

1 COLT-Viz: Interactive Visualization of Antibody Lineage Trees

2 Chenfeng He^{1,*}, Ben S. Wendel^{2,*}, Jun Xiao³, Keke Chen⁴, Ning Jiang^{1,5,\$}

3

4 Authors and Affiliations:

5 ¹Department of Biomedical engineering, Cockrell School of Engineering, The University of
6 Texas at Austin, Austin, TX 78712, USA

7 ²McKetta Department of Chemical Engineering, Cockrell School of Engineering, The University
8 of Texas at Austin, Austin, TX 78712, USA

9 ³ImmuDX LLC, Austin, TX 78750, USA

10 ⁴Department of Computer Science and Engineering, Wright State University, Dayton, OH 45431,
11 USA

12 ⁵Institute for Cellular and Molecular Biology, University of Texas at Austin, Austin, TX 78712,
13 USA.

14

15

16

17 *These authors contributed equally.

18

19 \$Correspondence should be addressed to:

20 Ning Jiang, Ph.D.

21 Email: jiang@austin.utexas.edu

22 Phone: 512-471-4860

23

24

Interactive antibody lineage structure visualization

25 **Abstract**

26 Many tools have been developed to visualize phylogenetic trees, which is a traditional technique
27 for evolutionary tree analysis. However, due to the unique characteristics of antibody lineage
28 trees, the phylogenetic method cannot adequately construct proper tree structures for antibody
29 lineages, and many other tools have been developed to address this problem. However, there still
30 lacks of an adequate tool to visualize the resulted antibody lineage structures that are more
31 complicated than phylogenetic trees. In addition, high-throughput sequencing-based antibody
32 repertoire profiling enables the counting of the number of transcripts associated with individual
33 antibody sequences, thus more dimensions need to be encoded in the tree structure visualization.
34 Further, users may wish to manually adjust the tree structure for a special context. When doing
35 so, they may wish to maintain some biological constraints that are applicable in antibody lineage
36 tree structure, such as isotype switching constraints or different sampling constraints. Here, we
37 report an interactive visualization tool (COLT-Viz) designed to display the number of RNA
38 copies, number of somatic hypermutations, and sample collection time associated with each
39 antibody sequence as well as the distance to neighboring sequences for each antibody sequence
40 in the lineage. COLT-Viz also allows users to interactively visualize and edit antibody lineage
41 structures while giving users the option to automatically check biological constraints on the
42 edited structures to ensure accuracy. COLT-Viz takes JSON text format as input files and can
43 easily be used to visualize networks with or without the biological constraints. We believe the
44 amount of information that can be displayed for complex antibody lineages, the interactive
45 interface, and the option of checking for biological constraints make COLT-Viz a versatile tool
46 for antibody lineage tree visualization that will guide further biological discoveries.

47

48

49 **Key words: bioinformatics software, antibody lineage structure, biological constraints,**
50 **interactive visualization, intuitive user interface**

Interactive antibody lineage structure visualization

51 **1. Introduction**

52 Lineage analysis is frequently used by researchers to analyze the evolution of different species of
53 organisms or to trace cellular development within an organism. Specifically in the immune
54 system, lineage analysis has been adopted to trace the mutational evolution of antibodies (1–3)
55 and understand selection pressure (4,5).

