1  # Real-time analysis and visualization of

2  # pathogen sequence data

3

4  Richard A. Neher[1,2] and Trevor Bedford[3]

5  [1]  Biozentrum, University of Basel, Switzerland

6  [2]  SIB Swiss Institute of Bioinformatics, Basel, Switzerland

7  [3]  Vaccine and Infectious Disease Division,

8  Fred Hutchinson Cancer Research Center, Seattle, USA

9  August 12, 2018

10  Abstract

11  The rapid development of sequencing technologies has to led to an

12  explosion of pathogen sequence data that are increasingly collected as

13  part of routine surveillance or clinical diagnostics. In public health,

14  sequence data is used to reconstruct the evolution of pathogens,

15  anticipate future spread, and target interventions. In clinical settings

16  whole genome sequences identify pathogens at the strain level, can be

17  used to predict phenotypes such as drug resistance and virulence, and

18  inform treatment by linking to closely related cases. While sequencing has

19  become cheaper, the analysis of sequence data has become an important

20  bottleneck. Deriving interpretable and actionable results for a large variety

21  of pathogens – each with their own complexities – from continuously

22  updated data is a daunting task and requires flexible bioinformatics

23  workflows and dissemination platforms. Here, we review recent

24  developments in real-time analysis of pathogen sequence data with a

25  particular focus on visualization and integration of sequence and

26  phenotypic data.

27  As pathogens replicate and spread, their genomes accumulate mutations.

28  These changes can now be detected via cheap and rapid whole genome

29  sequencing on unprecedented scale. Such sequence data are increasingly used

30  to track the spread of pathogens and predict their phenotypic properties. Both

31  applications have great potential to inform public health and treatment decisions

32  if sequencing data can be obtained and analyzed rapidly. Historically, however,

33  sequencing and analysis has lagged months-to-years behind sample collection.

34  The results from these studies have taught us much about pathogen molecular

35  evolution, genotype-phenotype maps, and epidemic spread, but have come

36  almost always too late to inform public health interventions or treatment

37  decisions.

38  The rapid development of sequencing technologies has made routine

39  sequencing of viral and bacterial genomes possible and tens of thousands of

40  whole genome sequences (WGS) are deposited in databases every year (see Fig.

41  1). Many more, regrettably, are sequenced and not shared. There are currently

42  two major directions in which high-throughput sequencing technologies are used

43  in public health and diagnostics: (i) to track outbreaks and epidemics to inform

44  public health response, and (ii) to characterize individual infections to tailor

45  treatment decisions.

46    *Sequencing in public health.* The utility of rapid sequencing and phylogenetic

47    analysis of pathogens is perhaps most evident for influenza viruses and food-

48    borne diseases. Due to rapid evolution of its viral surface proteins, the antigenic

49    properties of the circulating influenza viruses change every few years and the

50    seasonal influenza vaccine needs frequent updating (1). The WHO Global

51    Influenza Surveillance and Response System (GISRS) sequences hundreds of

52    viruses every month and many of these sequences are submitted to the GISAID

53    database (gisaid.org) within 4 weeks of sample collection. Phylogenetic analysis

54    of these data provide an accurate and up-to-date summary of the spread and

55    abundance of different viral variants that is crucial input to the biannual

56    consultations on seasonal influenza vaccine composition.

57    Such rapid turn-around and data sharing is considerably harder to achieve in

58    an outbreak setting in resource limited conditions. Nonetheless, Quick et al. (2)

59    achieved even shorter turn-around during the tail end of the 2014–2015 West

60    African Ebola outbreak. Similarly, Dyrdak et al. (3) analyzed an enterovirus

61    outbreak in Sweden and continuously updated the manuscript until publication

62    with sequences sampled within days of publication included in the analysis.

