

1 ***Phylo-Node: a molecular phylogenetic toolkit using Node.js***

2

3

4 Damien M. O'Halloran<sup>1,2\*</sup>

5

6 1. Department of Biological Sciences, The George Washington University, Science  
7 and Engineering Hall, Rm 6000, 800 22nd Street N.W., Washington DC 20052, USA.

8 2. Institute for Neuroscience, The George Washington University, 636 Ross Hall,  
9 2300 I St. N.W. Washington DC 20052, USA.

10

11

12

13 **\*Corresponding author address:**

14 Damien O'Halloran, The George Washington University, 636 Ross Hall, 2300 I St.

15 N.W. Washington DC 20052, USA.

16 Tel: 202-994-8955

17 Fax: 202-994-6100

18 EMAIL: [damienoh@gwu.edu](mailto:damienoh@gwu.edu)

19

20

21

22

23

24

25

26 **ABSTRACT**

27 **Background:** Node.js is an open-source and cross-platform environment that  
28 provides a JavaScript codebase for back-end server-side applications. JavaScript  
29 has been used to develop very fast, and user-friendly front-end tools for bioinformatic  
30 and phylogenetic analyses. However, no such toolkits are available using Node.js to  
31 conduct comprehensive molecular phylogenetic analysis.

32 **Results:** To address this problem, I have developed, *Phylo-Node*, which was  
33 developed using Node.js and provides a stable and scalable toolkit that allows the  
34 user to go from sequence retrieval to phylogeny reconstruction. Phylo-Node can  
35 execute the analysis and process the resulting outputs from a suite of software  
36 options that provides tools for sequence retrieval, alignment, primer design,  
37 evolutionary modeling, and phylogeny reconstruction. Furthermore, Phylo-Node  
38 enables the user to deploy server dependent applications, and also provides simple  
39 integration and interoperation with other Node modules and languages using Node  
40 inheritance patterns, and a customized piping module to support the production of  
41 diverse pipelines.

42 **Conclusions:** Phylo-Node is open-source and freely available to all users without  
43 sign-up or login requirements. All source code and user guidelines are openly  
44 available at the GitHub repository: <https://github.com/dohalloran/Phylo-Node>

45

46 **Keywords:** Node.js, JavaScript, phylogenetics

47

48

49

50

## 51 **BACKGROUND**

52 The cost of whole genome sequencing has plummeted over the last decade and as  
53 a consequence, the demand for genome sequencing technology has risen  
54 significantly [1]. This demand has meant that producing large complex datasets of  
55 DNA and RNA sequence information is common in small research labs, and in terms  
56 of human health this boom in sequence information and precipitous drop in  
57 sequencing costs has had a direct impact in the area of personalized medicine [2-5].  
58 However, once the sequence information becomes available, perhaps the greater  
59 challenge is then processing, analyzing, and interpreting the data. To keep pace with  
60 this challenge, the development of new, fast, and scalable software solutions is  
61 required to visualize and interpret this information.

62 JavaScript is a lightweight programming language that uses a web browser as  
63 its host environment. JavaScript is cross-platform and supported by all modern  
64 browsers. Because JavaScript is client-side, it is very fast, as it doesn't have to  
65 communicate with a server and wait for a response in order to run some code. Web  
66 browsers are ubiquitous and require no dependencies to deploy and operate, and so  
67 JavaScript represents an obvious solution for visualizing sequence information.  
68 Front-end developments using JavaScript have proven to be extremely efficient in  
69 providing fast, easy-to-use, and embeddable solutions for data analysis [6-14]. A  
70 very active community of developers at BioJS (<http://www.biojs.io/>) provides diverse  
71 components for parsing sequence data types, data visualization, and bioinformatics  
72 analysis in JavaScript [6, 7, 15-19].

