

Hundo: a Snakemake workflow for microbial community sequence data

Joseph Brown¹, Nicole Zavoshy¹, Colin J Brislawn¹, and Lee Ann McCue¹

¹Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland, WA USA

Corresponding author:

Joseph Brown¹

Email address: joseph.brown@pnnl.gov

ABSTRACT

Hundo is a software package that performs quality control and annotation of ribosomal RNA and internal transcribed spacer sequence reads. It provides a comprehensive suite of tools and options that aim to reduce the effort of performing robust sequence annotation while obtaining the maximum amount of data provenance to ensure replicability. The software package performing annotation is implemented in Python, the analytical workflow is implemented in Snakemake, and dependencies are installed via Bioconda. Extensive documentation and the full source code are available under the MIT license at: <https://github.com/pnnl/hundo>.

INTRODUCTION

Many packages exist for processing 16S ribosomal RNA (rRNA) and internal transcribed spacer (ITS) amplicon sequence reads, each with nuances pertaining to quality control, annotation, output, and importantly, usage and install requirements. Existing packages like QIIME (Caporaso et al., 2010), QIIME 2 (QIIME 2 Development Team, 2018), and mothur (Schloss et al., 2009) offer vast toolsets to carry out analyses, often resulting in over-complication and workflow fragility arising from incompatible file formats required at various stages of complex and ad hoc, user-specified workflows.

When designing a workflow, one of the primary considerations is how to account for user choice and how much flexibility to allow. Existing annotation and analysis packages like QIIME 2 account for user preferences by providing both pipelines and individual methods within those pipelines. In addition, these tools, including QIIME 2 and mothur, provide many additional tools for downstream analysis after sequences have been annotated. By accounting for all use cases and by adding additional analytical and visualization tools, the commands and subcommands for each of these packages increase in number and complexity. For instance, QIIME lists 150 individual commands, QIIME 2 has 20 sub-commands and 119 sub-sub-commands, and mothur lists 146 commands.

The documentation for QIIME2 and mothur are both exceptionally detailed; include comprehensive installation instructions; offer multiple tutorials, each walking the reader through a slew of subcommands where the inputs for the current task are dependent upon the outputs of the previous; and provide example datasets for users to test workflows on their own installation. While detailed documentation is useful for experts, it can be daunting for first-time users as they may not know which settings are important and how altering parameters may affect the outcome.

We introduce Hundo, a 16S, 18S, and ITS annotation workflow. It incorporates a minimal set of tools to carry out quality control and annotation while adding important aspects of environment and data provenance that existing software packages lack. By reducing the toolset and streamlining processing from raw sequence data to annotated (representative) sequences down to a single command, we lessen the burden of understanding how and where to get started, what was performed to the sequence reads, and what outputs were generated. User flexibility is offered through command line options that affect individual commands or change commands within the workflow, like swapping VSEARCH (Rognes et al., 2016) for alignment rather than BLAST (Altschul et al., 1997). We provide documentation, example data,

46 outputs, and a sample report, to facilitate rapid evaluation and adoption of our software. By restricting
47 use cases to quality control and annotation, and by taking advantage of cloud services and a workflow
48 manager, Hundo represents a well-documented, lightweight software package capable of performing
49 robust rRNA sequence annotation.

50 APPROACH

51 Hundo is a combination of a command line application written in Python 3 and Snakemake (Köster and
52 Rahmann, 2012). The command line interface (CLI) works to parse options and format the sub-command
53 to Snakemake and only requires specifying a directory containing paired-end FASTQ files.

54 Hundo provides environment definitions which are used by Snakemake to install tested versions
55 of binaries from Bioconda (Grüning et al., 2018) before starting to execute sequence processing tasks
56 defined in the workflow. In addition to handling prerequisite software installs, Hundo's version controlled
57 reference annotations (Brown, 2017) are downloaded during execution from Zenodo (Zenodo, 2018).

58 All output files generated from the workflow have associated provenance data including attributes like
59 the command used in their generation, the input file, and the date and time of execution. Additionally,
60 Hundo builds upon this metadata by compiling a comprehensive HTML report that includes quality control
61 statistics, interactive plots with taxonomic breakdowns across samples, the workflow's configuration, the
62 environment definitions, and explanations of output file contents. Providing comprehensive provenance
63 increases replicability, which is why both Hundo and QIIME2 include it by default.

64 RESULTS

65 Hundo is comprised of custom methods as well as wrappers for existing software. Sequence processing is
66 initiated by calling `hundo annotate`. Command line options are parsed, validated, and sent to the
67 accompanying workflow. Paired-end sequence reads are trimmed by quality and are optionally filtered of
68 sequence content like PhiX, a known contaminant among bacterial assemblies (Mukherjee et al., 2015),
69 using BBduk of the BBMap toolset (Bushnell, 2018). Passing reads are merged using VSEARCH then
70 aggregated into a single FASTA file with headers describing the origin and count of the sequences. Merged
71 sequence reads are evaluated by their expected error rates and filtered a final time before clustering into
72 operational taxonomic units (OTU).

