# Extending MetAMOS - new methods and new integrations

Marian Siwiak, Albert Bogdanowicz, Agnieszka Hajduk, Michał Krassowski, Paweł Jankowski, Mariia Savenko, Adam AP Pyzik, Paweł Szczęsny

Biodiversity analysis of metagenomic and metatranscriptomic data acquired from next-generation sequencing (NGS) requires following multiple analytic steps, often independent from each other with exception of passing output files of previous step as input for the following. If parameterization of steps following one after another is independent from one another, they may be pipelined. There are three most popular pipelines used for NGS analyses: QIIME, mothur and MetAMOS. In this work we describe our extensions to the latter. One is supplementing MetAMOS' default modes with taxonomic and metabolic biodiversity using metagenomics and metatranscriptomics data and the other provides a web-based interface to run predefined analyses that is easy to integrate with laboratory information management systems.

# Extending MetAMOS - new methods and new integrations

Marian Siwiak[1,2,3#], Albert Bogdanowicz[2#], Agnieszka Hajduk[4], Michał Krassowski[4], Paweł Jankowski[4], Mariia Savenko[4], Adam Pyzik[2], Paweł Szczęsny[2,5*]

1 Oxford Data Science Centre, 154 Oxford Road, Oxford, Oxfordshire, OX4 2EB, United Kingdom

2 Institute of Biochemistry and Biophysics Polish Academy of Sciences, ul. Pawinskiego 5a, 02-106 Warsaw, Poland

3 Applied Research Institute for Prospective Technologies, Vismaliukų str. 34, Vilnius, Lithuania

4 University of Warsaw, Krakowskie Przedmiescie 26/28, Warsaw, Poland

5 Faculty of Biology, University of Warsaw, Miecznikowa, Warsaw, Poland

# both authors equally contributed to this study

* Correspondence should be addressed to:

Pawel Szczesny

Department of Bioinformatics

Institute of Biochemistry and Biophysics Polish Academy of Sciences

ul. Pawinskiego 5A

02-106 Warsaw

Poland

Email: szczesny@ibb.waw.pl

29

# Abstract

31  Biodiversity analysis of metagenomic and metatranscriptomic data acquired from next-
32  generation sequencing (NGS) requires following multiple analytic steps, often independent from
33  each other with exception of passing output files of previous step as input for the following. If
34  parameterization of steps following one after another is independent from one another, they may
35  be pipelined. There are three most popular pipelines used for NGS analyses: QIIME, mothur and
36  MetAMOS. In this work we describe our extensions to the last package. One extension is
37  supplementing MetAMOS' default modes with taxonomic and metabolic analyses on
38  metagenomics and metatranscriptomics data and the other extension provides a web-based
39  interface to run predefined analyses that is easy to integrate with laboratory information
40  management systems.