56 Typically, an antibody lineage is represented as a tree structure, which describes the
57 possible development path (i.e., edges in a tree) among all unique antibody sequences (i.e., nodes
58 in a tree). Traditionally, researchers used phylogenetic trees for antibody lineages; however, a
59 phylogenetic tree is a binary tree with all nodes presented as leaf nodes, which is not suitable for
60 antibody lineages (3). Much work have been done to construct tree structures specific to
61 antibody lineages (3,6,7); however, there is still a need for a tool to efficiently visualize these
62 trees. The visualization of antibody lineage trees has a few unique needs: (a) there are multiple
63 features specific to antibody sequences (e.g., number of transcripts, number of somatic
64 hypermutations, isotype, etc), integrating all these properties together for tree visualization will
65 facilitate lineage analysis; (b) interactively visualize the antibody structure by showing antibody
66 sequence information and associated properties when users click on a node of interest. This
67 allows user to have access to detailed information of nodes (antibody sequences) without losing
68 the global view of a complex antibody lineage structure; and (c) interactively modify antibody
69 lineage structure. Oftentimes, multiple valid lineage tree structures exist for an antibody lineage.
70 Based on special considerations, researchers may want to manually alter their tree structures by
71 changing or adding connections between nodes. When doing so, they may want to make sure that
72 these modified connections do not violate certain biological constraints. These constraints can
73 either be biological constraints (8) or constraints determined by experiment design (9).

74 Although many network visualization tools with powerful functions have been developed
75 (e.g., Cytoscape (10), Graphviz (11)), they were not designed for antibody lineage structure
76 visualization and thus lack the consideration of biological properties and constraints. To our
77 knowledge, there is no tool specifically designed to visualize antibody lineage tree structures
78 with the functionalities mentioned above. For example, LymAnalyzer (12) incorporates FigTree
79 (13) to visualize antibody lineage tree structures. However, FigTree is only designed for
80 phylogenetic tree which is not suitable for antibody lineage structure and does not have an

Interactive antibody lineage structure visualization

81 interactive user interface. Another program, Change-O (6), modifies phylogenetic tree analysis
82 for antibody lineage application, but its visualization function can only generate non-interactive
83 figures. In addition, none of these tools can easily be used to visualize the multiple features of
84 antibody lineages, which are unique to antibody lineage tree visualization.

85 Bearing the drawbacks of current visualization tools in mind, we developed COLT-Viz
86 for antibody lineage tree structure visualization. COLT-Viz enables researchers to interactively
87 explore lineage trees with multiple antibody sequence properties encoded and automatically
88 checked when the tree structure is manually altered. At the same time, researchers have the
89 option to override error messages and proceed with the desired changes. This tool enables users
90 to discover interesting patterns of antibody evolution process and select interesting nodes
91 (antibody sequences) for further evaluation.

92 Specifically, COLT-Viz projects antibody sampling time, total mutations for a sequence,
93 mutation distances between two neighboring sequences, and antibody RNA copies of unique
94 sequences onto different elements of each node in the lineage visualization, which allows users
95 to intuitively understand the lineage evolution over time, mutation distances from a germline
96 sequence, mutation distances between neighboring nodes, and the clonal size associated with
97 each unique antibody sequence. Furthermore, users can fine tune the lineage by removing edges
98 and reconnecting nodes with all the constraints automatically checked by the system with an
99 option to either override the error message or make other alternative changes. When users edit
100 the connections, COLT-Viz checks two types of constraints: (a) isotype constraint: as described
101 in previous studies (8,14–16), antibody isotype switch can only proceed according to the order of
102 IGH constant genes located on the chromosome (17), thus users may want to maintain this
103 constraint while editing the lineage; (b) time constraint: as described in Schramm et al. (9),
104 longitudinal information offers another dimension in some biological experiments (18).
105 Although in many cases, sequences from an earlier time point may not necessarily be ancestors
106 of sequences from later time point, in some cases, such as transplantation, where donor samples
107 and recipient samples are experimentally separated, this constraint is critical. Thus, users may
108 want to re-enforce the time constraint when manually changing node connections. Therefore,
109 COLT-Viz also checks if any manual changes on the tree structure comply with the time
110 constraint that sequences from an earlier time point should be ancestors of sequences from a later

Interactive antibody lineage structure visualization

111 time point while give users the option of override this error message and proceed with their
112 intended changes.

