

# **MerCat: a versatile k-mer counter and diversity estimator for database-independent property analysis obtained from metagenomic and/or metatranscriptomic sequencing data**

Richard Allen White III<sup>1\*</sup>, Ajay Panyala<sup>2</sup>, Kevin Glass<sup>3</sup>, Sean Colby<sup>1</sup>, Kurt R Glaesemann<sup>2</sup>, Christer Jansson<sup>3</sup>, Janet K Jansson<sup>1\*</sup>

<sup>1</sup>Biological Sciences Division, Pacific Northwest National Laboratory, Richland, Washington 99352, USA

<sup>2</sup>Information technology, High Performance Computing (HPC) and Cloud Services, Pacific Northwest National Laboratory, Richland, Washington 99352, USA

<sup>3</sup>Environmental and Molecular Sciences Laboratory, Pacific Northwest National Laboratory, Richland, WA 99352, USA

\*To whom correspondence should be addressed

## 22 Abstract

23 **Summary:** MerCat ("Mer - Catenate") is a parallel, highly scalable and modular property  
24 software package for robust analysis of features in next-generation sequencing data. Using  
25 assembled contigs and raw sequence reads from any platform as input, MerCat performs  
26 k-mer counting of any length k, resulting in feature abundance counts tables. MerCat allows  
27 for direct analysis of data properties without reference sequence database dependency  
28 commonly used by search tools such as BLAST for compositional analysis of whole  
29 community shotgun sequencing (e.g., metagenomes and metatranscriptomes).

30 **Availability and implementation:** MerCat is written in Python and distributed under a BSD  
31 license. The source code of MerCat is freely available at <https://github.com/pnnl/mercat>  
32 MerCat is compatible with Python 2 and 3 and works on both Mac OS X and Linux. MerCat  
33 can also be easily installed using bioconda: `conda install mercat`

34 **Contact:** Richard Allen White III and Janet Jansson, Biological Sciences Division, Pacific  
35 Northwest National Laboratory, Richland, Washington 99352, USA. Email:  
36 [Richard.white@pnnl.gov](mailto:Richard.white@pnnl.gov) or [raw937@gmail.com](mailto:raw937@gmail.com), [Janet.jansson@pnnl.gov](mailto:Janet.jansson@pnnl.gov)

37 **Keywords:**  
38 K-mer counting  
39 Database-independent property analysis (DIPA)  
40 Metagenomic analysis  
41 Metatranscriptomic analysis  
42 Diversity-estimation

43

44

45

46

47

# 48 1 Introduction

49 Whole community sequencing of total DNA (i.e., metagenomics) and total RNA (i.e.,  
50 metatranscriptomics) have provided windows into the composition, functions and potential  
51 roles of microbial communities residing in complex ecosystems (e.g., soil) (White III *et al.*,  
52 2016a). The throughput of next generation sequencing (NGS) technologies is continuously  
53 increasing: sequence data currently requires terabytes of storage (White III *et al.*, 2016a) and  
54 read lengths can exceed 90 kbp (Laver *et al.*, 2015). Therefore, developments of robust  
55 bioinformatics tools are needed to analyze these data.

56 Reference sequence databases and tools that classify sequences are critical  
57 bottlenecks in metagenomics and metatranscriptomics. For example, tools that search  
58 reference sequence databases against query data such as homology-based BLAST are  
59 computationally slow against large databases (e.g., KEGG) (Silva *et al.*, 2016). In addition,  
60 metagenome assembly approaches, although recently improved for complex data types  
61 (Howe *et al.*, 2014; Li *et al.*, 2015; White III *et al.*, 2016b), are not able to assemble all data.  
62 Open-source reference sequence databases are facing a number of challenges, including  
63 finding lasting funding, many are moving to a subscription-based access (e.g., KEGG  
64 [www.kegg.jp/kegg/](http://www.kegg.jp/kegg/)), slowed development (e.g., COG <https://www.ncbi.nlm.nih.gov/COG/>), or  
65 discontinuation (e.g., CAMERA <http://camera.calit2.net/>).

