

1 AncesTree: an interactive immunoglobulin lineage tree visualizer

2

3 Mathilde Foglierini^{1,2*}, Leontios Pappas¹, Antonio Lanzavecchia¹, Davide Corti³ and Laurent Perez¹

4

5 ¹Università della Svizzera italiana, Faculty of Biomedical Sciences, Institute for Research in
6 Biomedicine, Via Vincenzo Vela 6, CH-6500 Bellinzona, Switzerland.

7 ²Swiss Institute of Bioinformatics, Lausanne, Switzerland.

8 ³Humabs Biomed SA, Vir Biotechnology, CH-6500 Bellinzona, Switzerland.

9

10 *Corresponding author

11 Mathilde Foglierini (mathilde.perez@irb.usi.ch)

12

13 Abstract

14 High-throughput sequencing of human immunoglobulin genes allows analysis of antibody repertoires
15 and the reconstruction of clonal lineage evolution. Phylip, an algorithm that has been originally
16 developed for applications in ecology and macroevolution, can also be used for the phylogenic
17 reconstruction of antibodies maturation pathway. The study of antibodies (Abs) affinity maturation
18 is of specific interest to understand the generation of Abs with high affinity or broadly neutralizing
19 activities. Phylogenic analysis enables the identification of the key somatic mutations required to
20 achieve optimal antigen binding. To complement Phylip algorithm, we developed AncesTree, a
21 graphic user interface (GUI) that aims to give researchers the opportunity to interactively explore
22 antibodies clonal evolution. AncesTree displays interactive immunoglobulins (Ig) phylogenic tree, Ig
23 related mutations and sequence alignments using additional information coming from specialized
24 antibody tools (such as IMGT®). The GUI is a Java standalone application allowing interaction with
25 Ig-tree that can run under Windows, Linux and Mac OS.

26

27 **Keywords:** Antibodies, Phylogenic tree, Sequence alignment, Immune repertoire exploration

28

29 **Introduction**

30 Development of Next Generation Sequencing (NGS) methodology and its use for high-throughput
31 sequencing of the Adaptive Immune Receptor Repertoire (AIRR-seq) has provided unprecedented
32 molecular insight into the complexity of the humoral adaptive immune response by generating Ig data
33 sets of 100 million to billions of reads. Different computational methods have been developed to
34 exploit and analyze these data (1). Retracing the antigen-driven evolution of Ig repertoires by
35 inferring antibody evolution lineages is a powerful method to understand how vaccines or pathogens
36 shape the humoral immune response (2-5). Indeed, Abs maturation is the result of clonal selection
37 during B cell expansion. A clonal lineage is defined as immunoglobulin sequences originating from
38 the same recombination event occurring between the V, D and J segments (6). Upon B cell receptor
39 (BCR) engagement by a given antigen, somatic hypermutations (SHMs) events will generate a large
40 BCR diversity, leading to antibodies with mutated Ig variable regions, thus forming a specific B-cell
41 lineage that extends from the naive unmutated B-cells, to somatically hypermutated and class
42 switched memory B or plasma-cells (7). Lineage tree building requires a common preprocessing step,
43 where all sequences with identical V, J genes and CDR3 length (with a high CDR3 similarity) are
44 grouped together (8-12). However, there is no consensus as to which phylogenetic method is optimal
45 to infer the ancestral evolutionary relationships among Ig sequences (13, 14). Actually, several
46 methods have been used, such as Levenshtein distance (LD), neighbor joining (NJ), maximum
47 parsimony (MP), maximum likelihood (ML), and Bayesian inference (BEAST) (9, 15-17). DNA
48 Maximum Likelihood program (Dnaml) of the PHYLIP package (18), is a ML method that has been
49 originally developed for applications in ecology. It is also commonly used to infer B cell clonal
50 lineages (19-24). Visualization of the phylogeny is performed using Dendroscope (25, 26). Currently
51 there is no efficient bioinformatics tool allowing an interactive display of phylogenic tree inferred