113 **2. Implementation**

114 COLT-Viz takes the lineage tree from tree-generating tools (e.g., COLT (7), Change-O
115 (6)) as input and outputs the tree visualization. The input of COLT-Viz contains one node file
116 and one edge file. The node and edge files are encoded in JSON (JavaScript Object Notation)
117 format (<http://json.org/>), which is a lightweight and human-readable text format. JSON expresses
118 data objects in attribute-value pairs, which makes it an ideal data format for visualizing antibody
119 sequences with multiple attributes. Outputs from tree-generating tools/algorithms can be
120 converted to the JSON text format to be compatible with COLT-Viz. Users can also manually
121 edit the JSON text files to change the lineage tree structure. The node file contains a list of nodes,
122 or unique antibody sequences. Each node element includes the node identity, node description,
123 associated sequence abundance (NRNAs, number of RNA copies, and NREADs, number of
124 sequencing reads), antibody isotype (i.e., antibody heavy chain families, IgM, IgD, IgG, IgA, and
125 IgE), timestamp (i.e., sampling time point, for example, when doing research on vaccine,
126 sequences obtained before vaccination are all assigned to a smaller number (e.g. 20) while
127 sequences obtained 7 days post vaccination are all assigned to a larger number (e.g. 40)), and the
128 number of somatic hypermutations from the germline sequence. Each edge in the edge file is
129 defined as a 3-tuple (parent node identity, weight, child node identity), which are also encoded in
130 the JSON format.

131 Users can explore the tree visualization, e.g., checking the detailed information of each
132 node, and fine tune the tree structure by editing and moving the edges. Before making changes,
133 COLT-Viz automatically checks the constraints and displays a warning message if either the
134 timepoint or isotype constraint is violated. However, user can choose to override this warning
135 message and proceed with the move. Once old edges are removed and new edges are added, the
136 system will automatically update the entire graph. Finally, the updated tree structure can be
137 saved back to the edge files in the JSON format for later usage.

138

139

Interactive antibody lineage structure visualization

140 3. Results

141 3.1 Lineage Tree Visualization

142 We visually encode multiple properties of nodes and edges into the lineage tree visualization, so
143 users can understand these properties intuitively. Figure 1 shows an antibody lineage we found in
144 a young African child who experienced acute malaria twice in two consecutive malaria seasons
145 from our previous study (19). There are a total of 23 unique sequences obtained from 4
146 timepoints in this lineage. The height of each node represents the number of RNA copies
147 associated with that unique sequence. Heatmap is used to color nodes – the warmer the color, the
148 greater the number of somatic hypermutations to the germline sequence. The border of each node
149 is colored to represent different timestamps, so users can easily identify the evolution of the
150 antibody lineage over time.

151 The tree layout algorithm uses the well-known force-directed graph drawing algorithm
152 (20). It progressively changes the nodes' positions, as if there are physical forces between the
153 nodes. Imagine that neighboring nodes are connected by springs: the force of the spring will
154 prevent the nodes from drifting too far apart or collapsing too close together. This algorithm
155 ideally places the root antibody sequence towards the center of the visualization and unfolds the
156 tree outward to maximize the use of the display area. Thus the directionality of the tree is in
157 general from center to periphery. Knowing this makes it easier for user to identify the root node,
158 also if user assign the root as a specific ID (e.g., 0, same as in Figure 1), then the root node
159 should be easily identified in the tree. Then, by following the directionality, users can easily
160 understand the development of the whole lineage.

161 Using COLT-Viz, features of the antibody lineages that are not obvious by mining the
162 data can be quickly perceived visually. Using the lineage in Figure 1 as an example, users can
163 quickly identify the root sequence and see that, for this lineage, the sequences at later time points
164 have accumulated more somatic hypermutations compared to sequences collected at earlier time
165 points, and they lie in the outer leaves of the lineage. In addition, the relatively even distribution
166 of node heights indicates that the RNA copies of these unique sequences do not differ widely.
167 These visual features help immunologists quickly glean an overview of the evolution of this
168 antibody lineage over time.