41

42  Keywords: metagenomics, metatranscriptomics, NGS, pipeline, web service

# Introduction

44  There is a number of standardized protocols available for NGS data dubbed Standard Operating
45  Procedures (SOPs). However, it is not uncommon for researchers to perform own sets of
46  analyses, depending on specific research topic needs. The most popular pipelines for NGS
47  analyses are QIIME (Caporaso et al. 2010), mothur (Schloss et al. 2009), and MetAMOS
48  (Schloss et al. 2009; Treangen et al. 2013). The first two are focused on biodiversity assessment
49  while the latter is aiming at assembly. Proper choice, deployment and pipelining of tools used in
50  such analyses is a task non-trivial even for an experienced bioinformatician, and may pose a big
51  problem for non-specialists.
52  The most easy to extended is definitely QIIME, as is it already a set of Python scripts. However
53  its installation is so challenging, that the installation method recommended by QIIME authors is
54  through VirtualBox or to use cloud computing instances with QIIME pre-installed. Mothur,
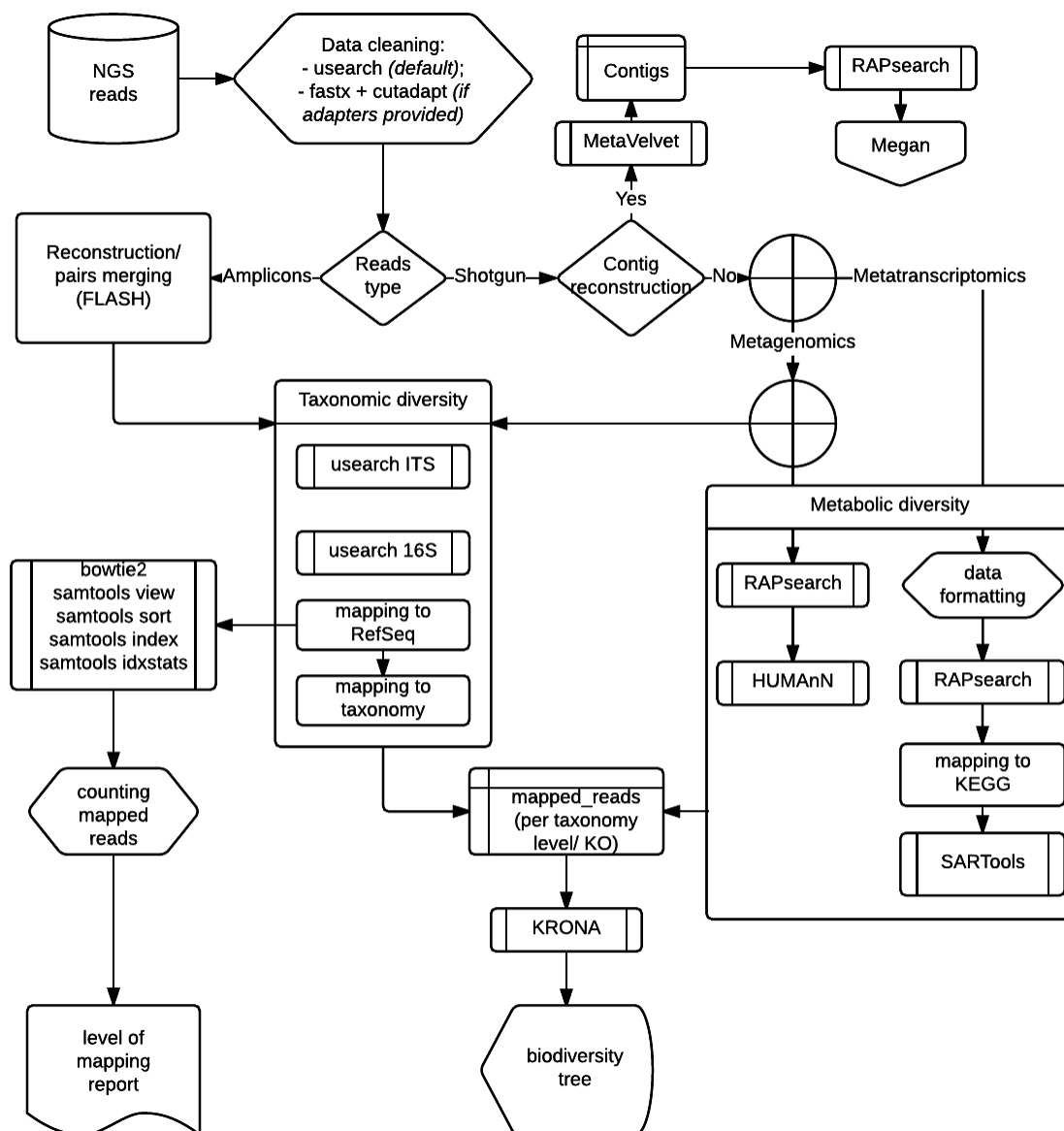55  being extremely portable due implementation in binary files, is hard to extend. In our case,

56    requiring work with multiple instances of the software in a dynamic hardware environment, we

57    set out to work with MetAMOS.

58    MetAMOS focus - namely creation of automated, reproducible, assembly & analysis pipeline

59    was in full alignment with our project. Moreover, it enabled fast and efficient implementation of

60    extension of functionalities, developed by our team.

61    This work presents an example of the solution to an effective cooperation between computational

62    and experimental biologists, merging complex analysis pipeline with user-friendly interface

63    through combination of MetAMOS with custom pipeline and web-interface.


64    **Custom pipeline**

65    Synthesis of research needs and technical capabilities of the research infrastructure required

66    development of in-house tool called bipype (Fig 1). "bipype" stands for bioinformatics-python-

67    pipe.