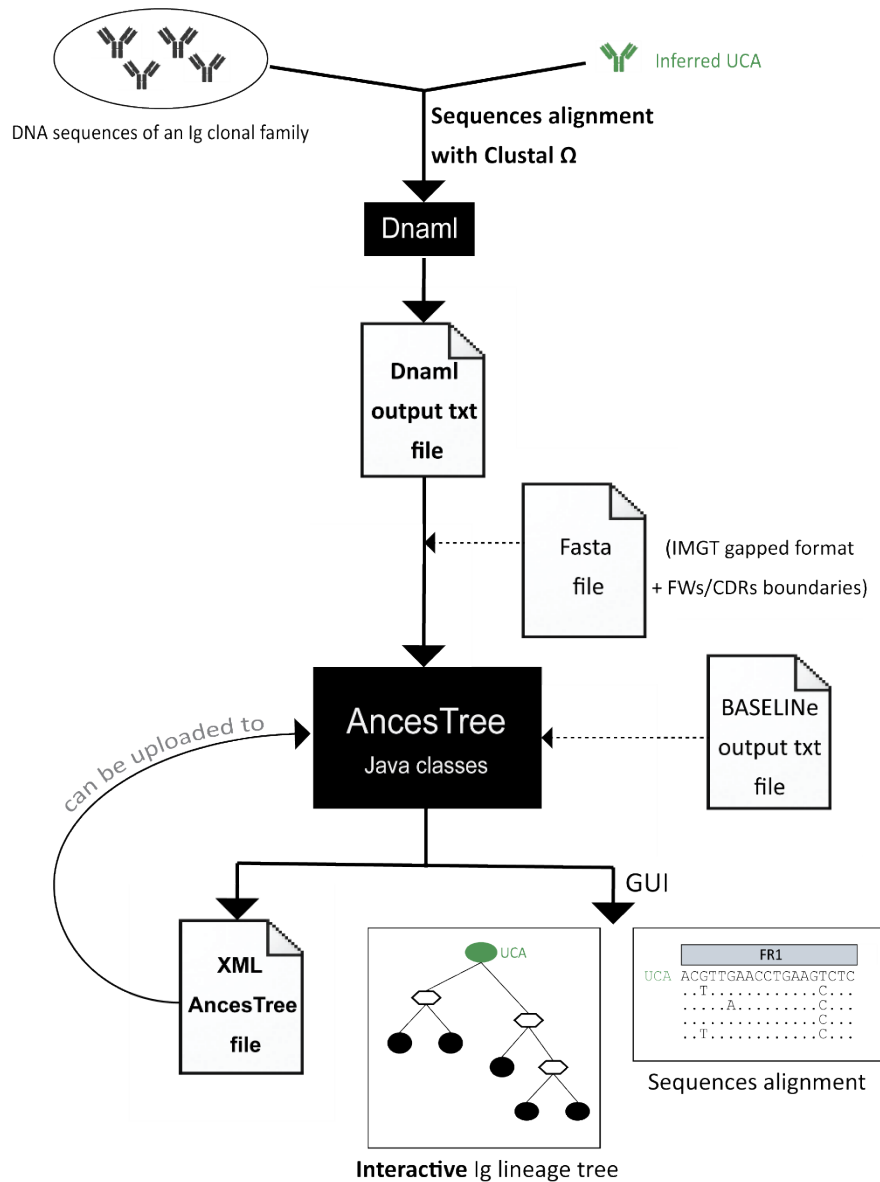
52 from Ig sequences. Here we developed AncesTree, a Dnaml Ig lineage tree visualizer that also
53 integrates information coming from most used antibody bioinformatics tools: IMGT® (27), Kabat
54 numbering (28) and BASELINE (29). AncesTree enables users to interact with the tree generated by
55 Dnaml via the GUI, which is a standalone application that is platform independent and only need
56 JAVA JRE 12 or higher as prerequisite software installed.

57

58 **Design and implementation**

59 The AncesTree workflow is presented in **Fig 1**, it consists of three different main steps: Input,
60 Processing and Outputs.

61



62

63 **Figure 1. AncesTree workflow.** DNA sequences of the variable region of an Ig clonal family of
64 interest are aligned with Clustal Ω . Dnaml processes the sequences and generates a phylogenetic tree.
65 The Dnaml output text file is then used as input for AncesTree. A fasta file with the UCA in gapped
66 IMGT format can be provided (with the FWs and CDRs nucleotide positions in the fasta identifier).
67 AncesTree processes the different inputs and reconstructs the phylogenetic tree with all information
68 related to Ig. BASELINE can be processed separately and its output saved in a text file and then
69 uploaded into AncesTree. The tree is displayed in a GUI and an Extensible Markup Language (XML)
70 file is produced (that could be used as direct input into AncesTree). Dashed arrows indicate optional
71 features.

72

73 **Input**

74 The required input for AncesTree usage is the output text file generated by Dnaml. Optionally, a fasta
75 file with data obtained from IMGT® can also be used to have full AncesTree features.

76 A clonal family is composed of heavy (or light) V(D)J sequences and their related unmutated common
77 ancestor (UCA). The UCA can be inferred with Antigen Receptor Probabilistic Parser (ARPP) UA
78 Inference software (30) or Cloanalyst (31). Then, sequences are aligned with Clustal Ω (32) and the
79 generated file in PHYLIP format can be provided for Dnaml. Dnaml is launched with the following
80 settings: ‘O’ for the outgroup root with the number corresponding to the UCA position provided in
81 the PHYLIP input text file and ‘5’ to reconstruct hypothetical sequences. The generated ‘outfile’ text
82 file can be used as input for AncesTree.