Interactive antibody lineage structure visualization

169

170 **3.2 Tree Editing and Constraint Maintaining**

171 COLT-Viz allows users to alter the network connections by removing and adding edges between
172 nodes. When a new edge is added, the system will check two types of constraints: (a) time
173 constraint, an edge $x \rightarrow y$ is valid only if x 's timestamp is earlier (smaller) than y 's; and (b)
174 isotype constraint, an edge $x \rightarrow y$ is valid only if the isotype is either unchanged or following the
175 class-switching rule: a class-switched sequence cannot switch back to IgM. If either of the
176 constraints is violated, COLT-Viz will display a warning. However, users can override this
177 warning and proceed to add the new edge or heed the warning and cancel the new edge. The
178 constraint checking helps avoid artificial mistakes and maintain interesting time features when
179 users manually modify antibody lineage trees.

180 **3.3 Intuitive graphic visualization user interface**

181 COLT-Viz employs an intuitive graphic visualization user interface (Figure 1). Within the
182 antibody lineage structure displaying window, users can drag on any node to reposition the tree.
183 Users can also click on any node to view the full antibody sequence and annotated information in
184 the sequence display window on the right. In addition, users can change the lineage structure by
185 clicking on a node and specifying the new parent node that they wish to establish a link in the
186 edge editing window.

187 **4. Discussion**

188 COLT-Viz is an interactive tool designed for antibody lineage tree visualization and editing. By
189 using COLT-Viz, researchers can understand lineage trees, validate algorithmic results, and fine-
190 tune the tree structures with additional domain knowledge that automated algorithms cannot
191 capture. Further, COLT-Viz allows researcher to identify nodes (antibody sequences) that bear
192 specific features that are only obvious in graphic settings. Constraints are automatically checked
193 to avoid manual editing errors. Although COLT-Viz is specifically designed for antibody lineage
194 analysis, we believe with small tweaks it can be easily applied to any other type of lineage
195 analysis in biomedical research.

196

197

Interactive antibody lineage structure visualization

198 **Availability of data and software**

199 Project name: COLT-Viz

200 Project home page: <https://github.com/immudx/paper> (sample data is included).

201 Operation systems: Platform independent

202 Programming language: Java

203 Any restrictions to use by non-academics: None

204

205 **Abbreviations**

206 IgM (D,G,A,E): Immunoglobulin M (D,G,A,E)

207 **Funding:** This work was supported by NIH grants R00AG040149 (N.J.) and by the Welch
208 Foundation grant F1785 (N.J.). NJ is a Cancer Prevention and Research Institute of Texas
209 (CPRIT) Scholar and a Damon Runyon-Rachleff Innovator. BW is a recipient of the Thrust 2000
210 - George Sawyer Endowed Graduate Fellowship in Engineering.

211

212 **Author contributions:** NJ conceived the idea and directed study; CH, BW and NJ participated
213 in the algorithm design. KC and JX wrote the computer software. CH, BW contributed to
214 software testing and optimization. CH, KC, JX and NJ contributed to the writing of the
215 manuscript. All of the authors read and approved the final manuscript.

216 **Conflict of Interest Statement:** Ning Jiang is a scientific advisor of ImmuDX LLC.