66 Database-independent property analysis (i.e., DIPA) which utilizes counting of  
67 k-mers subsequences (of length k) from sequence reads obtained from NGS platforms  
68 without a reference sequence database for matching query data. DIPA-based k-mer counting  
69 provides rapid and robust microbial community analysis and characterization without the

biases or limitations of sequence databases (Jiang *et al.*, 2012) and/or *de novo* assembly in order to compare and contrast sequence datasets. K-mers are critical to assembly (Li *et al.*, 2015), counting (Zhang *et al.*, 2014), partitioning (Howe *et al.*, 2014), genomic binning (Wu *et al.*, 2015) and classification (Jiang *et al.*, 2012). K-mer based counting is amongst the fastest approaches for profiling metagenomic and/or metatranscriptomic data (Lindgreen *et al.*, 2015).

There are many k-mer counters (Zhang *et al.*, 2014), and even database dependent k-mer profilers (Koslicki and Falush, 2016). MerCat provides only k-mer counting tool for assembled contigs (.fna), translated protein-coding ORFs (.faa) and NGS reads (.fastq) for any size k-mer. Alpha diversity metrics for microbial ecology including chao1, ace, simpson, goods coverage, dominance and fishers alpha are generated by MerCat. Nucleotide properties (e.g., %G+C, %A+T) and protein properties of translated protein-coding open reading frames (ORFs) (e.g., protein isoelectric point, pI, and hydrophobicity metrics) are also generated.

Here we describe MerCat, a tool that can accommodate any size sequence file by utilizing a 'divide and conquer' approach and then performs k-mer analysis. MerCat can be employed for rapid, robust, versatile analysis of NGS microbial community data using DIPA.

## 2 DESCRIPTION OF THE TOOL

MerCat is a modular and highly-scalable Python-based open-source software package. MerCat computes k-mer frequency counting to any length k on assembled contigs as nucleotide fasta, raw reads (e.g., fastq), and translated protein-coding ORFs (e.g., protein

91 fasta). The package also allows for user-defined custom analyses. Although raw read inputs  
 92 can be used in MerCat, it is not recommended due to low quality and sequencing errors, thus  
 93 we utilize Trimmomatic (Bolger *et al.*, 2014) for quality control trimming of low-quality data  
 94 obtained by fastq formats (default trimming is base pair quality score  $>Q_{30}$ ). K-mer counting in  
 95 MerCat has two modes: DNA mode, which can analyze nucleotide contigs directly, and  
 96 Protein mode, wherein nucleotide contigs are translated into protein-coding ORFs with  
 97 Prodigal (Hyatt *et al.*, 2012), using the metagenomic option (default) (Figure 1). Individual  
 98 sequence files or many files within a single folder can be analyzed by MerCat. Tabular outputs  
 99 include overall feature files (e.g., many files within a folder), or per-file feature analysis based  
 100 on k-mer frequency counts tables for either nucleotides and/or proteins fasta files. Tabular file  
 101 outputs are stored as comma-separated files for downstream analysis. MerCat can also  
 102 calculate Alpha diversity metrics for each file in both Protein and DNA mode. As a default, we  
 103 provide k-mer frequency stacked bar plots for individual samples and MDS plots for many  
 104 samples.

105           MerCat can handle input files >10 gigabyte by splitting them into multiple files.  
 106 MerCat computes on the individual files, then combines the resulting data, analyzes data and  
 107 produces the final output (as mentioned previously) for large input files. The combined overall  
 108 output generated may be too large to fit in the available memory of a standard computer. For  
 109 this reason, we used Dask, a Python-based parallel-computing library that enables  
 110 processing data that does not fit into available memory. Dask stores the data on a hard disk,  
 111 then loads portions of it back and forth into memory as needed for analysis. This enables  
 112 MerCat to scale from laptop to high-performance computing resources, all within the same

113 user friendly-package.

### 114 **3 SUMMARY**

115 MerCat provides DIPA for metagenomic and metatranscriptomic data, starting from  
116 nucleotide and protein sequence files and ending with tabular files for downstream analysis  
117 and visualization. MerCat is scalable, accommodating for large input files, is user-friendly,  
118 easy to install and is user customizable. MerCat fills a major computational bottleneck by  
119 enabling rapid analysis of many datasets and large datasets in a database independent  
120 manner.