83 To visualize the different frameworks (FW) and complementary-determining (CDR) regions that
84 composed the Ig variable region, a fasta file can be uploaded. The user provides a fasta file containing
85 the following information: the UCA V(D)J sequence in IMGT format including gaps, and the end
86 positions of each region included in the fasta identifier (separated by a space). This information is
87 easily retrieved using IMGT/V-QUEST (33) with the UCA nucleotide sequence as input.

88

89 **Processing**

90 AncesTree parses the Dnaml output file, and does not required a tree in Newick format. Indeed, the
91 relationship between the different nodes of the tree is already stored, in addition to the sequence of
92 each node, in the Dnaml output text file. The theoretical intermediate reconstructed sequences are
93 renamed branch points (BPs) and in the case of ambiguous nucleotide notation (IUPAC
94 nomenclature), AncesTree selects the nucleotide with the highest probability based on the Ig
95 sequences retrieved after this BP. AncesTree has the ability to collapse a node if the sequences are
96 identical, for example in the case of a theoretical BP correspond to an existing Ig. Moreover,

97 AncesTree will also draw different nodes clustered together in the case of identical Ig sequences, thus
98 providing a clear topology view of the tree.

99

100 **Outputs**

101 After running AncesTree, a sub-folder is automatically created in the ‘output’ folder, it uses the name
102 of the Dnaml output file. The folder will contain all produced files such as a XML file that can be
103 used for direct loading into the GUI.