73 Within Hundo, clusters are created from de-replicated merged sequences that have passed quality
74 control using VSEARCH. Sequences are pre-clustered into centroids using VSEARCH to accelerate
75 chimera filtering. Chimera filtering is completed in two steps: *de novo* and then reference-based.
76 Following chimera filtering, sequences are placed into clusters using distance-based, greedy clustering
77 with VSEARCH based on the allowable percent difference of the configuration.

78 After OTU sequences have been determined, BLAST or VSEARCH is used to align sequences to
79 the reference database. Sequence alignments are filtered relative to the best hit before the least common
80 ancestor (LCA) is calculated. Reference databases for 16S were curated by the CREST team and
81 Hundo incorporates the CREST LCA method (Lanzén et al., 2012) for 16S. ITS reference databases and
82 taxonomies are maintained by UNITE (Köljalg et al., 2013).

83 For downstream processing, counts are assigned to OTUs using the global alignment method of
84 VSEARCH, which outputs the final feature-abundance table as a tab-delimited text file. The Biom
85 tool (McDonald et al., 2012) is used to convert the tab-delimited table to biom. Multiple alignment of
86 sequences is completed using MAFFT (Nakamura et al., 2018) in order to generate a tree based on the
87 OTU sequences using FastTree2 (Price et al., 2010).

88 Hundo finalizes processing by generating a comprehensive HTML report (Figure 1) with summary
89 tables, interactive plots, explanation of the workflow, details of output files, and an embedded archive
90 containing the most relevant output files. The summary table includes counts per sample input, quality
91 controlled count, merged, and finally the count assigned to OTUs. Shannon, Simpson, and inverse
92 Simpson indexes are included alongside the counts. Sequence quality plots report qualities per base
93 per sequence before and after filtering. Observed taxonomies per sample are plotted by absolute count
94 and relative abundance across phylum, class, and order to facilitate a cursory analysis of taxonomic
95 composition across an experiment. To ease replication, the report includes the user's configuration and
96 environment setup with absolute versions of applications used in their protocol. The report concludes by