73 Node.js provides server-side back-end JavaScript. Node.js is written in C,  
74 C++, and JavaScript and uses the Google Chrome V8 engine to offer a very fast  
75 cross-platform environment for developing server side Web applications. Node is a

76 single-threaded environment, which means that only one line of code will be  
77 executed at any given time; however, Node employs non-blocking techniques for I/O  
78 tasks to provide an asynchronous ability, by using *callback* functions to permit the  
79 parallel running of code. Node holds much potential for the bioinformatic analysis of  
80 molecular data. A community of Node developers provides modules for bioinformatic  
81 sequence workflows (<http://www.bionode.io/>) which in time will likely parallel the  
82 BioJS community for the number of modules versus components. However, as of  
83 now there are no robust tools for phylogenetic analysis pipelines currently available  
84 using the Node.js codebase. To fill this void I have developed, *Phylo-Node*, which  
85 provides a Node.js toolkit that goes from sequence retrieval, to primer design, to  
86 alignment, to phylogeny reconstruction, all from a single toolkit. MolPhylo is fast,  
87 easy to use, and offers simple customization and portability options through various  
88 inheritance patterns. The Node package manager, *npm* (<https://www.npmjs.com/>),  
89 provides a very easy and efficient way to manage dependencies for any Node  
90 application. *Phylo-Node* is available at GitHub ([https://github.com/dohalloran/Phylo-](https://github.com/dohalloran/Phylo-Node)  
91 [Node](https://github.com/dohalloran/Phylo-Node)), *npm* (<https://www.npmjs.com/package/phylo-node>), and also BioJS  
92 (<http://www.biojs.io/d/phylo-node>).

93

## 94 **IMPLEMENTATION**

95 *Phylo-Node* was developed using the Node.js codebase. The *Phylo-Node* core  
96 contains methods for remote sequence retrieval, and phylogenetic analysis using a  
97 suite of popular software tools. A base wrapper object is used to prepare the  
98 arguments and directory prior to program execution. The base wrapper module is  
99 contained within the 'Wrapper\_core' directory. An individual software tool can be  
100 easily accessed and executed by importing the module for that tool so as to get

101 access to the method properties on that object (Figure 1). These method properties  
102 are available to the user by using the 'module.exports' reference object. Inside a  
103 driver script file, the user can import the main module object properties and variables  
104 by using the 'require' keyword which is used to import a module in Node.js. The  
105 'require' keyword is actually a global variable, and a script has access to its context  
106 because it is wrapped prior to execution inside the 'runInThisContext' function (for  
107 more details, refer to the Node.js source code: <https://github.com/nodejs>). Once  
108 imported, the return value is assigned to a variable which is used to access the  
109 various method properties on that object. For example: a method property on the  
110 'phym1' object is 'phym1.getphym1()', which invokes the 'getphym1' method on the  
111 'phym1' object to download and decompress the PhyML executable. For a complete  
112 list of all methods, refer to the 'README' file at the GitHub repository  
113 (<https://github.com/dohalloran/Phylo-Node/blob/master/README.md>). In order to  
114 correctly wrap and run each executable, new shells must be spawned so as to  
115 execute specific command formats for each executable. This was achieved by using  
116 'child.process.exec', which will launch an external shell and execute the command  
117 inside that shell while buffering any output by the process. Binary files and  
118 executables were downloaded and executed in this manner and the appropriate file  
119 and syntax selected by determining the user's operating system. Phylo-Node was  
120 validated on Microsoft Windows 7 Enterprise ver.6.1, MacOSX El Capitan  
121 ver.10.11.5, and Linux Ubuntu 64-bit ver.14.04 LTS.

122

## 123 **RESULTS AND DISCUSSION**

124 Phylo-Node is a toolkit to interface with key applications necessary in building a  
125 phylogenetic pipeline (Figure 2). Firstly, Phylo-Node allows the user to remotely