68

69

70 Fig. 1. Chart with bipype functionalities

71 Bipype accepts three types of inputs: amplicons, WGS (whole genome sequences) and

72 metatranscriptomic data. Bipype may work with paired-end and single read sequences. For

73 amplicon (prokaryotic or fungal) data, if needed, paired-end reads are merged and sequence

74 reconstruction is performed. Reconstructed 16S or ITS sequences are searched in proper

75 reference databases and hits to related taxonomic units are counted.

76 For WGS reads three paths are available. In one they will either be used for reconstruction of

77 contigs, which will be further used for reference database search (outside of the pipeline). In the

78   second they will be used directly in taxonomic diversity search, in third they will be compared to

79   sequences related to metabolic pathways.

80   The biodiversity data at taxonomical and functional levels are similar in their tree structure - we

81   provide display of results in a common form of html interactive comparative (multiple samples

82   presented in single file) tree (Fig. 3c,). Bipype is a portable solution with a few dependencies, but

83   to simplify analyses for the end-users, GUI presenting available analysis pathways and

84   (optional!) parameters of each step was required. Below we present MetAMOS integration

85   process, meeting this requirement. Bipype source code is available at:

86   https://github.com/krassowski/bipype/

## Integration with MetAMOS

88   As the number of customisable steps in MetAMOS is limited, we needed to overcome certain

89   inflexibility of this system by inserting the whole bipype run in the first possible step and

90   terminating the workflow afterwards. In order to do so, we added bipype to the list of tools used

91   by metAMOS as a fake assembler. We added a script that created empty files required by

92   metAMOS when it is checking if the assembly step finished successfully. Sample configuration

93   file for such fake assembler is shown in Fig. 2. Command line options for metAMOS allow

94   skipping most of the steps in any workflow, making it possible to finish immediately after

95   assembly.

96

```
[CONFIG]
name bipype_amplicons
input FASTQ
scaffoldOutput [RUNDIR]/Scaffolds.fasta
location python/bipype
output 16S_ITS.krona
threads --threads
unpaired [FIRST]
paired_interleaved [FIRST] [SECOND]
paired [FIRST] [SECOND]
commands mkdir [RUNDIR] && \
      bipype --out_dir [RUNDIR] --cutadapt use_paths both --mode run -ITS -16S -ot ITS
16S --input [INPUT] && \
      bipype_cheat [RUNDIR]
```

97   Fig. 2. Sample configuration file

### Web interface to MetAMOS

We have developed a web interface to analyses run on MetAMOS to let others not only access results but to run preconfigured analyses on their own. Fig. 3a shows the webpage where user can choose the type of analysis and library on which it would be run. Each type of analysis corresponds to one configuration file with appropriate bipype options. Status of analyses for each sample and workflow is stored in an SQL database. When user selects a workflow and a sample, web service checks if the corresponding job is finished and either shows the result (Fig. 3c), or its computation progress (Fig. 3b). If the job is not even started, the web service adds it to the database. A background service script periodically queries the database and runs queued jobs.

Additionally, we have added a simple results management, that is possibility to remove results in case re-computation is needed.