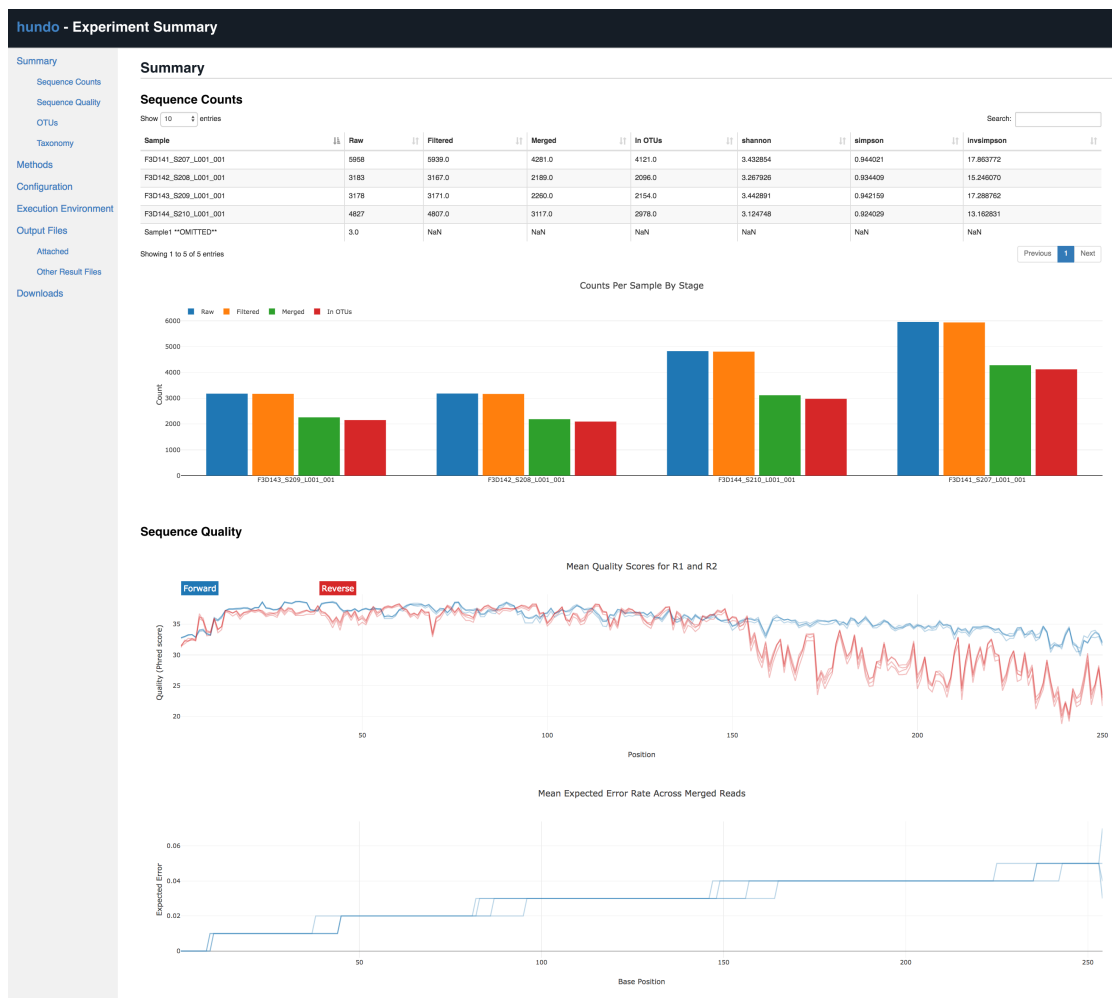


Figure 1. Example HTML report generated by Hundo. An interactive version is available at: <https://pnnl.github.io/hundo/>.

97 defining output files and providing a link to download a subset of output files that are most relevant for use
 98 in dedicated analysis packages like VEGAN (Dixon, 2003) or phyloseq (McMurdie and Holmes, 2013).

99 DISCUSSION

100 Using 16S rRNA sequence data from Dohnalkova et al. (Dohnalkova et al., 2017), we aimed to reproduce
 101 their findings among rhizospheric microbial communities of *Pinus resinosa* grown in controlled, laboratory
 102 environments, using Hundo and Phyloseq. After observing complex extracellular soil organic matter
 103 forming associations with minerals, Dohnalkova et al. performed 16S amplicon sequencing on 30 samples
 104 to identify microbes capable of contributing to the mineral weathering process and those capable of
 105 nitrogen fixation.

106 Hundo was installed and executed as `hundo annotate` across paired-end sequence data with
 107 filtering for Illumina adapters and PhiX contamination, SILVA rRNA reference database for taxonomic
 108 annotation and chimera removal, and both `--jobs` and `--threads` set to 24. Hundo generated the
 109 feature abundance table and phylogenetic tree necessary for downstream processing in 12 minutes with
 110 progress messages, a clear message of a successful run, and no additional user input. The feature
 111 abundance table and taxonomic tree were imported into Phyloseq to plot abundances and their phylogeny.

112 Our analysis of the microbial communities replicates and confirms taxonomic composition of microbial
 113 communities among the soil rhizosphere includes high percentages of *Burkholderiales* and *Rhizobiales* as
 114 well as unidentified microbes phylogenetically similar to *Leptospiraceae*. Relative to the published work,

115 our replication study shows that among OTUs classified as *Rhizobiales*, we see an increase in genus level
116 assignments which could allow finer interpretation among these community members. These data, our
117 work, and recreated figures are publicly available in a version controlled code and data archive (Brislawn,
118 2018).

119 CONCLUSION

120 Hundo efficiently processes raw paired-end rRNA and ITS sequence data from a single command into
121 standardized outputs. By tailoring its functionality to a well-defined workflow and by using Snakemake to
122 handle dependency installs and track provenance, Hundo serves as a tool for replicable 16S, 18S, and ITS
123 sequence processing.

124 FUNDING

125 This research was supported by the Microbiomes in Transition Initiative LDRD Program at the Pacific
126 Northwest National Laboratory, a multi-program national laboratory operated by Battelle for the DOE
127 under Contract DE-AC06-76RL01830.