126 download sequences by building a unique URL and passing this string to the NCBI  
127 e-utilities API (<http://www.ncbi.nlm.nih.gov/books/NBK25501/>). Any number of genes  
128 can be supplied as command-line arguments to Phylo-Node by accessing the  
129 *fetch\_seqs.fasta* method on the *fetch\_seqs* object in order to retrieve sequence  
130 information in FASTA format. The module for remote sequence retrieval is contained  
131 within the 'Sequence' directory. Phylo-Node also provides methods on specific  
132 objects to download various executable files using the 'download' module. Any  
133 binary can be downloaded using the base module 'get\_executable' contained within  
134 the 'Download' directory, however objects pertaining to specific tools such as PhyML  
135 also contain methods for downloading and unpacking binaries (see README.md file  
136 for details). Phylo-Node then provides modules to execute the following programs  
137 from within the './Tool/Run' directory: Primer3 [20-22] to facilitate primer design;  
138 Clustal Omega [23], K-align [24], and MUSCLE [25, 26] for multiple sequence  
139 alignments; Codeml [27], PAL2NAL [28], and Slr [29] for selection analysis;  
140 jModelTest2 [30] and ProtTest3 [31] to determine the best-fit model of evolution; and  
141 PhyML [32, 33] for phylogeny reconstruction. The PhyML executable is also  
142 employed by jModelTest2 and ProtTest3. Primer3 is the most popular software for  
143 primer design, and takes a very lengthy list of input variables to optimize primer  
144 selection. Clustal Omega, K-align, and MUSCLE are very fast and accurate multiple  
145 sequence alignment tools that are commonly used to build robust DNA, RNA, or  
146 protein alignments. Codeml is part of the PAML suite [27], and alongside PAL2NAL  
147 [28] and Slr [29] can be used to determine rates of selection. ProtTest3 determines  
148 the best-fit model of evolution for protein sequences across 120 different potential  
149 models, while jModelTest2 determines best-fit models of nucleotide substitution from  
150 DNA sequence alignments, and PhyML is a popular program for building

151 phylogenies using maximum likelihood. Together, Phylo-Node provides a toolkit that  
152 allows the user to go from raw sequence to phylogeny using Node.

153 Phylo-Node is highly scalable and customizable, and was inspired by projects  
154 such as BioPerl [34] which provides very diverse tools that include Perl modules for  
155 many bioinformatic tasks and also parsers and wrappers for diverse sequence  
156 formats and applications. BioPerl's open source structure and architecture allows  
157 users to plug new modules into BioPerl pipelines to design new applications. Node.js  
158 implements prototypal inheritance as per JavaScript but also provides access to the  
159 'module.exports' object which permits easy portability between the Phylo-Node  
160 toolkit and any other modules, and also interoperation between different languages  
161 by using the 'child.process.exec' process. Therefore, Phylo-Node can be integrated  
162 with existing Node.js bioinformatics tools [35, 36] or software written in other  
163 languages. For example, both jModelTest2 and ProtTest3 require a Java runtime  
164 environment ([http://www.oracle.com/technetwork/java/javase/downloads/jre8-](http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html)  
165 [downloads-2133155.html](http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html)), and by using 'require' to import each module, the user  
166 can execute the analysis of jModelTest2 and ProtTest3. The 'protest' and  
167 'jmodeltest2' modules and driver scripts (index.js) are contained within the respective  
168 'Protest3' and 'jModelTest2' directories and sample input is provided in the  
169 'COX2\_PF0016' sub-directory of the 'Input\_examples' folder for ProtTest3 and the  
170 sample FASTA file aP6.fas (also contained within the 'Input\_examples' folder) can  
171 be used for testing jModelTest2.

172 To further facilitate the ease of interoperation between various applications and  
173 components, the Phylo-Node package also contains a module called 'phylo-  
174 node\_pipes' inside the 'Pipes' directory. The 'phylo-node\_pipes' module allows the  
175 user to easily pipe data between different applications by requiring the

176 'child\_process' module which provides the ability to spawn child processes. Through  
177 'phylo-node\_pipes', the user can chain commands together that will be executed in  
178 sequence to build consistent, and extensive pipelines. The 'Pipes' directory contains  
179 sample driver scripts for using the 'phylo-node\_pipes' module. Finally, Phylo-Node  
180 also enables the user to create a web server to deploy embeddable applications  
181 such as JBrowse [37] which provides genome visualization from within a web  
182 browser.

