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

2

3

4    Damien M. O'Halloran[1,2*]

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

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 provides simple integration and interoperation with other Node modules to develop

39 workflows across multiple components and languages using Node inheritance

40 patterns and a customized piping module to support the production of diverse

41 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

**BACKGROUND**

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

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

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

3

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-

91    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

4

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  'phyml' object is 'phyml.getphyml()', which invokes the 'getphyml' method on the

111  'phyml' 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

5

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; jModelTest2 [27] and ProtTest3 [28] to determine the best-fit model of

140    evolution, and PhyML [29, 30] for phylogeny reconstruction. The PhyML executable

141    is also employed by jModelTest2 and ProtTest3. Primer3 is the most popular

142    software for primer design, and takes a very lengthy list of input variables to optimize

143    primer selection. Clustal Omega, K-align, and MUSCLE are very fast and accurate

144    multiple sequence alignment tools that are commonly used to build robust DNA,

145    RNA, or protein alignments. PhyML is a popular program for building phylogenies

146    using maximum likelihood, and ProtTest3 determines the best-fit model of evolution

147    for protein sequences across 120 different potential models, while jModelTest2

148    determines best-fit models of nucleotide substitution from DNA sequence

149    alignments. Together, Phylo-Node provides a novel toolkit that allows the user to go

150    from raw sequence to phylogeny using Node.

151    Phylo-Node is highly scalable and customizable, and was inspired by projects

152    such as BioPerl [31] which provides very diverse tools that include Perl modules for

153    many bioinformatic tasks and also parsers and wrappers for diverse sequence

154    formats and applications. BioPerl's open source structure and architecture allows

155    users to plug new modules into BioPerl pipelines to design new applications. Node.js

156    implements prototypal inheritance as per JavaScript but also provides access to the

157    'module.exports' object which permits easy portability between the Phylo-Node

158    toolkit and any other modules, and also interoperation between different languages

159    by using the 'child.process.exec' process. Therefore, Phylo-Node can be integrated

160    with existing Node.js bioinformatics tools [32, 33] or software written in other

161    languages. For example, both jModelTest2 and ProtTest3 require a Java runtime

162    environment (http://www.oracle.com/technetwork/java/javase/downloads/jre8-

163    downloads-2133155.html), and by using 'require' to import each module, the user

164    can execute the analysis of jModelTest2 and ProtTest3. The 'prottest' and

165    'jmodeltest2' modules and driver scripts (index.js) are contained within the respective

166    'Prottest3' and 'jModelTest2' directories and sample input is provided in the

167    'COX2_PF0016' sub-directory of the 'Input_examples' folder for ProtTest3 and the

168    sample FASTA file aP6.fas (also contained within the 'Input_examples' folder) can

169    be used for testing jModelTest2.

170    To further facilitate the ease of interoperation between various applications and

171    components, the Phylo-Node package also contains a module called 'phylo-

172    node_pipes' inside the 'Pipes' directory. The 'phylo-node_pipes' module allows the

173    user to easily pipe data between different applications by requiring the

174    'child_process' module which provides the ability to spawn child processes. Through

175    'phylo-node_pipes', the user can chain commands together that will be executed in

7

176    sequence to build consistent, and extensive pipelines. The 'Pipes' directory contains

177    sample driver scripts for using the 'phylo-node_pipes' module.

178

179    **CONCLUSIONS**

180    In conclusion, Phylo-Node is a novel package that leverages the speed of Node.js to

181    provide a robust and efficient toolkit for researchers conducting molecular

182    phylogenetics. Phylo-Node can be easily employed to develop complex but

183    consistent workflows, and integrated with existing bioinformatics tools using the

184    Node.js codebase.

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

**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

**ACKNOWLEDGMENTS**

216

**AUTHOR CONTRIBUTIONS**

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

219 manuscript.

220

**COMPETING INTERESTS**

222 The author declares no competing interests.

223

224

225

9

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: **HeatMapViewer:**
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.

11

296   25. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy and high throughput.** Nucleic Acids Res 2004, **32**(5):1792-1797.

298   26. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced time and space complexity.** BMC Bioinformatics 2004, **5**:113.

300   27. Darriba D, Taboada GL, Doallo R, Posada D: **jModelTest 2: more models, new heuristics and parallel computing.** Nat Methods 2012, **9**(8):772.

302   28. Darriba D, Taboada GL, Doallo R, Posada D: **ProtTest 3: fast selection of best-fit models of protein evolution.** Bioinformatics 2011, **27**(8):1164-1165.

304   29. Guindon S, Gascuel O: **A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood.** Syst Biol 2003, **52**(5):696-704.

306   30. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O: **New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0.** Syst Biol 2010, **59**(3):307-321.

309   31. Stajich JE, Block D, Boulez K, Brenner SE, Chervitz SA, Dagdigian C, Fuellen G, Gilbert JG, Korf I, Lapp H, Lehvaslaiho H, Matsalla C, Mungall CJ, Osborne BI, Pocock MR, Schattner P, Senger M, Stein LD, Stupka E, Wilkinson MD, Birney E: **The Bioperl toolkit: Perl modules for the life sciences.** Genome Res 2002, **12**(10):1611-1618.