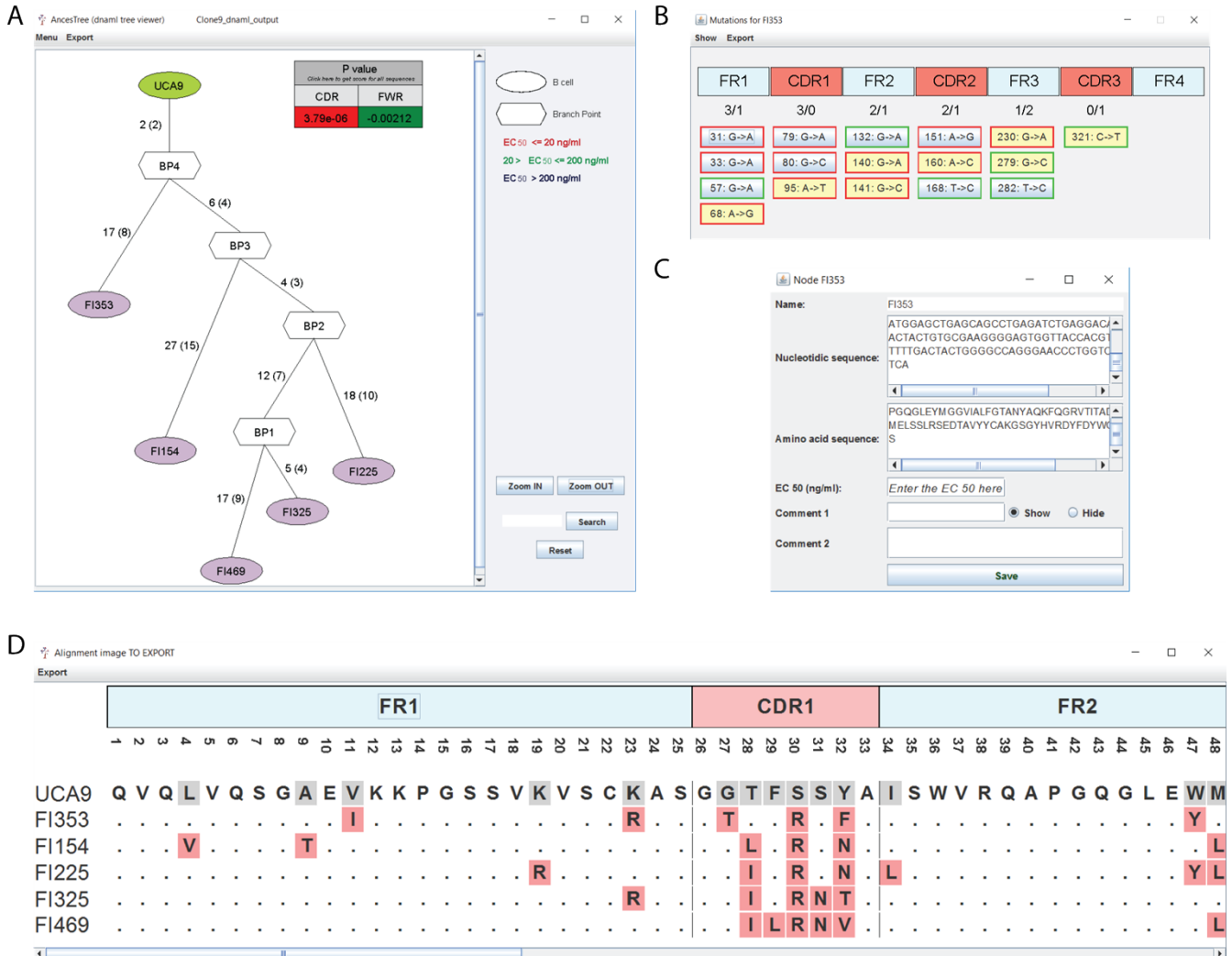
104

105 AncesTree displays the processed tree in the main panel of the GUI (**Fig 2A**). The number of
106 nucleotide and amino acid mutations written on the edge between each node/sequence (with amino
107 acid mutations shown in parenthesis) is clickable and enables the opening of a new window frame
108 that displays the detailed location of each mutation (**Fig 2B**). Of note, the color of the box around
109 each mutated codon indicates whether the mutation is replacement (R) in red or silent (S) in green.
110 This information is also available as R/S numbers under each region. The user can view the amino
111 acid mutations, and have access by default to the Kabat numbering of the related amino acid
112 position (without internet access, AncesTree will use the absolute position). To obtain the
113 nucleotide or protein sequence of a node, the user can click on it (**Fig 2C**). The user has also the
114 possibility to enter the EC50 for the specified Ig. The sequence alignments (DNA or protein) are
115 also accessible in a new frame via the ‘Menu’ button on the top (**Fig 2D**). The alignment view is
116 customizable: the sequences can be selected or deselected, as well as the different positions or
117 regions. Different color modes can be chosen.

118 If the user is interested in a BASELINE analysis of its clonal family of interest, and if the optional
119 input fasta file (with the UCA VDJ sequence including gaps) was provided, AncesTree generates
120 automatically the fasta input file needed for this software (<http://selection.med.yale.edu/baseline/>).

121 Once BASELINE is processed, its output can be loaded into AncesTree to have a nice graphic view

122 of antigen-driven selection occurring for this particular clonal family. All generated graph can be
 123 exported in PNG or EPS format, the alignment can also be exported in a Tab-separated Values (TSV)
 124 file.



125
 126 **Figure 2. Snapshot of Ancestree GUI.** (A) The tree generated by Dnaml is displayed in the main
 127 panel. The BASELINE analysis for the clonal family is displayed in the right upper corner. (B) The
 128 mutations between two nodes can be displayed in a separate window and they are positioned using
 129 IMGT® sequence annotation. (C) The user can have access to each specific node to obtain the related
 130 sequences (DNA or protein) and add comments. (D) An alignment is generated with the UCA
 131 appearing in the first lane, and a ruler indicates the different regions that compose an Ig sequence.

132

133 **Results**

134 To demonstrate the utility of Ancestree we analyzed a case study by performing the analysis of an
135 Ig lineage tree targeting the fusion protein (F) of the Respiratory Syncytial Virus (RSV). RSV is an
136 enveloped RNA virus belonging to the recently defined *Pneumoviridae* family (34). Infection of
137 healthy adults by RSV typically results in mild respiratory symptoms. However, viral infection of
138 infants and older adults, accounts for a substantial hospitalization burden in both age groups (35).
139 Indeed, RSV infection is the second cause of infant mortality worldwide after malaria (36).
140 Understanding the immunological basis for the development of potent neutralizing antibodies is a key
141 step for the development of an effective vaccine for RSV.

142

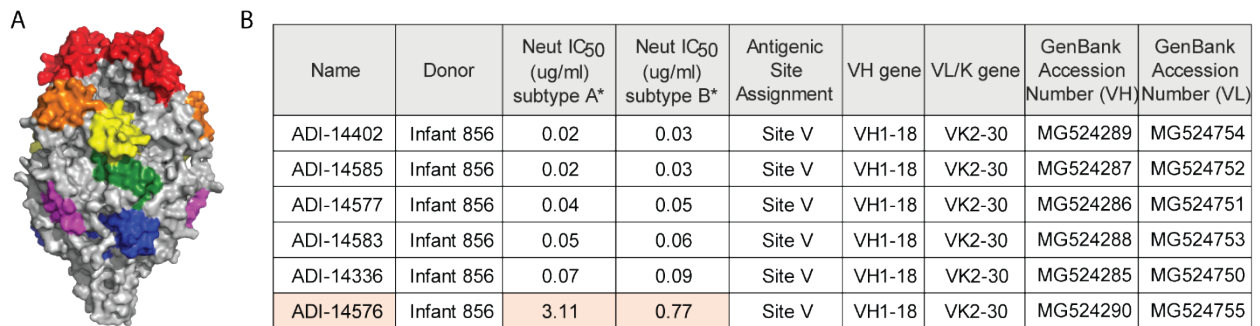
143 **Case study: Exploration of Ig lineage targeting the Fusion protein of the Respiratory Syncytial** 144 **Virus (F-RSV)**