63    Molecular epidemiology techniques can reconstruct the temporal and spatial

64    spread of an outbreak. In this case, the accumulation of mutations alongside a

65    molecular clock estimate can be used to date the origin of an outbreak. Similarly,

66    by linking samples that originate from different geographic locations,

67    phylogeographic methods can reconstruct geographic spread and differentiate

68    distinct introductions. The resolution of these inferences critically depends on

69    the rate at which mutations accumulate in the sequenced locus, which increases

70    with the per site evolutionary rate and the length of the locus.

3

71    RNA viruses accumulate changes in their genome with a typical rate of 0.0005

72    to 0.005 changes per site and year (4). Rate estimates vary from virus to virus

73    and depend on the time scale of observation or whether measured within or

74    between hosts. Ebola virus and Zika virus, for example, evolve at a rate of $\mu \approx$

75    0.001 per site per year. The expected time interval without a substitution along

76    a transmission chain is $1/(\mu L)$, which corresponds to approximately 5 weeks for

77    Zika virus ($L \approx 10$kb) and 3 weeks for Ebola virus ($L \approx 19$kb). Hence evolution and

78    spread of such RNA viruses can be resolved on the scale of a month. While this

79    temporal resolution is typically insufficient to resolve individual transmissions, it

80    is high compared to the duration of outbreaks. Rapid sequencing and analysis

81    therefore has the potential to inform intervention efforts as outbreaks are

82    unfolding. In particular, they rule out direct transmission and differentiate

83    different introductions or zoonosis.

84    Phylodynamic and phylogeographic methods are best established for viral

85    pathogens with high evolutionary rates and small genomes for which large scale

86    sequencing has been possible for years. The evolutionary rates of bacteria are

87    many orders of magnitudes lower than those of RNA viruses. However, bacteria

88    also have about 100 to 1000-fold larger genomes and it is now possible to

89    sequence entire bacterial genomes at low cost. Substitution rate estimates in

90    bacteria come with substantial uncertainty but they tend to be on the order of

91    one substitution per megabase per year (with about one to two orders of

92    magnitude of variation between species (5)). With a typical genome size of 5

93    megabases, this translates into 5–10 substitutions per genome and year —

94    similar to many RNA viruses. The substitution rate in the core genome of MRSA,

95    for example, was estimated to be $1.3 \times 10^{-6}$ per site and year (6). The core

96    genome of *Listeria monocytogenes* evolves more slowly at about one

97    substitution every 2.5y (7). Hence real-time phylogenetics for bacterial outbreak

98    tracking is possible in much the same way as for RNA viruses. Analysis of bacterial

99    genomes, however, is vastly more complicated than that of RNA viruses with

100   short genomes. Bacteria frequently exchange genetic material via horizontal

101   transfer, take up genes from the environments and rearrange their genome.

102   Recombination can blur phylogenetic signal and recombinant sequence is often

103   difficult to remove. Furthermore, strong selection within hosts, for example

104   through drug therapy, can accelerate evolution by an order of magnitude (8). If

105   not properly accounted for, these processes can blur any temporal signal and

106   obscure links between closely related isolates.

107   Even with whole genomes, phylogenetic resolution typically is insufficient to

108   make the case for a direct transmission, but transmission can be confidently

109   ruled out for divergent sequences, seemingly unrelated cases can be grouped

110   into outbreaks (e.g. an outbreak of drug resistant MtB among migrants arriving

111   in multiple European countries (9)), predominant routes of transmission and

112   likely sources in the environment or animal reservoirs can be identified.

113   GenomeTrackr and PulseNet, for example, are a large federated efforts to

114   sequence tens of thousands of genomes from food-borne outbreaks and clinical

115   samples (10; 11). All sequence data from these projects are publicly available on

116   NCBI with little delay and are analyzed in real-time to track outbreaks. The

117   recently    released    Pathogen    Detection    system    by    NCBI

118   (www.ncbi.nlm.nih.gov/pathogens/) provides    convenient    access    to    the

119   sequence and metadata generated by these projects as well as phylogenetic

120   analysis.