314   32. Kim J, Levy E, Ferbrache A, Stepanowsky P, Farcas C, Wang S, Brunner S, Bath T, Wu Y, Ohno-Machado L: **MAGI: a Node.js web service for fast microRNA-Seq analysis in a GPU infrastructure.** Bioinformatics 2014, **30**(19):2826-2827.

317   33. Page M, MacLean D, Schudoma C: **blastjs: a BLAST+ wrapper for Node.js.** BMC Res Notes 2016, **9**:130-016-1938-1.

319
320

321

322

323

324

325

326

327

12

328    **FIGURE LEGENDS**

329    **Figure 1. Workflow for Phylo-Node.**

330    Phylo-Node is organized into a workflow of connected modules and driver scripts. In

331    order to interface with a phylogenetic tool, the base wrapper module is invoked to

332    process command-line requests that are then passed into the software specific

333    module. The input for the specific software can be passed into the base wrapper

334    from: a) sample input directory; b) from a folder specified by the user; or c) by using

335    the sequence retrieval module which is contained within the 'Sequence' directory.

336    The 'Pipes' folder contains a module for easy piping of data between applications in

337    Phylo-Node. Binaries can be downloaded using the 'get_executable' module from

338    within the 'Download' folder.

339

340    **Figure 2. Graphical overview of Phylo-Node applications within the *./Tools/Run***

341    **directory.**

342    Phylo-Node provides a toolkit for interacting with various applications including: the

343    sequence alignment software Clustal Omega [23], K-align [24], and MUSCLE [25,

344    26]; the primer design software, Primer3 [20-22]; software for determining the best-fit

345    models of evolution: jModelTest2 [27] and ProtTest3 [28]; and also the phylogeny

346    reconstruction software, PhyML [29, 30]. Phylo-Node also enables the user to

347    retrieve sequences remotely from the NCBI database using Entrez Programming

348    Utilities. A key feature of Phylo-Node is interoperability between languages and other

349    Node modules, which can be easily leveraged to form stable and scalable pipelines.

350    This concept of interoperation and inheritance is highlighted by the brown cog at the

351    bottom of Figure 2 that represents the potential to integrate any other module(s)
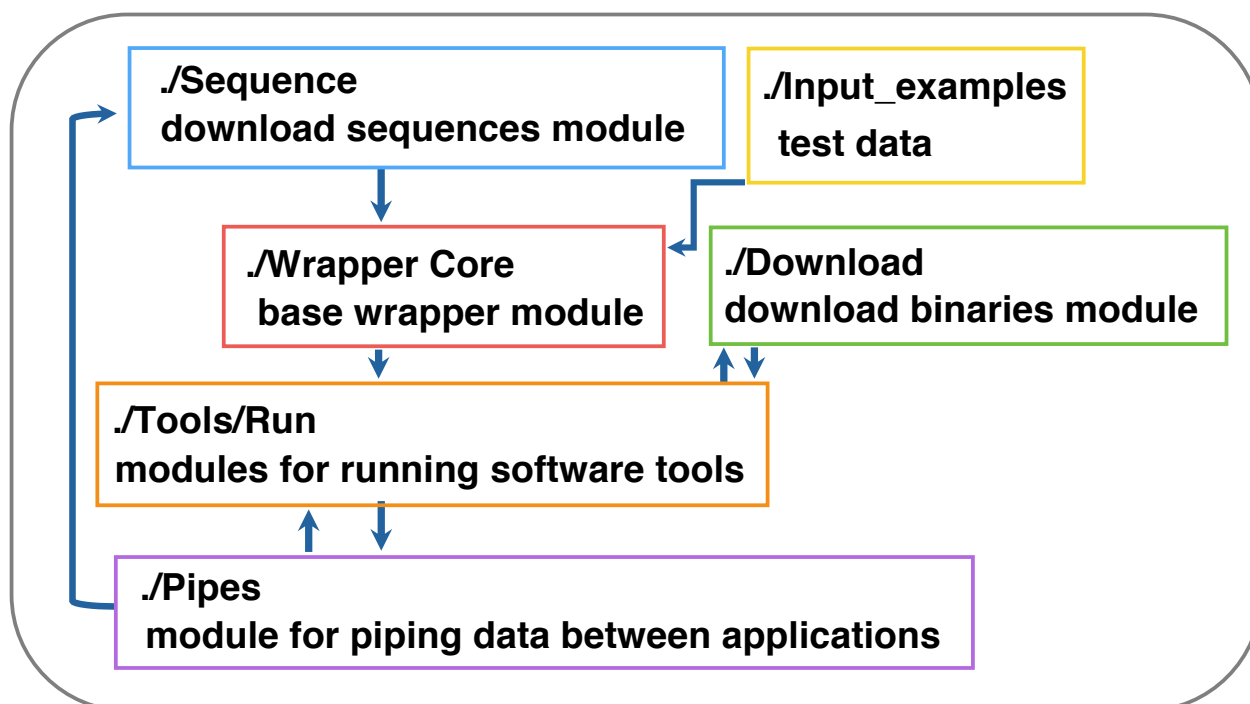
352    [*require('./module');*] with Phylo-Node.

13

**Figure1**

**Figure 2**