Fig. 3a. screenshot of the interface: sample and workflow choice

112

113    Fig 3b. screenshot of the interface: job progress
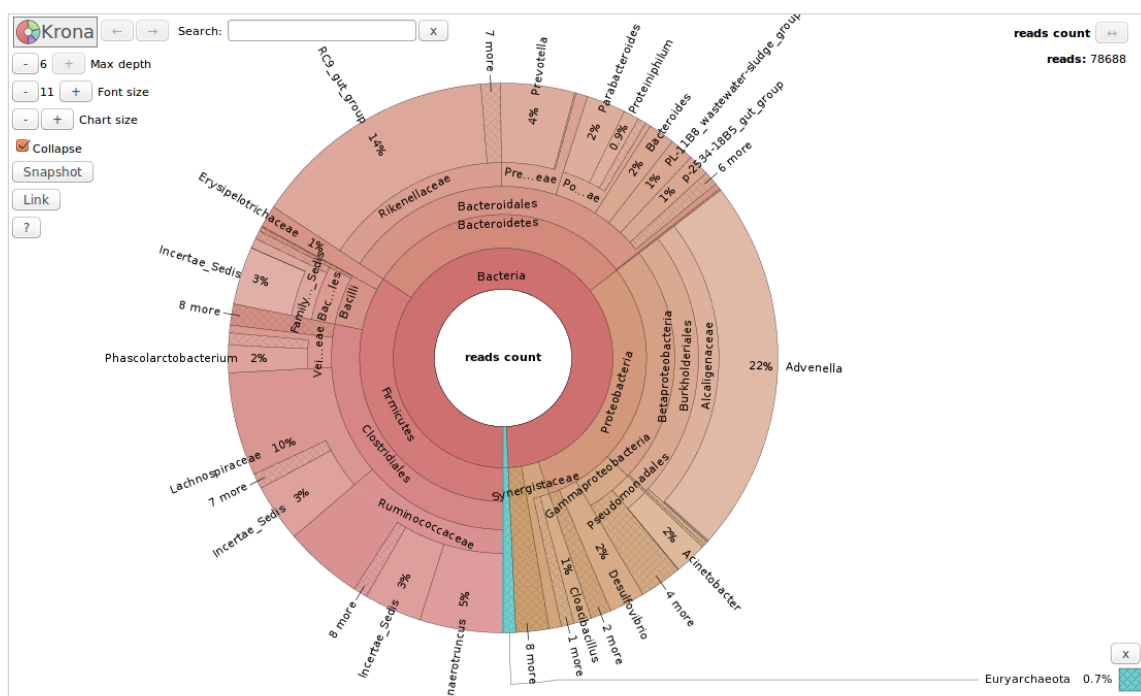


114

115    Fig 3c. screenshot of the interface: job results

116    Web interface source code is available at:

117    https://bitbucket.org/Serpens/metamos_web_interface/src

# Materials and Methods

bipype is written in Python v2.7.3 and links following tools: fastx-toolkit v0.0.13 ("FASTX-Toolkit" 2016), usearch v7.0.959_i86linux32 (Edgar 2010), FLASH v1.2.7 (Magoč and Salzberg 2011), bowtie2 v2.2.4 (Langmead et al. 2009), samtools v0.1.18 (Li et al. 2009), RAPsearch 2.12_64bits (Zhao, Tang, and Ye 2012), MetaVelvet v1.2.01 (Namiki et al. 2012), HUMAnN v0.99 (Abubucker et al. 2012), SARTools v1.2.0 (Varet et al. 2015), KRONA 2.0 (Ondov, Bergman, and Phillippy 2013). Presented workflow was performed on Illumina reads with varying insert lengths, either provided in form of parameter or read from filename.

Databases used as reference include: SILVA rRNA database (Griffith, Malachi, and Griffith 2004; Quast et al. 2013), Unified system for the DNA based fungal species linked to the classification (UNITE) (Kõljalg et al. 2013), The Reference Sequence (RefSeq) Database (Griffith, Malachi, and Griffith 2004), NCBI Taxonomy Database (Wheeler 2004), Kyoto Encyclopaedia of Genes and Genomes database (Kanehisa et al. 2014),

Web interface is build using MetAMOS v.1.5rc3, Python v2.7.3, Django v1.4.5, jQuery v1.11 and Bootstrap v3.1.1. It was tested with SQLite v3.7.13 and Apache v2.2.22.

# Conclusions

We present a way to develop time and resource efficient customizable pipeline serving for metagenomic and metatranscriptomic analyses, by extending MetAMOS workflow engine.

The proposed solution, besides abovementioned efficiencies presents following benefits:

- modularity allowing insertion of more custom tools and analyses (use of different assembly and display tools, search engines etc.),
- user-friendly web interface enabling easy access and steep learning curve for new team members, also enabling quick and repeated complete analyses by personnel not focused on software development,
- unified display methodology for both types of diversity data.

Our extensions are obviously available as open source under GNU GPLv2 license, allowing other researchers to build upon our work.

146

## Acknowledgements

## References

151 Abubucker, Sahar, Nicola Segata, Johannes Goll, Alyxandria M. Schubert, Jacques Izard, Brandi L.
152     Cantarel, Beltran Rodriguez-Mueller, et al. 2012. "Metabolic Reconstruction for Metagenomic Data
153     and Its Application to the Human Microbiome." *PLoS Computational Biology* 8 (6): e1002358.
154 Caporaso, J. Gregory, Justin Kuczynski, Jesse Stombaugh, Kyle Bittinger, Frederic D. Bushman,
155     Elizabeth K. Costello, Noah Fierer, et al. 2010. "QIIME Allows Analysis of High-Throughput
156     Community Sequencing Data." *Nature Methods* 7 (5): 335–36.
157 Edgar, Robert C. 2010. "Search and Clustering Orders of Magnitude Faster than BLAST." *Bioinformatics*
158     26 (19): 2460–61.
159 "FASTX-Toolkit." 2016. Accessed January 27. http://hannonlab.cshl.edu/fastx_toolkit/.
160 Griffith, Malachi, Griffith Malachi, and Obi L. Griffith. 2004. "RefSeq (the Reference Sequence
161     Database)." In *Dictionary of Bioinformatics and Computational Biology*.
162 Kanehisa, Minoru, Susumu Goto, Yoko Sato, Masayuki Kawashima, Miho Furumichi, and Mao Tanabe.
163     2014. "Data, Information, Knowledge and Principle: Back to Metabolism in KEGG." *Nucleic Acids*
164     *Research* 42 (Database issue): D199–205.
165 Kõljalg, Urmas, R. Henrik Nilsson, Kessy Abarenkov, Leho Tedersoo, Andy F. S. Taylor, Mohammad
166     Bahram, Scott T. Bates, et al. 2013. "Towards a Unified Paradigm for Sequence-Based Identification
167     of Fungi." *Molecular Ecology* 22 (21): 5271–77.
168 Langmead, Ben, Cole Trapnell, Mihai Pop, and Steven L. Salzberg. 2009. "Ultrafast and Memory-
169     Efficient Alignment of Short DNA Sequences to the Human Genome." *Genome Biology* 10 (3): R25.
170 Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo
171     Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. 2009. "The
172     Sequence Alignment/Map Format and SAMtools." *Bioinformatics* 25 (16): 2078–79.
173 Magoč, Tanja, and Steven L. Salzberg. 2011. "FLASH: Fast Length Adjustment of Short Reads to
174     Improve Genome Assemblies." *Bioinformatics* 27 (21): 2957–63.
175 Namiki, Toshiaki, Tsuyoshi Hachiya, Hideaki Tanaka, and Yasubumi Sakakibara. 2012. "MetaVelvet: An
176     Extension of Velvet Assembler to de Novo Metagenome Assembly from Short Sequence Reads."
177     *Nucleic Acids Research* 40 (20): e155.
178 Ondov, Brian D., Nicholas H. Bergman, and Adam M. Phillippy. 2013. "Krona: Interactive Metagenomic
179     Visualization in a Web Browser." In *Encyclopedia of Metagenomics*, 1–8.
180 Quast, Christian, Elmar Pruesse, Pelin Yilmaz, Jan Gerken, Timmy Schweer, Pablo Yarza, Jörg Peplies,
181     and Frank Oliver Glöckner. 2013. "The SILVA Ribosomal RNA Gene Database Project: Improved
182     Data Processing and Web-Based Tools." *Nucleic Acids Research* 41 (Database issue): D590–96.
183 Schloss, Patrick D., Sarah L. Westcott, Thomas Ryabin, Justine R. Hall, Martin Hartmann, Emily B.
184     Hollister, Ryan A. Lesniewski, et al. 2009. "Introducing Mothur: Open-Source, Platform-
185     Independent, Community-Supported Software for Describing and Comparing Microbial
186     Communities." *Applied and Environmental Microbiology* 75 (23): 7537–41.
187 Treangen, Todd J., Sergey Koren, Daniel D. Sommer, Bo Liu, Irina Astrovskaya, Brian Ondov, Aaron E.
188     Darling, Adam M. Phillippy, and Mihai Pop. 2013. "MetAMOS: A Modular and Open Source
189     Metagenomic Assembly and Analysis Pipeline." *Genome Biology* 14 (1): R2.

190   Varet, Hugo, Varet Hugo, Coppée Jean-Yves, and Dillies Marie-Agnès. 2015. "SARTools: A DESeq2-
191        and edgeR-Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq Data."
192        doi:10.1101/021741.
193   Wheeler, D. L. 2004. "Database Resources of the National Center for Biotechnology Information."
194        *Nucleic Acids Research* 33 (Database issue): D39–45.
195   Zhao, Yongan, Haixu Tang, and Yuzhen Ye. 2012. "RAPSearch2: A Fast and Memory-Efficient Protein
196        Similarity Search Tool for next-Generation Sequencing Data." *Bioinformatics*  28 (1): 125–26.

197