183

## 184 **CONCLUSIONS**

185 In conclusion, Phylo-Node is a novel package that leverages the speed of Node.js to  
186 provide a robust and efficient toolkit for researchers conducting molecular  
187 phylogenetics. Phylo-Node can be easily employed to develop complex but  
188 consistent workflows, and integrated with existing bioinformatics tools using the  
189 Node.js codebase.

190

191

192

193

194

195

196

197

198

199

200



201

## 202 **AVAILABILITY AND REQUIREMENTS**

- 203 • Project name: Phylo-Node
- 204 • Project home page: <https://github.com/dohalloran/phylo-node>
- 205 • Operating system(s): Platform independent
- 206 • Programming language: Node.js
- 207 • Other requirements: none
- 208 • License: MIT
- 209 • Any restrictions to use by non-academics: no restrictions or login requirements

210

## 211 **ACKNOWLEDGMENTS**

212 I thank members of the O'Halloran lab for critical reading of the manuscript, and  
213 would like to thank The George Washington University (GWU) Columbian College of  
214 Arts and Sciences, GWU Office of the Vice-President for Research, and the GWU  
215 Department of Biological Sciences for Funding.

216

## 217 **AUTHOR CONTRIBUTIONS**

218 D.O'H. conceived the idea for *Phylo-Node*, wrote and tested the code, and wrote the  
219 manuscript.

220

## 221 **COMPETING INTERESTS**

222 The author declares no competing interests.

223

224

225

226

## References

- 227 1. Shaffer C: **Next-generation sequencing outpaces expectations.** Nat Biotechnol  
228 2007, **25**(2):149.
- 229 2. Wade N: **The quest for the \$1,000 human genome: DNA sequencing in the**  
230 **doctor's office? At birth? It may be coming closer.** N Y Times Web 2006, :F1,  
231 F3.
- 232 3. Mardis ER: **Anticipating the 1,000 dollar genome.** Genome Biol 2006, **7**(7):112.
- 233 4. Service RF: **Gene sequencing. The race for the \$1000 genome.** Science 2006,  
234 **311**(5767):1544-1546.
- 235 5. Hayden EC: **The \$1,000 genome.** Nature 2014, **507**(7492):294-295.
- 236 6. Yachdav G, Goldberg T, Wilzbach S, Dao D, Shih I, Choudhary S, Crouch S,  
237 Franz M, Garcia A, Garcia LJ, Gruning BA, Inupakutika D, Sillitoe I, Thanki AS,  
238 Vieira B, Villaveces JM, Schneider MV, Lewis S, Pettifer S, Rost B, Corpas M:  
239 **Anatomy of BioJS, an open source community for the life sciences.** Elife 2015,  
240 **4**:10.7554/eLife.07009.
- 241 7. Gomez J, Garcia LJ, Salazar GA, Villaveces J, Gore S, Garcia A, Martin MJ,  
242 Launay G, Alcantara R, Del-Toro N, Dumousseau M, Orchard S, Velankar S,  
243 Hermjakob H, Zong C, Ping P, Corpas M, Jimenez RC: **BioJS: an open source**  
244 **JavaScript framework for biological data visualization.** Bioinformatics 2013,  
245 **29**(8):1103-1104.
- 246 8. Salazar GA, Meintjes A, Mulder N: **PPI layouts: BioJS components for the**  
247 **display of Protein-Protein Interactions.** F1000Res 2014, **3**:50-50.v1. eCollection  
248 2014.
- 249 9. Gomez J, Jimenez R: **Sequence, a BioJS component for visualising**  
250 **sequences.** F1000Res 2014, **3**:52-52.v1. eCollection 2014.
- 251 10. Cui Y, Chen X, Luo H, Fan Z, Luo J, He S, Yue H, Zhang P, Chen R:  
252 **BioCircos.js: an interactive Circos JavaScript library for biological data**  
253 **visualization on web applications.** Bioinformatics 2016, **32**(11):1740-1742.
- 254 11. Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, Goodstein DM,  
255 Elsik CG, Lewis SE, Stein L, Holmes IH: **JBrowse: a dynamic web platform for**  
256 **genome visualization and analysis.** Genome Biol 2016, **17**(1):66-016-0924-1.
- 257 12. Salavert F, Garcia-Alonso L, Sanchez R, Alonso R, Bleda M, Medina I, Dopazo  
258 J: **Web-based network analysis and visualization using CellMaps.** Bioinformatics  
259 2016, .