145 To demonstrate the practical use of Ancestree, we re-analyzed an Ig dataset generated post infection
146 by Respiratory Syncytial Viral infection (HRSV). The dataset was collected by isolating antibodies
147 direct against the RSV F protein, a class I fusion protein mediating viral entry into host cells (37).
148 The Ig sequences were clustered by grouping antibodies sharing the same VH and VL gene usage,
149 HCDR3 length and identity (at least 85% for HCDR3). Among the clusters generated, we chose Igs
150 targeting the antigenic site V of RSV F located near amino acid 447 between the α 3 helix and β 3/ β 4
151 hairpin of F-RSV in prefusion (**Fig 3A**). About 70% of the mAbs targeting this site use the same VH
152 and VL germline pair (VH1–18 and VK2–30) (37-39). We identified an Ig family of interest
153 containing potent neutralizers targeting site V with one outlier, the mAb ADI-14576, being less potent
154 and with a 10-fold decrease in binding affinity (**Fig 3B**). We used Dnaml to generate VH sequences
155 phylogenic tree and launched Ancestree to analyze and interact with the produced phylogenic tree
156 (**Fig 4A**). The EC50 (ng/ml) related to the neutralization assay against RSV subtype A are reported
157 in each node (of note, EC50 against subtype B are in the same range for each Ig). Surprisingly, a
158 common mutation 92:G->A (kabat position 31: S ->N) is shared between all the Igs, except for ADI-

159 14576 that does not share this mutation. The alignment of the Ig protein sequences highlights clearly
160 this shared mutation (**Fig 4B**). A result suggesting that ADI-14576 underwent less affinity maturation
161 and therefore diverges from all the other family members. Interestingly, the 31:S->N mutation is
162 located in the HCDR1 and asparagine residues are often involved in protein binding sites. It is
163 tempting to speculate that the Serine to Asparagine substitution is in part responsible for the higher
164 potency and binding titer of the antibodies.

165

166



167

168

169 **Figure 3. Clonal family against RSV-F protein antigenic site V.** (A) Shown is the prefusion
170 conformation of RSV F trimer. The antigenic sites are colored, site Ø (red), I (blue), II (yellow), III
171 (green), IV (purple) and V (orange). Representation was done using PDB ID 4mmu (40) and
172 prepared using PyMOL software (The PyMOL Molecular Graphics System, Version 4.5
173 Schrödinger, LLC). (B) Table showing the different characteristic of a mAbs clonal family isolated
174 from an infant (≥ 6 months) after RSV infection. The Igs neutralization titers are shown as well as
175 their related Germline annotations. ADI-14576 is highlighted because of its lower neutralization
176 value in comparison to the other mAbs of the same clonal family.