### 121 **Acknowledgments**

122 We thank Nathan Johnson for his assistance in preparing excellent figures.

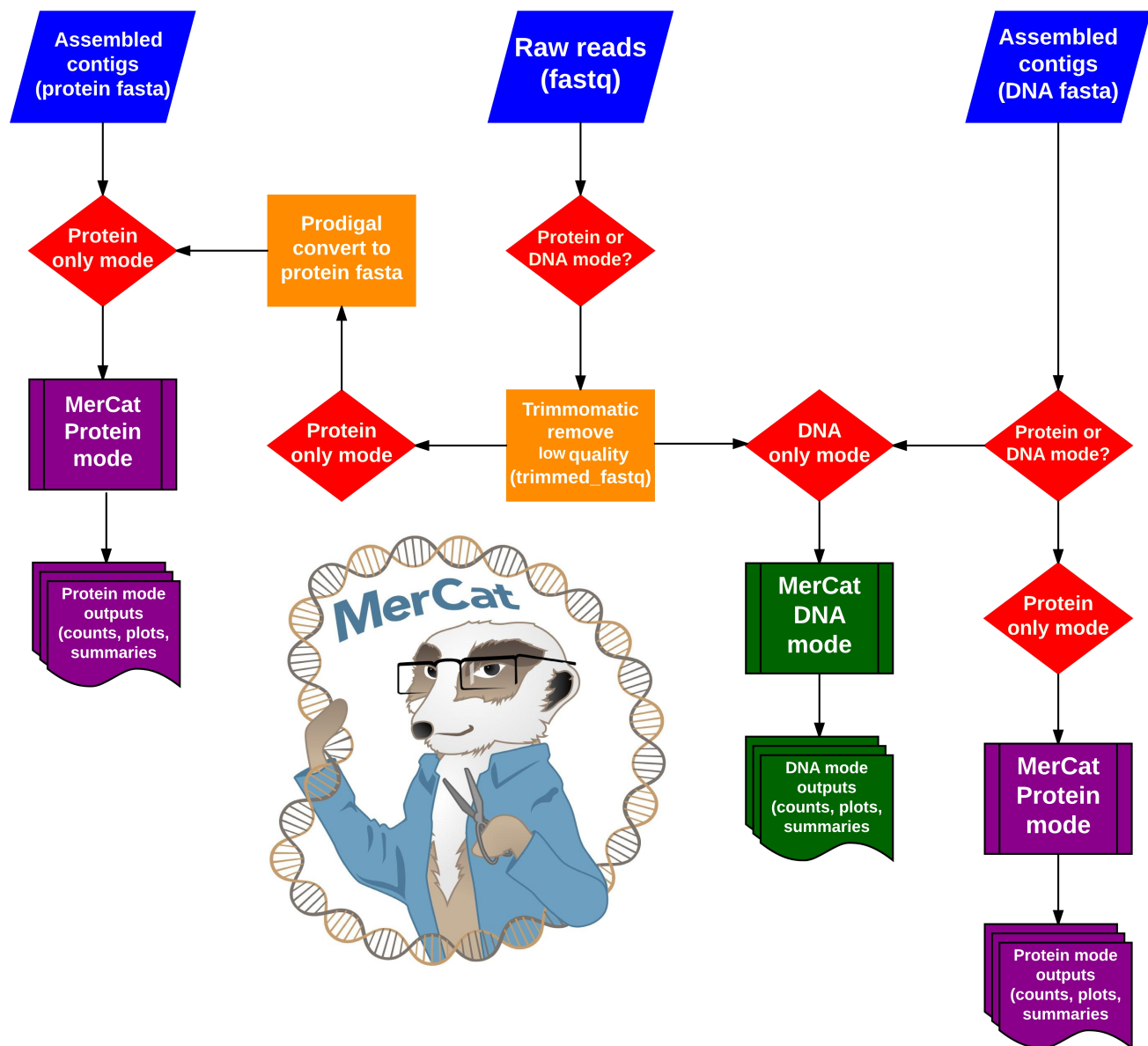
### 123 **Funding**

124 This research was provided by PNNL Laboratory-Directed Research and  
125 Development (LDRD) Program at PNNL; Microbiomes in Transition (MinT) Initiative, and  
126 PNNL Initiative integrated Plant-Atmosphere-Soil System (iPASS; PNNL Project # 204412), a  
127 multiprogram national laboratory operated by Battelle for the Department of Energy (DOE)  
128 under contract DE-AC06-76RL01830. A portion of the research was conducted using PNNL  
129 Institutional Computing (PIC) at PNNL and at EMSL, a national scientific user facility  
130 sponsored by the DOE Office of Biological and Environmental Research and located at  
131 PNNL.

132

# References

- 134 Bolger, A.M. *et al.* (2014) Trimmomatic: a flexible trimmer for Illumina sequence data.  
135 *Bioinformatics*, **30**, 2114-2120.
- 136 Hyatt, D. *et al.* (2012) Gene and translation initiation site prediction in metagenomic  
137 sequences. *Bioinformatics*, **28**, 2223-2230.
- 138 Howe, A.C. *et al.* (2014) Tackling soil diversity with the assembly of large, complex  
139 metagenomes. *Proc Natl Acad Sci USA.*, **111**, 4904-4909.
- 140 Jiang, B. *et al.* (2012) Comparison of metagenomic samples using sequence signatures.  
141 *BMC Genomics*, **13**, 730. doi: 10.1186/1471-2164-13-730.
- 142 Koslicki, D. and Falush, D. (2016) MetaPalette: a k-mer Painting Approach for Metagenomic  
143 Taxonomic Profiling and Quantification of Novel Strain Variation. *Msystems*, **3**, pii: e00020-16.
- 144 Laver, T. *et al.* (2015) Assessing the performance of the Oxford Nanopore Technologies  