- 260 13. Franz M, Lopes CT, Huck G, Dong Y, Sumer O, Bader GD: **Cytoscape.js: a**  
261 **graph theory library for visualisation and analysis**. *Bioinformatics* 2016,  
262 **32(2):309-311**.
- 263 14. Vanderkam D, Aksoy BA, Hodes I, Perrone J, Hammerbacher J: **pileup.js: a**  
264 **JavaScript library for interactive and in-browser visualization of genomic data**.  
265 *Bioinformatics* 2016, .
- 266 15. Garcia L, Yachdav G, Martin MJ: **FeatureViewer, a BioJS component for**  
267 **visualization of position-based annotations in protein sequences**. *F1000Res*  
268 2014, **3:47-47.v2**. eCollection 2014.
- 269 16. Kalderimis A, Stepan R, Sullivan J, Lyne R, Lyne M, Micklem G: **BioJS**  
270 **DAGViewer: A reusable JavaScript component for displaying directed graphs**.  
271 *F1000Res* 2014, **3:51-51.v1**. eCollection 2014.
- 272 17. Villaveces JM, Jimenez RC, Habermann BH: **KEGGViewer, a BioJS**  
273 **component to visualize KEGG Pathways**. *F1000Res* 2014, **3:43-43.v1**.  
274 eCollection 2014.
- 275 18. Villaveces JM, Jimenez RC, Habermann BH: **PsicquicGraph, a BioJS**  
276 **component to visualize molecular interactions from PSICQUIC servers**.  
277 *F1000Res* 2014, **3:44-44.v1**. eCollection 2014.
- 278 19. Yachdav G, Hecht M, Pasmanik-Chor M, Yeheskel A, Rost B: **HeatMapView:**  
279 **interactive display of 2D data in biology**. *F1000Res* 2014, **3:48-48.v1**. eCollection  
280 2014.
- 281 20. You FM, Huo N, Gu YQ, Luo MC, Ma Y, Hane D, Lazo GR, Dvorak J, Anderson  
282 OD: **BatchPrimer3: a high throughput web application for PCR and sequencing**  
283 **primer design**. *BMC Bioinformatics* 2008, **9:253-2105-9-253**.
- 284 21. Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen  
285 SG: **Primer3--new capabilities and interfaces**. *Nucleic Acids Res* 2012,  
286 **40(15):e115**.
- 287 22. Untergasser A, Nijveen H, Rao X, Bisseling T, Geurts R, Leunissen JA:  
288 **Primer3Plus, an enhanced web interface to Primer3**. *Nucleic Acids Res* 2007,  
289 **35(Web Server issue):W71-4**.
- 290 23. Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H,  
291 Remmert M, Soding J, Thompson JD, Higgins DG: **Fast, scalable generation of**  
292 **high-quality protein multiple sequence alignments using Clustal Omega**. *Mol*  
293 *Syst Biol* 2011, **7:539**.
- 294 24. Lassmann T, Sonnhammer EL: **Kalign--an accurate and fast multiple**  
295 **sequence alignment algorithm**. *BMC Bioinformatics* 2005, **6:298**.