128 REFERENCES

- 129 Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D. J. (1997).
130 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic acids*
131 *research*, 25(17):3389–3402.
- 132 Brislawn, C. J. (2018). Minimal replication of Dohnalkova 2017 using Hundo. *Open Science Framework*.
133 DOI: 10.17605/OSF.IO/T4NC7.
- 134 Brown, J. (2017). 16S and ITS Reference Annotation Databases [Data set]. *Zenodo*. DOI: 10.5281/zen-
135 odo.1043977.
- 136 Bushnell, B. (2018). BBMap. <https://sourceforge.net/projects/bbmap/>, last accessed on 09/08/2018.
- 137 Caporaso, J Gregoand Huttley, G. A., Kelley, S. T., Knights, D., Koenig, J. E., Ley, R. E., Lozupone,
138 C. A., McDonald, D., Muegge, B. D., Pirrung, M., Reeder, J., Sevinsky, J. R., Turnbaugh, P. J., Walters,
139 W. A., Widmann, J., Yatsunencko, T., Zaneveld, J., and Knight, R. (2010). QIIME allows analysis of
140 high-throughput community sequencing data. *Nature methods*, 7(5):335–336.
- 141 Dixon, P. (2003). VEGAN, a package of R functions for community ecology. *Journal of Vegetation*
142 *Science*, 14(6):927–930.
- 143 Dohnalkova, A. C., Tfaily, M. M., Smith, A. P., Chu, R. K., Crump, A. R., Brislawn, C. J., Varga, T.,
144 Shi, Z., Thomashow, L. S., Harsh, J. B., et al. (2017). Molecular and microscopic insights into the
145 formation of soil organic matter in a red pine rhizosphere. *Soils*, 1(1):4.
- 146 Grüning, B., Dale, R., Sjödin, A., Chapman, B. A., Rowe, J., Tomkins-Tinch, C. H., Valieris, R., Köster,
147 J., and Bioconda Team (2018). Bioconda: sustainable and comprehensive software distribution for the
148 life sciences. *Nature methods*, 15(7):475–476.
- 149 Kõljalg, U., Nilsson, R. H., Abarenkov, K., Tedersoo, L., Taylor, A. F. S., Bahram, M., Bates, S. T., Bruns,
150 T. D., Palme, J. B., Callaghan, T. M., Douglas, B., Drenkhan, T., Eberhardt, U., Dueñas, M., Grebenc,
151 T., Griffith, G. W., Hartmann, M., Kirk, P. M., Kohout, P., Larsson, E., Lindahl, B. D., Lücking, R.,
152 Martín, M. P., Matheny, P. B., Nguyen, N. H., Niskanen, T., Oja, J., Peay, K. G., Peintner, U., Peterson,
153 M., Põldmaa, K., Saag, L., Saar, I., Schübler, A., Scott, J. A., Senés, C., Smith, M. E., Suija, A.,
154 Taylor, D. L., Telleria, M. T., Weiss, M., and Larsson, K. H. (2013). Towards a unified paradigm for
155 sequence-based identification of fungi. *Molecular Ecology*, 22(21):5271–5277.
- 156 Köster, J. and Rahmann, S. (2012). Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics*
157 (*Oxford, England*), 28(19):2520–2522.
- 158 Lanzén, A., Jørgensen, S. L., Huson, D. H., Gorfer, M., Grindhaug, S. H., Jonassen, I., Øvreås, L.,
159 and Urich, T. (2012). CREST—classification resources for environmental sequence tags. *PLoS one*,
160 7(11):e49334.
- 161 McDonald, D., Clemente, J. C., Kuczynski, J., Rideout, J. R., Stombaugh, J., Wendel, D., Wilke, A., Huse,
162 S., Hufnagle, J., Meyer, F., Knight, R., and Caporaso, J. G. (2012). The Biological Observation Matrix
163 (BIOM) format or: how I learned to stop worrying and love the ome-ome. *GigaScience*, 1(1):7.
- 164 McMurdie, P. J. and Holmes, S. (2013). phyloseq: an R package for reproducible interactive analysis and
165 graphics of microbiome census data. *PLoS one*, 8(4):e61217.

- 166 Mukherjee, S., Huntemann, M., Ivanova, N., Kyrpides, N. C., and Pati, A. (2015). Large-scale con-
167 tamination of microbial isolate genomes by Illumina PhiX control. *Standards in genomic sciences*,
168 10:18.
- 169 Nakamura, T., Yamada, K. D., Tomii, K., and Katoh, K. (2018). Parallelization of MAFFT for large-scale
170 multiple sequence alignments. *Bioinformatics (Oxford, England)*, 34(14):2490–2492.
- 171 Price, M. N., Dehal, P. S., and Arkin, A. P. (2010). FastTree 2—approximately maximum-likelihood trees
172 for large alignments. *PloS one*, 5(3):e9490.
- 173 QIIME 2 Development Team (2018). QIIME 2. <https://qiime2.org/>, last accessed on 09/08/2018.
- 174 Rognes, T., Flouri, T., Nichols, B., Quince, C., and Mahé, F. (2016). VSEARCH: a versatile open source
175 tool for metagenomics. *PeerJ*, 4(17):e2584.
- 176 Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., Lesniewski, R. A.,
177 Oakley, B. B., Parks, D. H., Robinson, C. J., Sahl, J. W., Stres, B., Thallinger, G. G., Van Horn,
178 D. J., and Weber, C. F. (2009). Introducing mothur: open-source, platform-independent, community-
179 supported software for describing and comparing microbial communities. *Applied and Environmental*
180 *Microbiology*, 75(23):7537–7541.
- 181 Zenodo (2018). Zenodo - Research. Shared. <https://zenodo.org/>, last accessed on 09/08/2018.