121    These examples illustrate the potential and feasibility of obtaining actionable

122    information from pathogen sequence data for both viral and bacterial infections.

123    However, with rapidly increasing data volumes, efficient processing pipelines and

124    tools that help with interpretation – e.g. visualizations – increasingly become the

125    bottleneck.

126    *Sequencing in diagnostics and therapy.* For some pathogens like Zika virus,

127    sequencing the genome has no implications for treatment. In the case of HIV,

128    however, drug resistance profiles derived from sequence data have directly

129    informed treatment for years (12). As the genetic basis of drug resistance

130    phenotypes are better understood, rapid whole genome sequencing will

131    increasingly be used to diagnose and phenotype pathogens directly from the

132    clinical specimen. Such culture-free methods are particularly important for

133    tuberculosis, in which culture based susceptibility testing takes many weeks.

134    Votintseva et al. (13) have recently shown that high-throughput sequencing

135    directly from respiratory samples can provide drug resistance profiles of *M.*

136    *tuberculosis* within a day.

137    Sequencing for diagnostic purposes or for public health surveillance have

138    different aims and requirements, but can complement each other. Public health

139    response typically requires recent data with an emphasis on dynamics.

140    Surveillance data provides context for the individual case in the clinics requires a

141    stable database with validated content to make reliable predictions on drug

142    susceptibility, phylogenetic context, and protective measures. Clinical

143    sequencing data, however, should be fed into surveillance databases

144    immediately whenever ethically possible. Only with rapid and open sharing of

6

145    sequencing data can the full potential of molecular epidemiology be realized

146    (11).

147        The challenges involved in sample collection, processing, sequencing and

148    data sharing have been discussed at length elsewhere (14). Here, we focus on

149    software developments that facilitate the implementation of real-time analysis

150    with an emphasis on web-based visualization, as a full review of general tools for

151    genomic analysis and visualization is not easily encompassed.


## Rapid and interpretable analysis of genomic data

153    A typical molecular epidemiological analysis aims to identify transmission

154    clusters, date the introduction of the pathogen, detail geographic spread, and in

155    some cases identify potential phenotypic change of a pathogen from sequence

156    data. The rapidly increasing numbers of sequenced genomes make

157    comprehensive analysis computationally challenging. While 1000s of viral

158    genomes can be aligned within minutes (e.g. by MAFFT) and the reconstruction

159    of a basic phylogenetic tree typically takes less than one hour (e.g. using IQ-TREE,

160    RAxML or FastTree), the most popular tool for phylodynamic inference (BEAST)

161    (15) will often take weeks to finish.

162        To overcome these hurdles, several tools that use simpler heuristics have

163    been developed to infer time-stamped phylogenies (16; 17; 18). Rather than

164    sampling a large number of tree topologies, these tools use the topology of an

165    input tree with little or no modification. Dating of ancestral events tends to be of

166    comparable accuracy to BEAST (16; 17; 18). However, these tools do not

167    integrate uncertainty of tree reconstruction and provide limited flexibility to

168    infer demographic models. Furthermore, the heuristics used by these program

169    are based on assumptions (for example that sequences are closely related) and

170    they are not expected to be accurate in all scenarios. The computational cost of

171    Bayesian phylodynamics could be mitigated if methods for continuous updating

172    and augmenting of the Markov chain with additional data were developed. For

173    the time being, however, efficient heuristics and sensible approximations deliver

174    sufficiently accurate and reliable results when near real-time analysis is required.

175    Viral genomes: Nextflu and Nextstrain

176    The number of influenza viruses that are sequenced and phenotyped per month

177    has increased sharply to a point that a comprehensive and timely manual analysis

178    and annotation of the results is no longer feasible. In 2014, we developed an

179    automated phylodynamic analysis pipeline that operates on an up-to-date

180    database of sequences and serological information. The results of this pipeline

181    are available were made available at nextflu.org and included a phylogeny,

182    branch-specific mutations, frequency trajectories of mutations and variants, and

183    a model of antigenic evolution.