217 **Acknowledgements:** Not applicable

218

219

220

221

222

Interactive antibody lineage structure visualization

223 **References**

- 224 1. Jiang N, He J, Weinstein J a, Penland L, Sasaki S, He X-S, Dekker CL, Zheng N-Y,
225 Huang M, Sullivan M, et al. Lineage structure of the human antibody repertoire in
226 response to influenza vaccination. *Sci Transl Med* (2013) **5**:171ra19.
227 doi:10.1126/scitranslmed.3004794
- 228 2. Liao H-X, Lynch R, Zhou T, Gao F, Alam SM, Boyd SD, Fire AZ, Roskin KM, Schramm
229 CA, Zhang Z, et al. Co-evolution of a broadly neutralizing HIV-1 antibody and founder
230 virus. *Nature* (2013) **496**:469–76. doi:10.1038/nature12053
- 231 3. Barak M, Zuckerman NS, Edelman H, Unger R, Mehr R. IgTree©: Creating
232 Immunoglobulin variable region gene lineage trees. *J Immunol Methods* (2008) **338**:67–74.
233 doi:10.1016/j.jim.2008.06.006
- 234 4. Yaari G, Vander Heiden JA, Uduman M, Gadala-Maria D, Gupta N, Joel JN, O'Connor
235 KC, Hafler DA, Laserson U, Vigneault F, et al. Models of somatic hypermutation
236 targeting and substitution based on synonymous mutations from high-throughput
237 immunoglobulin sequencing data. *Front Immunol* (2013) **4**:
238 doi:10.3389/fimmu.2013.00358
- 239 5. Hershberg U, Uduman M, Shlomchik MJ, Kleinstein SH. Improved methods for detecting
240 selection by mutation analysis of Ig V region sequences. *Int Immunol* (2008) **20**:683–694.
241 doi:10.1093/intimm/dxn026
- 242 6. Gupta NT, Vander Heiden JA, Uduman M, Gadala-Maria D, Yaari G, Kleinstein SH.
243 Change-O: A toolkit for analyzing large-scale B cell immunoglobulin repertoire
244 sequencing data. *Bioinformatics* (2015) **31**:3356–3358. doi:10.1093/bioinformatics/btv359
- 245 7. Chen K, Sai V, Gogu A, Wu D, Ning J. COLT: Constrained Lineage Tree Generation
246 from Sequence Data. *Proc IEEE Int Conf Bioinforma Biomed 2016* (2016)
- 247 8. Horns F, Vollmers C, Croote D, Mackey SF, Swan GE, Dekker CL, Davis MM, Quake
248 SR. Lineage tracing of human B cells reveals the in vivo landscape of human antibody
249 class switching. *Elife* (2016) **5**: doi:10.7554/eLife.16578.001
- 250 9. Schramm CA, Sheng Z, Zhang Z, Mascola JR, Kwong PD, Shapiro L. SONAR: A high-

Interactive antibody lineage structure visualization

- 251 throughput pipeline for inferring antibody ontogenies from longitudinal sequencing of B
252 cell transcripts. *Front Immunol* (2016) **7**: doi:10.3389/fimmu.2016.00372
- 253 10. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski
254 B, Ideker T. Cytoscape: A software Environment for integrated models of biomolecular
255 interaction networks. *Genome Res* (2003) **13**:2498–2504. doi:10.1101/gr.1239303
- 256 11. Gansner ER, North SC. An open graph visualization system and its applications to
257 software engineering. *Softw Pract Exp* (2000) **30**:1203–1233. doi:10.1002/1097-
258 024X(200009)30:11<1203::AID-SPE338>3.0.CO;2-N
- 259 12. Yu Y, Ceredig R, Seoighe C. LymAnalyzer: A tool for comprehensive analysis of next
260 generation sequencing data of T cell receptors and immunoglobulins. *Nucleic Acids Res*
261 (2015) **44**: doi:10.1093/nar/gkv1016
- 262 13. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti
263 and the BEAST 1.7. *Mol Biol Evol* (2012) **29**:1969–1973. doi:10.1093/molbev/mss075
- 264 14. Iwasato T, Shimizu A, Honjo T, Yamagishi H. Circular DNA is excised by
265 immunoglobulin class switch recombination. *Cell* (1990) **62**:143–9. doi:0092-
266 8674(90)90248-D [pii]
- 267 15. von Schwedler U, Jack HM, Wabl M. Circular DNA is a product of the immunoglobulin
268 class switch rearrangement. *Nature* (1990) **345**:452–456. doi:10.1038/345452a0
- 269 16. Yoshida K, Matsuoka M, Usuda S, Mori A, Ishizaka K, Sakano H. Immunoglobulin
270 switch circular DNA in the mouse infected with *Nippostrongylus brasiliensis*: evidence
271 for successive class switching from mu to epsilon via gamma 1. *Proc Natl Acad Sci U S A*
272 (1990) **87**:7829–7833. doi:10.1073/pnas.87.20.7829
- 273 17. Murphy, K., & Weaver C. *Janeway's immunobiology*. Garland Science
- 274 18. Wu X, Zhang Z, Schramm CA, Joyce MG, Do Kwon Y, Zhou T, Sheng Z, Zhang B,
275 O'Dell S, McKee K, et al. Maturation and diversity of the VRC01-antibody lineage over
276 15 years of chronic HIV-1 infection. *Cell* (2015) **161**:480–485.
277 doi:10.1016/j.cell.2015.03.004

