open source metagenomic assembly
226 and analysis pipeline. *Genome Biol.*, **14**, R2.
- 227 White,R.A. III, *et al.* (2016a) Moleculo Long-Read Sequencing Facilitates Assembly and
228 Genomic Binning from Complex Soil Metagenomes. *Msystems*, **1**, e00045-16.
- 229 White,R.A. III, *et al.* (2016b) The past, present and future of microbiome analyses. *Nature*
230 *Protocols*, **11**, 2049-2053.
- 231 White,R.A. III, *et al.* (2017) MerCat: a versatile k-mer counter and diversity estimator for
232 database-independent property analysis obtained from metagenomic and/or
233 metatranscriptomic sequencing data. *PeerJ Preprints*, **5**, e2825v1
- 234 Wu,Y.W. *et al.* (2016) MaxBin 2.0: an automated binning algorithm to recover genomes from
235 multiple metagenomic datasets. *Bioinformatics*, **32**, 605–607.
- 236 Yin,Y. *et al.* (2012) dbCAN: a web resource for automated carbohydrate-active enzyme
237 annotation. *Nucleic Acids Res.* **40**, 445-451.
- 238
- 239
- 240
- 241
- 242



244 **Figure 1: ATLAS workflow.** MerCat and Lambda based (DNA-DNA database) search
 245 modules will be added to a future version of ATLAS.