184        Nextflu is now integrated in the more general platform Nextstrain, that provides

185    an online platform for outbreak investigations of diverse viruses and is available at

186    nextstrain.org (19). Nextstrain uses TreeTime (18) to infer time-scaled phylogenies

187    and conduct ancestral sequence inference. In addition, Nextstrain uses the discrete

188    ancestral character inference of TreeTime to infer the likely geographic state of

189    ancestral nodes. Since this approach applies "mutation" models to "migration", it is

190    often called a "mugration model". A phylodynamic/phylogeographic analysis of 1000

8

191    sequences of length 10kb takes on the order of an hour on a standard laptop
192    computer.

193    Bacterial WGS data

194    Bacterial WGS data typically comes in the form of millions of short reads that can
195    either be assembled into contigs, mapped against reference sequences, or
196    classified based on kmer distributions. A large number of tools have been
197    developed for rapid species identification, typing, and variant calling. WGSA by
198    the , for example, allows the user to upload an assembly and WGSA will detect
199    the species and infer the multi-locus sequence type within a few seconds. In
200    addition, WGSA predicts antibiotic resistance profiles for a number of species.
201    WGSA was developed by the Center for Genomic Pathogen Surveillance and is
202    available at wgsa.net.

203        Bacterial genomes are very dynamic and frequently gain or lose genes. Even
204    closely related bacteria can differ in the presence or absence of dozens of genes.
205    To track transmission and detect clusters, genomes are typically compared at a
206    set of *core genes* present in all bacteria of a species. Genes present in only a
207    fraction of individuals are referred to as *accessory genes*.

208        Clinically important genes such antibiotic resistance determinants or
209    virulence factors are often not part of the core genome and are horizontally
210    transferred between strains and species. Collections of bacterial genomes are
211    therefore analyzed using pan-genome tools that aim to cluster all genes in the
212    collection of genomes into orthologous groups. Early methods for pan-genome
213    analysis scaled poorly with the number of genomes that are analyzed since every
214    gene in every genome needed to be compared to every other gene. The first tool

215    capable of analyzing 100s of bacterial genomes was Roary (20). Roary is designed

216    to work with very similar genes (as is common in outbreak scenarios) and

217    accelerates inference of orthologous gene clusters by pre-clustering genomes. A

218    more recent pan-genome analysis pipeline capable of large scale analysis is panX

219    (21) that speeds up clustering by hierarchically building up the complete pan-

220    genome from sub-pan-genomes inferred from smaller batches of genomes. PanX

221    is coupled to a web-based visualization platform discussed below.

222    While the pan-genome tools cluster annotated genes in the collection of

223    genomes, they are of little help to assess the origin and distribution of a particular

224    sequence. Traditional tools for homology search in NCBI only index assembled

225    sequence, but today the majority of sequence data are stored in short read archives

226    rather than in Genbank. Bradley et al. (22) developed a method to search the entire

227    collection of microbial sequence data including metagenomic samples from a wide

228    variety of environment. The ability to search this vast treasure trove of data will

229    likely be transformative in assessing spread and prevalence of novel resistance

230    determinants. The recently discovered mobile colistin resistance gene *mcr-1*, for

231    example, was found in more than 100 datasets where it wasn't previously

232    described(22).

233    Outlook

234    Most current analysis pipelines require rerunning the entire analysis even when

235    only a single sequence is added. While this strategy is still feasible today, this will

236    likely become unsustainable in the future. Applications that support cheap

237    updating of datasets and on-line addition of user data will likely replace current

238    versions.

10

## Visualization and interpretation

239

240    With increasing dataset sizes, interpretation and exploration of data become

241    increasingly challenging. Phylogenetic trees can be visualized as familiar planar

242    graphs, but the tree alone only shows genetic similarity between isolates and

243    becomes quickly unintelligible as the number of sequences increases. To make

244    pathogen sequence data truly useful, it needs to be integrated with other types

245    of information, ideally in an interactive way. A suitable platform to do so is the

246    web-browser and several powerful web applications have emerged over the last

247    few years. In addition, browser-based visualizations are naturally disseminated

248    online.