145 MinION. *Biomol. Detect Quantif.*, **3**, 1-8.
- 146
- 147 Lindgreen, S. *et al.* (2016) An evaluation of the accuracy and speed of metagenome analysis  
148 tools. *Sci. Rep.*, **6**, 19233.
- 149
- 150 Li, D. *et al.* (2015) MEGAHIT: an ultra-fast single-node solution for large and complex  
151 metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**, 1674-1676.
- 152
- 153 Silva, G.G. *et al.* (2016) SUPER-FOCUS: a tool for agile functional analysis of shotgun  
154 metagenomic data. *Bioinformatics*. **32**, 354-361.
- 155
- 156 White, R.A. III, *et al.* (2016) (a) The past, present and future of microbiome analyses. *Nature*  
157 *Protocols* **11**, 2049-2053.
- 158 White, R.A. III, *et al.* (2016) (b) Molecule long-read sequencing facilitates assembly and  
159 genomic binning from complex soil metagenomes. *Msystems*, **3**, e00045-16.  
160 doi:10.1128/mSystems.00045-16.
- 161 Wu, Y.W. *et al.* (2016) MaxBin 2.0: an automated binning algorithm to recover genomes from  
162 multiple metagenomic datasets. *Bioinformatics* **32**, 605-607.
- 163 Zhang, Q. *et al.* (2014) These Are Not the K-mers You Are Looking For: Efficient Online K-mer  
164 Counting Using a Probabilistic Data Structure. *PLoS ONE* **9**, e101271.
- 165
- 166



167 **Figure 1: MerCat workflows.** Inputs include assembled contigs (.fna), assembled contigs  
 168 previously translated protein-coding ORFs (.faa) and NGS reads (.fastq) for any size k-mer.  
 169 Outputs include tabular count tables for individual mers, stacked bar and MDS plots, alpha  
 170 diversity statistics. Prodigal uses metagenomic mode as default for translating assembled  
 171 nucleotide contigs into protein-coding ORFs (.faa). Trimmomatic as default requires base  
 172 quality  $>Q_{30}$ .