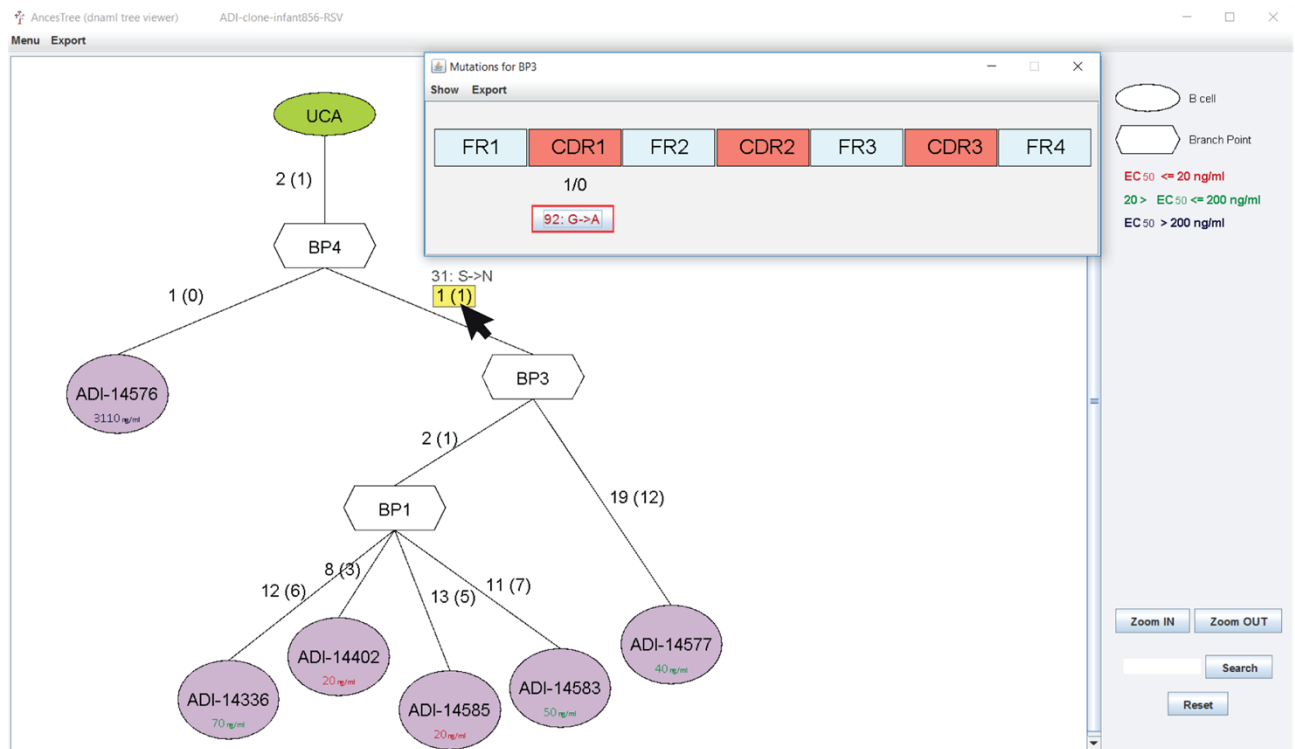
177

178

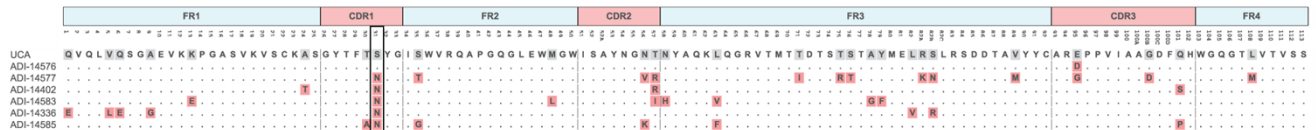
179

180

A



B



181

182 **Figure 4. Phylogenetic analysis of the VH chain of a clonal family RSV-F specific.** (A) Phylogenetic
 183 tree displayed in AncestryTree where the user clicked on the shared mutation for all Igs below BP3 node
 184 (31: S->N). (B) Protein alignment of the different Ig sequences, the mutation 31: S->N is boxed.

185

186 Concluding remarks

187 To summarize, we developed an intuitive, easy and interactive GUI allowing the visualization and
 188 exploration of antibody clonal evolution. Our application is open access and only needs the file
 189 produced by Dnaml and restricted information specific to antibody sequence analysis.

190

191 Availability

192 Ancestry is open-source software implemented in Java and freely available from
193 <https://bitbucket.org/mathildefog/ancestry>. Documentation for installation and user tutorial are
194 provided.

195

196 **Authors' contributions**

197 MF developed the application and performed the analyses. LPa, AL, DC and LPe participated in the
198 design of the application. MF and LPe wrote the paper. All authors read and approved the final
199 manuscript.

200

201 **Competing interests**

202 The authors declare that they have no competing interests.

203

204 **Funding**

205 This research did not receive any specific grant from funding agencies in the public, commercial, or
206 not-for-profit sectors.

207

208 **Acknowledgements**

209 The authors acknowledge present and past members of the Lanzavecchia's group for comments and
210 feedback on the software.

211

212 **References**

213 1. Miho E, Yermanos A, Weber CR, Berger CT, Reddy ST, Greiff V. Computational Strategies for
214 Dissecting the High-Dimensional Complexity of Adaptive Immune Repertoires. *Frontiers in Immunology*.
215 2018;9(224).

- 216 2. Wang C, Liu Y, Cavanagh MM, Le Saux S, Qi Q, Roskin KM, et al. B-cell repertoire responses to
217 varicella-zoster vaccination in human identical twins. *Proc Natl Acad Sci U S A*. 2015;112(2):500-5.
- 218 3. Hoehn KB, Fowler A, Lunter G, Pybus OG. The Diversity and Molecular Evolution of B-Cell Receptors
219 during Infection. *Mol Biol Evol*. 2016;33(5):1147-57.
- 220 4. Zhu J, Ofek G, Yang Y, Zhang B, Louder MK, Lu G, et al. Mining the antibodyome for HIV-1-neutralizing
221 antibodies with next-generation sequencing and phylogenetic pairing of heavy/light chains. *Proceedings of*
222 *the National Academy of Sciences*. 2013;110(16):6470-5.
- 223 5. Jackson KJ, Liu Y, Roskin KM, Glanville J, Hoh RA, Seo K, et al. Human responses to influenza
224 vaccination show seroconversion signatures and convergent antibody rearrangements. *Cell host & microbe*.
225 2014;16(1):105-14.
- 226 6. Tonegawa S. Somatic generation of antibody diversity. *Nature*. 1983;302(5909):575-81.
- 227 7. Adler R. Janeway's immunobiology. *Choice: Current Reviews for Academic Libraries*.
228 2008;45(10):1793-4.
- 229 8. Boyd SD, Marshall EL, Merker JD, Maniar JM, Zhang LN, Sahaf B, et al. Measurement and clinical
230 monitoring of human lymphocyte clonality by massively parallel VDJ pyrosequencing. *Sci Transl Med*.
231 2009;1(12):12ra23.
- 232 9. Stern JN, Yaari G, Vander Heiden JA, Church G, Donahue WF, Hintzen RQ, et al. B cells populating the
233 multiple sclerosis brain mature in the draining cervical lymph nodes. *Sci Transl Med*. 2014;6(248):248ra107.
- 234 10. Tsioris K, Gupta NT, Ogunniyi AO, Zimnisky RM, Qian F, Yao Y, et al. Neutralizing antibodies against
235 West Nile virus identified directly from human B cells by single-cell analysis and next generation sequencing.
236 *Integr Biol (Camb)*. 2015;7(12):1587-97.
- 237 11. Jiang N, He J, Weinstein JA, Penland L, Sasaki S, He XS, et al. Lineage structure of the human antibody
238 repertoire in response to influenza vaccination. *Sci Transl Med*. 2013;5(171):171ra19.
- 239 12. Glanville J, Kuo TC, von Budingen HC, Guey L, Berka J, Sundar PD, et al. Naive antibody gene-segment
240 frequencies are heritable and unaltered by chronic lymphocyte ablation. *Proc Natl Acad Sci U S A*.
241 2011;108(50):20066-71.

- 242 13. Greiff V, Miho E, Menzel U, Reddy ST. Bioinformatic and Statistical Analysis of Adaptive Immune
243 Repertoires. *Trends in immunology*. 2015;36(11):738-49.
- 244 14. Yaari G, Vander Heiden JA, Uduman M, Gadala-Maria D, Gupta N, Stern JN, et al. Models of somatic
245 hypermutation targeting and substitution based on synonymous mutations from high-throughput
246 immunoglobulin sequencing data. *Front Immunol*. 2013;4:358.
- 247 15. Barak M, Zuckerman NS, Edelman H, Unger R, Mehr R. IgTree: creating Immunoglobulin variable
248 region gene lineage trees. *J Immunol Methods*. 2008;338(1-2):67-74.
- 249 16. Andrews SF, Kaur K, Pauli NT, Huang M, Huang Y, Wilson PC. High preexisting serological antibody
250 levels correlate with diversification of the influenza vaccine response. *J Virol*. 2015;89(6):3308-17.
- 251 17. Wu X, Zhang Z, Schramm CA, Joyce MG, Kwon YD, Zhou T, et al. Maturation and Diversity of the
252 VRC01-Antibody Lineage over 15 Years of Chronic HIV-1 Infection. *Cell*. 2015;161(3):470-85.
- 253 18. Felsenstein J. PHYLIP - phylogeny inference package (version 3.2). *Cladistics*. 1989;5:164-6.
- 254 19. Liao HX, Lynch R, Zhou T, Gao F, Alam SM, Boyd SD, et al. Co-evolution of a broadly neutralizing HIV-1
255 antibody and founder virus. *Nature*. 2013;496(7446):469-76.
- 256 20. Pappas L, Foglierini M, Piccoli L, Kallewaard NL, Turrini F, Silacci C, et al. Rapid development of broadly
257 influenza neutralizing antibodies through redundant mutations. *Nature*. 2014;516(7531):418-22.
- 258 21. Kallewaard NL, Corti D, Collins PJ, Neu U, McAuliffe JM, Benjamin E, et al. Structure and Function
259 Analysis of an Antibody Recognizing All Influenza A Subtypes. *Cell*. 2016;166(3):596-608.
- 260 22. Tan J, Pieper K, Piccoli L, Abdi A, Perez MF, Geiger R, et al. A LAIR1 insertion generates broadly
261 reactive antibodies against malaria variant antigens. *Nature*. 2016;529(7584):105-9.
- 262 23. Di Niro R, Lee SJ, Vander Heiden JA, Elsner RA, Trivedi N, Bannock JM, et al. Salmonella Infection
263 Drives Promiscuous B Cell Activation Followed by Extrafollicular Affinity Maturation. *Immunity*.
264 2015;43(1):120-31.
- 265 24. Revell LJ, Chamberlain SA. Rphylip: an R interface for PHYLIP. *Methods in Ecology and Evolution*.
266 2014;5(9):976-81.

- 267 25. Huson DH, Richter DC, Rausch C, DeZulian T, Franz M, Rupp R. Dendroscope: An interactive viewer
268 for large phylogenetic trees. *BMC Bioinformatics*. 2007;8:460.
- 269 26. Wu X, Zhang Z, Schramm CA, Joyce MG, Kwon YD, Zhou T, et al. Maturation and Diversity of the
270 VRC01-Antibody Lineage over 15 Years of Chronic HIV-1 Infection. *Cell*. 2015;161(3):470-85.
- 271 27. Lefranc MP, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, et al. IMGT(R), the
272 international ImMunoGeneTics information system(R) 25 years on. *Nucleic Acids Res*. 2015;43(Database
273 issue):D413-22.
- 274 28. Abhinandan KR, Martin AC. Analysis and improvements to Kabat and structurally correct numbering
275 of antibody variable domains. *Mol Immunol*. 2008;45(14):3832-9.
- 276 29. Yaari G, Uduman M, Kleinstein SH. Quantifying selection in high-throughput Immunoglobulin
277 sequencing data sets. *Nucleic Acids Res*. 2012;40(17):e134.
- 278 30. Kepler TB. Reconstructing a B-cell clonal lineage. I. Statistical inference of unobserved ancestors.
279 *F1000Res*. 2013;2:103.
- 280 31. Kepler TB, Munshaw S, Wiehe K, Zhang R, Yu JS, Woods CW, et al. Reconstructing a B-Cell Clonal
281 Lineage. II. Mutation, Selection, and Affinity Maturation. *Front Immunol*. 2014;5:170.
- 282 32. Sievers F, Higgins DG. Clustal omega. *Curr Protoc Bioinformatics*. 2014;48:3.13.1-6.
- 283 33. Brochet X, Lefranc MP, Giudicelli V. IMGT/V-QUEST: the highly customized and integrated system for
284 IG and TR standardized V-J and V-D-J sequence analysis. *Nucleic Acids Res*. 2008;36(Web Server issue):W503-
285 8.
- 286 34. Collins PL, Fearn R, Graham BS. Respiratory syncytial virus: virology, reverse genetics, and
287 pathogenesis of disease. *Curr Top Microbiol Immunol*. 2013;372:3-38.
- 288 35. Widmer K, Griffin MR, Zhu Y, Williams JV, Talbot HK. Respiratory syncytial virus- and human
289 metapneumovirus-associated emergency department and hospital burden in adults. *Influenza Other Respir
290 Viruses*. 2014;8(3):347-52.

- 291 36. Lozano R, Naghavi M, Foreman K, Lim S, Shibuya K, Aboyans V, et al. Global and regional mortality
292 from 235 causes of death for 20 age groups in 1990 and 2010: a systematic analysis for the Global Burden of
293 Disease Study 2010. *Lancet*. 2012;380(9859):2095-128.
- 294 37. Goodwin E, Gilman MSA, Wrapp D, Chen M, Ngwuta JO, Moin SM, et al. Infants Infected with
295 Respiratory Syncytial Virus Generate Potent Neutralizing Antibodies that Lack Somatic Hypermutation.
296 *Immunity*. 2018;48(2):339-49.e5.
- 297 38. Gilman MS, Castellanos CA, Chen M, Ngwuta JO, Goodwin E, Moin SM, et al. Rapid profiling of RSV
298 antibody repertoires from the memory B cells of naturally infected adult donors. *Science immunology*.
299 2016;1(6).
- 300 39. Mousa JJ, Kose N, Matta P, Gilchuk P, Crowe JE, Jr. A novel pre-fusion conformation-specific
301 neutralizing epitope on the respiratory syncytial virus fusion protein. *Nat Microbiol*. 2017;2:16271.
- 302 40. McLellan JS, Chen M, Joyce MG, Sastry M, Stewart-Jones GB, Yang Y, et al. Structure-based design of
303 a fusion glycoprotein vaccine for respiratory syncytial virus. *Science*. 2013;342(6158):592-8.
- 304
- 305