249    Microreact

250    Microreact is a web application based on React (a JavaScript framework for

251    interactive applications), D3.js (a JavaScript library for producing dynamic,

252    interactive data visualizations), Phylocanvas (a JavaScript flexible tree viewer),

253    and Leaflet (a JavaScript mapping toolkit) (23). Microreact allows exploration of

254    a phylogenetic tree, the geographic locations, and a time line of the samples. It

255    is available at microreact.org. Custom data sets can be loaded into the

256    application in the form of a Newick tree and a tabular file containing information

257    such as geographic location or sampling data.

258    Nextstrain

259    Nextstrain was developed as a more generic and flexible version of Nextflu (19)

260    which is available at nextstrain.org. Similar to Microreact, Nextstrain uses React,

11

261   D3.js, and Leaflet, but uses a custom made tree viewer that has flexible zooming

262   and annotation options. The tree can be decorated with any discrete or

263   continuous attribute, both on tips of the tree and inferred values for internal

264   nodes (for example geographic location). Nextstrain maps individual mutations

265   to branches in the tree and thereby allows to associate mutations with

266   phenotypes or geographic distributions. The map in Nextstrain shows putative

267   transmission events and a panel indicates genetic diversity across the genome

268   (Fig. 2).

269       The analyses by Nextstrain and Nextflu critically depend on timely and open

270   sharing of sequence information that many laboratories around the globe

271   contribute. To incentivize early pre-publication sharing of data, platforms like

272   Nextstrain need to explicitly acknowledge the individual contributions. Ideally, such

273   platforms should provide added value to authors, such as for example deep links

274   that show data by a particular group in the context of the outbreak.

275   Phandango

276   Phandango is an interactive viewer for bacterial whole genome sequencing data

277   (24) and combines a phylogenetic tree with metadata columns and gene

278   presenceabsence maps or recombination events. Phandango is available at

279   phandango.net and can ingest the output of a number of analysis tools

280   commonly used for the analysis of bacterial WGS data such Gubbins, Roary and

281   BRAT.

12

282   panX

283   PanX is a pan-genome analysis pipeline that is coupled to a web-browser based
284   visualization (21). Similar to Phandango, it displays a core genome SNP phylogeny
285   but is otherwise more centered on genetic variation in individual genes.
286   Pangenomes of about 100 bacterial species based on curated reference genomes
287   are available at pangenome.de. The tree and alignment of each gene in the
288   pangenome can be accessed rapidly by a searching a table of gene names and
289   annotations. PanX then displays gene and species tree side by side and maps
290   gene gain and loss events to branches in the core genome tree and mutations to
291   branches in the gene tree. Trees can be colored by arbitrary attributes such as
292   resistance phenotypes and associations between genetic variation and these
293   phenotypes can be explored.

294   Other tools

295   SpreaD3 allows of visualization of phylogeographic reconstructions from models
296   implemented in the software package BEAST (25). PhyloGeoTool is a
297   webapplication to interactively navigate large phylogenies and to explore
298   associated clinical and epidemiological data (26). TreeLink displays phylogenetic
299   trees alongside metadata in an interactive web application (27).

300   ## Challenges in data integration and visualization

301   With rapidly increasing volumes of sequence data, decisions as to how the data
302   are filtered and what analysis are shown become increasingly important.
303   Epidemiological investigations of a novel outbreak typically seek to identify the

13

304    sources, track the spread, and detect transmission chains. In this case, a generic

305    combination of map, tree, and time line will often be an appropriate and

306    sufficient visualization. Nextstrain and Microreact both follow this paradigm.

307        However, when analyzing established pathogens that continuously adapt to

308    treatments, vaccines, or pre-existing immunity, more specific applications will be

309    necessary since case data, phenotype data, and clinical parameters differ wildly

310    by pathogen. Such data will generally have a common core such as sample date

311    and location, but other parameters such as drug resistance phenotypes, disease

312    severity, host age, risk group, serology, etc., are pathogen specific. While these

313    data are often stored in non-standardized formats and ethical and technical

314    reasons can impede data sharing, these data are often at least as important for

315    interpretation of the epidemiological dynamics as phylodynamic inference from

316    sequence data. The value of either data is greatly increased by seamless

317    integration, but the idiosyncrasies require flexible analysis and visualization

318    frameworks that can be tailored to specific pathogens.

319    One such example is the serological characterization of influenza viruses via

320    hemagglutination inhibition (HI) titers using antisera raised in ferrets. Such titers

321    are routinely measured as part of GISRS to monitor the antigenic evolution of

322    influenza viruses are a good example how phenotype information can be

323    interactively integrated with phylogeny and molecular evolution. HI titers are

324    reported in large tables and have been traditionally visualized using

325    multidimensional scaling without any reference to the phylogeny. In Bedford et

326    al. (28) and Neher et al. (29), we developed methods to integrate the molecular

327    and antigenic evolution of influenza virus. This integration allows association of

328    genotypic changes with antigenic evolution and suggests an intuitive and

14

329    interactive visualization of HI titer data on the phylogeny. A screenshot of this

330    integration is shown in Fig. 3. Due to data sharing restrictions, most HI titer data

331    are not openly available, but historical data by McCauley and colleagues are

332    visualized along with the molecular evolution at hi.nextflu.org.

333    In addition to phenotype integration, it is crucial to choose the right level of

334    detail for a specific application. This is particularly true for bacteria where the

335    relevant    information    might    be    the    core    genome    phylogeny,    the

336    presence/absence of particular genes or plasmids, or individual mutations in

337    specific genes. If the analysis tool and the visualization does not provide a fine

338    grained analysis at the relevant level, the most important patterns might stay

339    hidden. On the other hand, sifting through every gene or mutation is prohibitive.

340    The primary aim should be to highlight the most important and robust patterns

341    upfront and provide flexible methods to filter and rank variants (e.g. by recent

342    rise in frequency, association with host, resistance or risk group, etc). The user

343    should then have the possibility to expose detail on demand when a deeper

344    exploration is required.

345    Similarly, parameter inferences and model abstractions are very useful to get

346    a concise summary of the data, but should be complemented by the ability of

347    interrogate the raw data (e.g. an estimate of the evolutionary rate should be

348    accompanied by a scatter plot of root-to-tip divergence and sampling time). This

349    is particularly important in outbreak scenarios when methods are applied to an

350    emerging pathogen in a developing situation.

351    For clinical applications, the presentation of the results of an analysis should

352    be focused on the sample in question and only provide reliable and actionable

353    information, while suggestive and correlative results tend to be a distraction (30).

15

## Conclusions

354

355  High-throughput and rapid sequencing is revolutionizing infectious disease
356  diagnostics and epidemiology. Sequence data can be used to unambiguously
357  identify pathogens, to link related cases, to reconstruct the spread of an
358  outbreak, and will soon allow detailed prediction of a pathogen's phenotype.

359      The Global Influenza Surveillance and Response System (GISRS) is a good
360  example of a near real-time surveillance system. Hundreds of viruses are
361  sequenced and phenotyped every month and the sequence data are shared in a
362  timely manner. A global comprehensive analysis of these data, updated about
363  once a week, is available at nextflu.org. These analyses directly inform the
364  influenza vaccine strain selection process (1).

365      Several public health agencies have adopted WGS as their primary tool for
366  outbreak investigation and many centers share these data openly with
367  commendable timeliness. The GenomeTrakr and PulseNet networks, for
368  example, now sequence and openly release about 5000 bacterial genomes per
369  month (11; 10). These data are accessible on NCBI through the recently released
370  Pathogen Detection system at www.ncbi.nlm.nih.gov/pathogens with analysis
371  results available via FTP.

372      These two examples clearly show that near real-time genomic surveillance is
373  possible and valuable and all the individual components to implement such
374  surveillance are in place. However, to realize this potential for many more
375  pathogens, sample collection and sequencing has to be streamlined, data need to
376  be shared along with the relevant metadata, and specific analysis methods and
377  visualizations need to be implemented and maintained.

16

## Acknowledgments

## References

385 [1]    Morris DH, Gostic KM, Pompei S, Bedford T, Łuksza M, Neher RA, Grenfell
386        BT, Lassig M, McCauley JW. 2017. Predictive Modeling of Influenza Shows
387        the¨ Promise of Applied Evolutionary Biology. Trends microbiology.

388 [2]    Quick J, Loman NJ, Duraffour S, Simpson JT, Severi E, Cowley L, Bore JA,
389        Koundouno R, Dudas G, Mikhail A, Ouédraogo N, Afrough B, Bah A, Baum
390        JH, Becker-Ziaja B, Boettcher JP, Cabeza-Cabrerizo M, Camino-Sanchez A,
391        Carter LL, Doerrbecker J, Enkirch T, Dorival IGG, Hetzelt N, Hinzmann J,
392        Holm T, Kafetzopoulou LE, Koropogui M, Kosgey A, Kuisma E, Logue CH2,10,
393        Mazzarelli A, Meisel S, Mertens M, Michel J, Ngabo D, Nitzsche K, Pallash E,
394        Patrono LV, Portmann J, Repits JG, Rickett NY, Sachse A, Singethan K,
395        Vitoriano I, Yemanaberhan RL, Zekeng EG, Trina R, Bello A, Sall AA, Faye O,
396        Faye O, Magassouba N, Williams CV, Amburgey V, Winona L, Davis E,
397        Gerlach J, Washington F, Monteil V, Jourdain M Bererd M, Camara A,
398        Somlare H, Camara A, Gerard M, Bado G, Baillet B, Delaune D, Nebie KY,
399        Diarra A, Savane Y, Pallawo RB, Gutierrez GJ, Milhano N, Roger I, Williams
400        CJ, Yattara F, Lewandowski K, Taylor J, Rachwal P, Turner D, Pollakis G,

401     Hiscox JA, Matthews DA, O'Shea MK, Johnston AM, Wilson D, Hutley E, Smit
402     E, Di Caro A, Woelfel R, Stoecker K, Fleischmann E, Gabriel M, Weller SA,
403     Koivogui L, Diallo B, Keita S, Rambaut A, Formenty P, Gunther S, Carroll MW.
404     2016. Real-time, portable genome sequencing for Ebola surveillance.
405     Nature. 530(7589):228–232.

406  [3]  Dyrdak R, Grabbe M, Hammas B, Ekwall J, Hansson KE, Luthander J, Naucler
407       P, Reinius H, Rotzen-Ostlund M, Albert J. 2016. Outbreak of enterovirus
408       D68¨ of the new B3 lineage in Stockholm, Sweden, August to September
409       2016. Eurosurveillance; 21(46).

410  [4]  Duffy S, Shackelton LA, Holmes EC. 2008. Rates of evolutionary change in
411       viruses: patterns and determinants. Nat. Rev. Genet. 9(4):267– 276.
412       https://www.nature.com/articles/nrg2323.

413  [5]  Duchene S, Holt KE, Weill FX, Le Hello S, Hawkey J, Edwards DJ, Fourment
414       M, Holmes EC. 2016. Genome-scale rates of evolutionary change in
415       bacteria. Microb. Genomics Nov; 2(11). https://www.ncbi.nlm.nih.
416       gov/pmc/articles/PMC5320706/.

417  [6]  Holden MT, Hsu LY, Kurt K, Weinert LA, Mather AE, Harris SR, Strommenger
418       B, Layer F, Witte W, de Lencastre H, Skov R, