

BioCicle: A Tool for Summarizing and Comparing Taxonomic Profiles out of Biological Sequence Alignments

Meili Vanegas-Hernandez*, Fabio Andres Lopez-Corredor*, Tiberio Hernandez*,
Alejandro Reyes* and John Alexis Guerra-Gomez**

*Universidad de los Andes, Colombia **Northeastern University, Silicon Valley, USA.

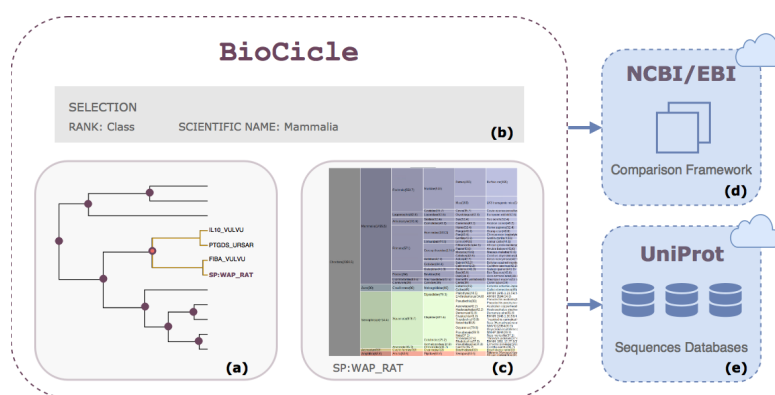


Figure 1: BioCicle key features. We present BioCicle, an open source and web-based application for summarizing and comparing taxonomic profiles for biological sequence alignments. BioCicle supports several input formats, as well as direct sequence comparisons using the NCBI/EBI's (d) and UniProt's (e) APIs. Once data has been loaded, BioCicle presents an overview of sequence comparisons results (a) and allows the user to filter regions of interest (b) to specifically analyze the most representative species that match to the each query according to the score of similarity (c).

Keywords: Information Visualization, Bioinformatics, Biological Sequence Comparison,

1. Introduction

2 A common practice in metagenomics consists in collecting genetic material
3 from environmental samples, for later classifying them by comparing them to
4 existing sequences in large known biological databases. BLAST and HMM are
5 among the most common tools used for achieving this goal. The result of the

6 comparison consists of a large output that includes information about the similarity
7 of the sampled sequences against the closest matches in the database. There is a
8 lack of effective tools for summarizing and comparing such output, which leads
9 overwhelmed analyst to keep only the first most probable result, ignoring all the
10 remaining ones. This can lead to miss-classification, as more meaningful matches
11 could be hidden in the remaining results. For example, think of a result that ranks
12 a frog as the first result with a score of 70%, but the remaining results are bacteria
13 with scores around 60%. An analysts could miss-classify the sequence as a frog,
14 when there are high chances of being a bacteria.

15 To address this problem, this paper presents BioCicle, an interactive visual
16 analytics system that summarizes all the results in a sequence comparison output,
17 highlighting the corresponding scores, and allowing analysts to make decisions
18 based on the whole input. Moreover, BioCicle also supports the summarization of
19 multiple outputs, to identify patterns in groups of samples. We validated BioCicle
20 by working closely with domain experts and by means of a case study presented in
21 the paper.

22 BioCicle was designed following Munzner’s visual analytics framework [1],
23 identifying the main analysis tasks that analysts usually performed when doing
24 sequence comparisons. From our close work with domain experts, we identified
25 that sequence comparisons usually output three type of results:

- 26 • (RS1) Sequence alignments
- 27 • (RS2) Taxonomic reports
- 28 • (RS3) Aligned sequences’ descriptions

29 Building on this types of results, and considering that domain experts can run
30 single or multiple queries depending on the number of samples they collected, this
31 paper contributes an analysis task taxonomy illustrated in 2.

32 **2. Related work**

33 After a task-driven analysis developed along with a group of bioinformaticians,
34 we proposed a taxonomy of the state of the art. The six tasks were either
35 *summarizing* or *comparing* for the three different outputs of interest: sequence
36 alignments (AT1), taxonomic reports (AT2) and sequences’ descriptions (AT3).
37 BioCicle focuses on AT2 allowing taxonomic profiles summarization (AT2a) and

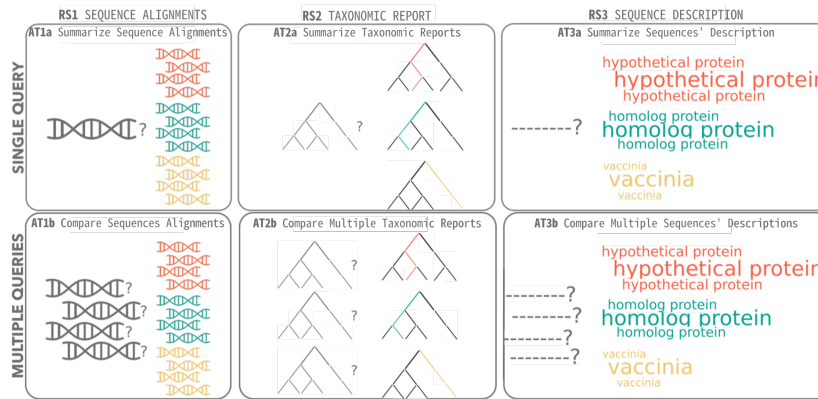


Figure 2: Taxonomy of Analysis tasks for biological sequence alignments.

	RS1 Sequence Alignments		RS2 Taxonomic Report		RS3 Sequences' Descriptions	
	AT1a Summarize Sequence Alignments	AT1b Compare Multiple Sequences' Alignments	AT2a Summarize Taxonomic Reports	AT2b Compare Multiple Taxonomic Reports	AT3a Summarize Sequences' Descriptions	AT3b Compare Multiple Sequences' Descriptions
Restrictive	Bov Blast2Go Artemis HnnEditor	Circoletto ClustalW Hmmer Geneclusterviz Megan BlastGrabber	MG-Rast	Megan Metarep BlastGrabber		
Non-Restrictive			AmphoraVizu MetaPhlan Krona BioCicle	BioCicle		OUR ONGOING WORK

Figure 3: Taxonomy of state-of-the-art implementations supporting the analysis tasks defined in Figure 2.

38 comparisons (AT2b) for single and multiple queries using visual analytics design
 39 principles [1] and withstanding multiple input formats.

40 We made a revision of 17 tools presenting sequence alignment results and
 41 classified each tool considering which analysis task they addressed. As a result, we
 42 present a taxonomy (Table 3).

43 The vast majority of tools are based in single-sequence alignment summarization
 44 and comparison (AT1). Nonetheless, some approaches tackle taxonomic reports
 45 summarization [2], [3], [4], [5] for single-query alignments (AT2a) and comparison
 46 [6], [7], [8] for multiple comparisons results (AT2b).

47 Most of the tools identified are focused in single-query displays. The main
 48 drawback is that each sequence alignment must be analyzed independently, which
 49 leads to a highly cost-intensive understanding of the results.

50 **3. Implementation**

51 *3.1. AT2a: Summarize Taxonomic Reports for a Single Sequence Comparison*

52 The visualization consisted in three different components: an icicle map, a
53 collapsible tree and a set of small multiples, as shown in Figure ??.

54 1. **Icicle Tree:** Each dimension of the taxonomy (i.e. class, genus, order, etc.)
55 was located in a different column of the icicle. Also, each dimension was
56 treated as a nominal variable and represented using the spatial region in the
57 limited column. The score value, as a numerical variable, was represented
58 using the length of the columns divisions. Therefore, the height of each level
59 was calculated with a linear scale having the score value as the domain and
60 the number of pixels as the range. As an usual icicle tree, the child nodes
61 score values/length contributed to the parent's value, meaning that an entire
62 column considered the 100% of the sequences displayed.

63 2. **Collapsible Tree:** The result set considered for this task was a grouping of
64 multiple comparison results. Each comparison's result was represented as a
65 list, as the one described in the previous subsection. As the explained before,
66 each of these lists represented a tree, having as leaves the species with which
67 the unknown was compared and the score of similarity. This resulting tree is
68 called a taxonomic report. The grouping of those taxonomic reports could
69 be interpreted as a conglomeration of trees.

70 3. **Small Multiples:** We presented the total amount of icicles as small multiples
71 for each of the unclassified sequences. indicating over which subgroup it
72 was being iterating.

73 Our implementation allows to dynamically explore taxonomic reports out of
74 multiple comparisons comparisons and compare general characteristics out of them,
75 supporting a considerable amount of unclassified species.