- 296 25. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy and**  
297 **high throughput**. Nucleic Acids Res 2004, **32**(5):1792-1797.
- 298 26. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced**  
299 **time and space complexity**. BMC Bioinformatics 2004, **5**:113.
- 300 27. Yang Z: **PAML: a program package for phylogenetic analysis by maximum**  
301 **likelihood**. Comput Appl Biosci 1997, **13**(5):555-556.
- 302 28. Suyama M, Torrents D, Bork P: **PAL2NAL: robust conversion of protein**  
303 **sequence alignments into the corresponding codon alignments**. Nucleic Acids  
304 Res 2006, **34**(Web Server issue):W609-12.
- 305 29. Massingham T, Goldman N: **Detecting amino acid sites under positive**  
306 **selection and purifying selection**. Genetics 2005, **169**(3):1753-1762.
- 307 30. Darriba D, Taboada GL, Doallo R, Posada D: **jModelTest 2: more models, new**  
308 **heuristics and parallel computing**. Nat Methods 2012, **9**(8):772.
- 309 31. Darriba D, Taboada GL, Doallo R, Posada D: **ProtTest 3: fast selection of**  
310 **best-fit models of protein evolution**. Bioinformatics 2011, **27**(8):1164-1165.
- 311 32. Guindon S, Gascuel O: **A simple, fast, and accurate algorithm to estimate**  
312 **large phylogenies by maximum likelihood**. Syst Biol 2003, **52**(5):696-704.
- 313 33. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O: **New**  
314 **algorithms and methods to estimate maximum-likelihood phylogenies:**  
315 **assessing the performance of PhyML 3.0**. Syst Biol 2010, **59**(3):307-321.
- 316 34. Stajich JE, Block D, Boulez K, Brenner SE, Chervitz SA, Dagdigian C, Fuellen G,  
317 Gilbert JG, Korf I, Lapp H, Lehvaslaiho H, Matsalla C, Mungall CJ, Osborne BI,  
318 Pocock MR, Schattner P, Senger M, Stein LD, Stupka E, Wilkinson MD, Birney E:  
319 **The Bioperl toolkit: Perl modules for the life sciences**. Genome Res 2002,  
320 **12**(10):1611-1618.
- 321 35. Kim J, Levy E, Ferbrache A, Stepanowsky P, Farcas C, Wang S, Brunner S,  
322 Bath T, Wu Y, Ohno-Machado L: **MAGI: a Node.js web service for fast microRNA-**  
323 **Seq analysis in a GPU infrastructure**. Bioinformatics 2014, **30**(19):2826-2827.
- 324 36. Page M, MacLean D, Schudoma C: **blastjs: a BLAST+ wrapper for Node.js**.  
325 BMC Res Notes 2016, **9**:130-016-1938-1.
- 326 37. Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, Goodstein DM,  
327 Elsik CG, Lewis SE, Stein L, Holmes IH: **JBrowse: a dynamic web platform for**  
328 **genome visualization and analysis**. Genome Biol 2016, **17**:66-016-0924-1.
- 329  
330

331

## 332 **FIGURE LEGENDS**

333

### 334 **Figure 1. Workflow for Phylo-Node.**

335 Phylo-Node is organized into a workflow of connected modules and driver scripts. In  
336 order to interface with a phylogenetic tool, the base wrapper module is invoked to  
337 process command-line requests that are then passed into the software specific  
338 module. The input for the specific software can be passed into the base wrapper  
339 from: a) sample input directory; b) from a folder specified by the user; or c) by using  
340 the sequence retrieval module which is contained within the 'Sequence' directory.  
341 The 'Pipes' folder contains a module for easy piping of data between applications in  
342 Phylo-Node. Binaries can be downloaded using the 'get\_executable' module from  
343 within the 'Download' folder to deploy modules within the 'Run' directory or to provide  
344 applications to a web server from within the 'Server' directory.

345

### 346 **Figure 2. Graphical overview of Phylo-Node applications.**

347 Phylo-Node provides a toolkit for interacting with various applications including: the  
348 sequence alignment software Clustal Omega [23], K-align [24], and MUSCLE [25,  
349 26]; the primer design software, Primer3 [20-22]; software for estimating selection:  
350 Codeml [27], PAL2NAL [28], and Slr [29]; software for determining the best-fit  
351 models of evolution: jModelTest2 [30] and ProtTest3 [31]; and also the phylogeny  
352 reconstruction software, PhyML [32, 33]. Phylo-Node also enables the user to  
353 retrieve sequences remotely from the NCBI database using Entrez Programming  
354 Utilities, and create a web server to deploy applications such as JBrowse [37]. A key  
355 feature of Phylo-Node is interoperability between languages and other Node

356 modules, which can be easily leveraged to form stable and scalable pipelines. This  
357 concept of interoperation and inheritance is highlighted by the brown cog at the  
358 bottom of Figure 2 that represents the potential to integrate any other module(s)  
359 [*require('./module');*] with Phylo-Node.