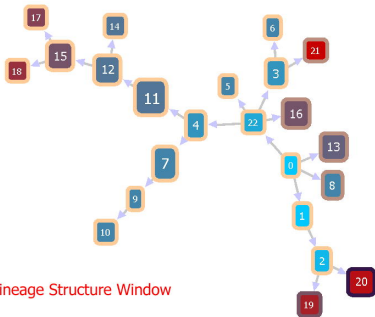
Interactive antibody lineage structure visualization

- 278 19. Wendel BS, He C, Qu M, Wu D, Hernandez SM, Ma K-Y, Liu EW, Xiao J, Crompton PD,
279 Pierce SK, et al. Accurate immune repertoire sequencing reveals malaria infection driven
280 antibody lineage diversification in young children. *Nat Commun* (2017) **8**:531.
281 doi:10.1038/s41467-017-00645-x
- 282 20. Kobourov SG. Spring Embedders and Force Directed Graph Drawing Algorithms. *arXiv*
283 *Prepr arXiv12013011* (2012)1–23.

284

285 **Figures:**

286 **Figure 1. User interface of COLT-Viz.** Within the ‘Lineage structure window’, height of each
287 node represents the RNA copy of each sequence, color of each node represents the mutation with
288 respect to the germline, color of the border of each node represents the time point and the edge
289 length represents the edit distance between neighboring nodes if user specifies the option. By
290 clicking on a node, the information of this node will be displayed in the ‘sequence display
291 window’. User can also edit the connections through ‘edge editing window’. The ‘lineage
292 structure window’ shows an antibody lineage found in a young child who experienced two
293 consecutive malaria infection. 23 sequences were found from 4 time points and visualized using
294 COLT-viz. Four time points are: year 1 pre-malaria infection (light orange node frame for most
295 of the nodes displayed), year 1 acute malaria infection (brown node frame for node 16, 13, and 8),
296 year 2 pre-malaria infection (light purple), and year 2 acute malaria infection (dark purple).
297 Sequence 0 is the root and arrows point to the progeny sequences. Sequences from later time
298 exist in the outer leaves of the lineage, which represents their evolution overtime.



Lineage Structure Window

Clicked Node:

NID: 0
 LABEL: >ITGCCGTGAGTTAT_1_2|ga|1|6|relpbmc
 TIMESTAMP: 0
 ISOTYPE: IGA
 NRNs: 1
 NREADs: 2
 NMutations: 1
 SEQUENCE: ACCTTCACCGGGCTACTATATGCCACTGGGGTGC
 GACAGGGCCCTGGACAAGGGCTTGAGTGGATGGGATGGA
 CAACCCCTAACAGTGGTGGCACAACTATGCACAGAAGTTT
 AGGGCAGGGTCCACCATGACCAGGGGACAGCTCCATCAGCA
 CAGCCTACATGGAGCTGAGCAGGGCTGAGATCTGACGACAC
 GDDCGTGTATTACTGTGGCAGCCGAACTGGGGAGTTACTTT
 GACTACTGGGGCCAGGSAAGCGCTGGTCACCGTCTCCGCA

Sequence display window

Change the clicked node's parent to node

Change

Edge Editing Window