76 **4. Results**

77 In order to validate the framework, this case study considered a set of 179
78 sequences. These sequences were obtained by sequencing a genetic marker over
79 a single sample. The gene was the 16S ribosomal rRNA, which is a molecular
80 genetic marker and it is used to detect variations in DNA. In this case, sequencing
81 such gene is meant to identify the diversity of the sample. The original dataset

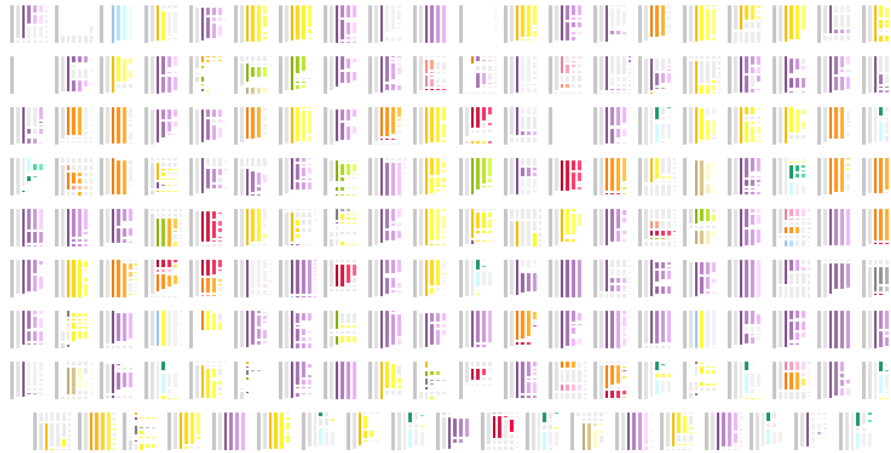


Figure 4: Overview of the taxonomic profiles for the entire 16S dataset using *small multiples*. This visualization shows an overview of the trees generated for each of the compared sequence. Colors and spatial regions are preserved for nodes all over the entire for ease of understanding. Prevailing tones are purple, yellow, and orange.

82 consisted in 23.000 sequences but, considering time limitations, only 179 were
 83 employed for this study. The results are shown in Figure4.

84 The small multiples allow the user to have a general overview of the results. In
 85 Figure 4 it is possible to see there is a diverse distribution of ranks over the sample,
 86 however, yellow, purple, and orange tones prevail in the general sight.

87 After analyzing the results for this specific sample, it was possible to quickly
 88 describe the dataset as sequences belonging mostly to class *Proteobacteria*, with a
 89 high amount (49,16%) having high similarities to sequences belonging to the order
 90 *Gammaproteobacteria*. Additional sequences had high similarities with sequences
 91 belonging to classes *Alphaproteobacteria* (29,05%), *Deltaproteobacteria* (7,82%)
 92 and *Flavobacteria* (12,29%).

93 5. Validation

94 To test the contribution of BioCicle in the performing of the identified tasks by
 95 bioinformaticians, potential users were asked to use the tool and validate the easier
 96 approach to common tasks. A recorded session with two computational biology
 97 students and a researcher in bioinformatics was done, in which we could evaluate
 98 both the usability and usefulness of the platform. They were also prompted to
 99 identify the potential tasks that BioCicle allows them to do. A standard usability
 100 scale was used to evaluate how well each task is performed by the web platform.

101 From the recordings of testing sessions we got the proof from actual users
 102 that BioCicle is able to help in a quick processing of the results from a BLASTp
 103 execution. The findings from the tests, aside from web development usability,
 104 allowed the team to define the respective upcoming developments around two
 105 fronts: Specific tasks and new functionalities.

106 5.1. Specific tasks

107 The users determined which tasks, in terms of their domain vocabulary, are
 108 able to be executed with the tool's assist. For each task they graded BioCicle in a
 109 difficulty scale from 1 to 7, where 1 was a task hard to perform with the web page,
 110 and 7 a task easily executed with the help of the platform.

111 The description in the terms of the platform of these tasks, with their user-assigned
 112 grades, is well explained in Table 5. In first instance, we recognized most identified
 113 tasks as similar to the original analysis tasks identified (AT2a and AT2b from
 114 Figure 2). User-defined tasks 1, 3, 4 and 6 make reference to sequence alignments,
 115 including both single and multiple queries. Task 2 is an analysis work that could
 116 be done after performing the identified analysis task, while task 5 is an additional
 117 function of filtering to give more control to users of the information entered to the
 118 application.