Figure 1

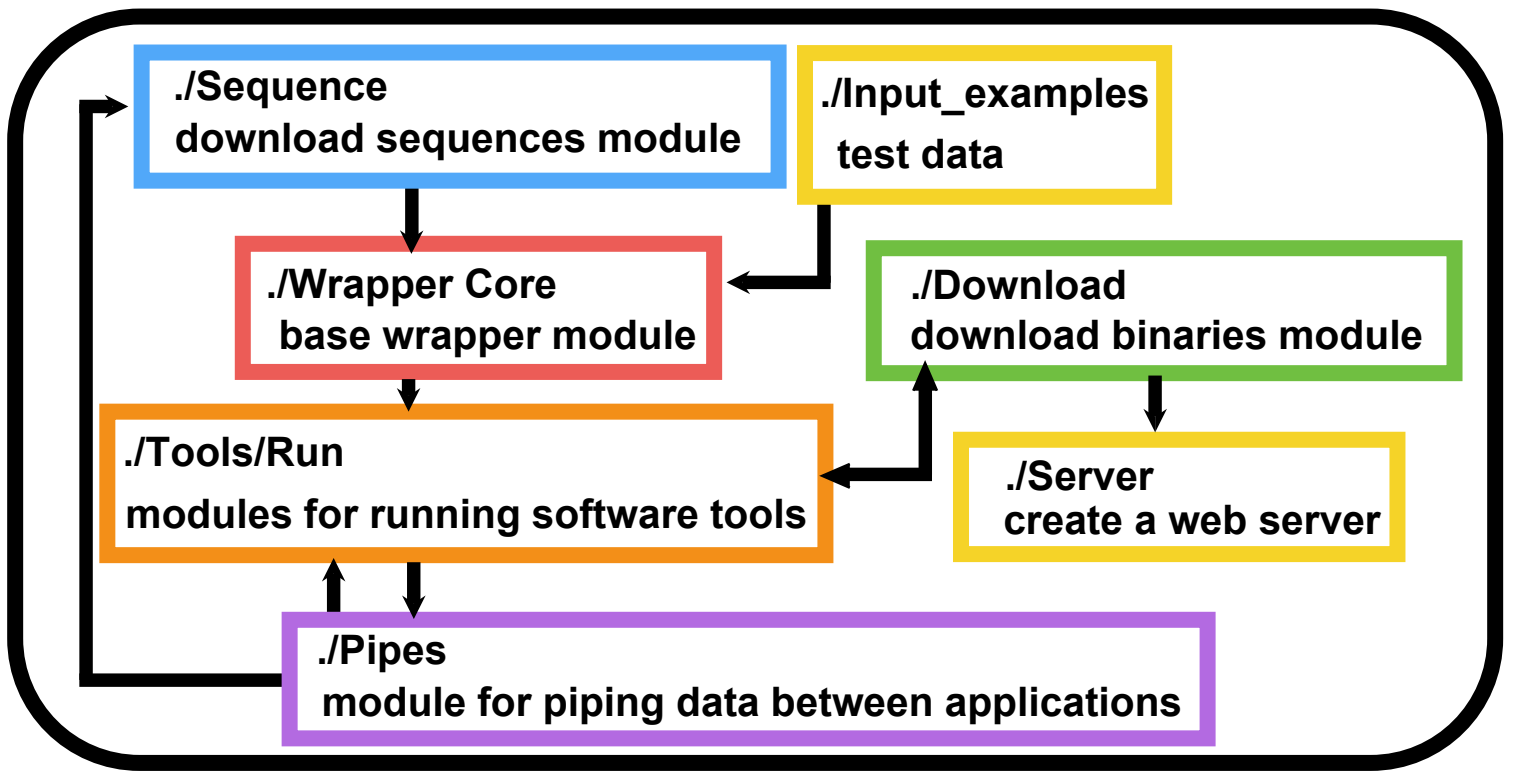


Figure 2