User Defined Task	Grade (1-7)	Description in terms of BioCicle	Related Analysis Tasks
UD1	5	Visualization of taxonomy from entered comparison file from sequences of virus proteins	AT2a & AT2b
UD2	6	Explore nodes in collapsible tree and Icicle and determine which nodes are from a different nature from the expected from the initially uploaded sequences	Analysis done after results from AT2b
UD3	2	For an Icicle leave, understand how it ranks among other leaves according to their score.	AT2a
UD4	6	Identify the nodes from the collapsible tree or Icicle with the highest scores	AT2b
UD5	6	Select a node from the collapsible tree, and download the filtered xml file.	Filtering results from AT2b
UD6	7	For a comparison file from a single sequence, visualize the taxonomy in the Icicle	AT2a

Figure 5: Comparison between User-defined tasks and Analysis tasks identified in Figure 2.

119 The grades for the user-defined tasks turned out to be favorable with the
 120 exception of user-defined task 3. In the session the user who identified and graded
 121 this task, a student mainly interested in assembling and classification of virus,
 122 stated that the results were displaying both the taxonomy for the virus and the
 123 bacteria where the virus is found at. As it is unclear if the taxonomy is calculated
 124 from the sequence from the virus or the bacteria containing it, it is hard for the user
 125 to understand if the sampled virus actually corresponds to the retrieved genome.

126 6. Conclusions and Future Work

127 Biological sequence comparisons are a widely used methodology for sequence
 128 identification and characterization. Such methodology assists detection of regions

129 of similarity between DNA, RNA or protein sequences, which may imply evolutionary
130 relationships between species. Sequence's comparison outputs have often a
131 considerable amount of information, thus, fast extraction of relevant information is
132 a very costly process. Besides, sequence misclassification can be easily achieved if
133 the sequence comparison output is misread. Comparable mistakes affect not only
134 the new individual classification, but also it ensures future misclassifications, as
135 such sequence will be part of the comparison set in future sequence alignments.

136 This project introduces BioCicle, a tool to summarize and compare either single
137 or multiple results displays for taxonomic profiles in sequence alignments (AT2). It
138 was developed following the visual analytics principles. In addition, the application
139 was directly connected to the NCBI/EBI and UniProt API's, allowing custom
140 comparison generation on demand. BioCicle was constantly tested and evaluated
141 along with a group of bioinformaticians and presented as a proof of concept.

142 Although BioCicle tackles an unaddressed problematic (AT2b) with non-restrictive
143 characteristics, there is still an untapped potential in sequence description analysis
144 for either single or multiple comparisons displays (AT3). Ideally, BioCicle could be
145 extended as a decision support system framework for sequence alignment analysis
146 for taxonomic reports and sequence's description. Methods such as text-analysis,
147 feature selection and data mining could ease sequence's description analysis and
148 decrease incorrect insertions rates in biological databases.

149 **7. References**

- 150 [1] T. Munzner, Visualization analysis and design, CRC press, 2014.
- 151 [2] D. Huson, MEGAN Community Edition - Interactive exploration and analysis
152 of large-scale microbiome sequencing data., PLoS Computational Biology 12
153 (2016) e1004957.
- 154 [3] B. D. Ondov, N. H. Bergman, A. M. Phillippy, Interactive metagenomic
155 visualization in a Web browser, BMC Bioinformatics 12 (2011) 385.
- 156 [4] B. A. Goll, Johannes; Rusch, Douglas B; Tanenbaum, David M;
157 Thiagarajan, Mathangi; Li, Kelvin; Methé, S. Yooseph, METAREP:
158 JCVI Metagenomics Reports - an open source tool for high- performance
159 comparative metagenomics, Bioinformatics 26 (2010) 2631–2632.
- 160 [5] F. Meyer, D. Paarmann, M. D'Souza, Etal., The metagenomics RAST
161 server—a public resource for the automatic phylo- genetic and functional
162 analysis of metagenomes, BMC bioinformatics 9 (2008) 386.

- 163 [6] Y. Zhai, J. Tchieu, M. H. Saier, JMMB Bioinformatics Corner A Web-Based T
164 ree V iew (TV) Program for the Visualization of Phylogenetic Trees, Journal
165 of molecular microbiology and 4 (2002) 69–70.
- 166 [7] N. Segata, L. Waldron, A. Ballarini, V. Narasimhan, O. Jousson,
167 C. Huttenhower, Metagenomic microbial community profiling using unique
168 clade-specific marker genes., Nature methods 9 (2012) 811–4.
- 169 [8] C. Kerepesi, B. Szalkai, V. Grolmusz, Visual Analysis of the Quantitative
170 Composition of Metagenomic Communities: the AmphoraVizu Webserver,
171 Microbial Ecology 69 (2015) 695–697.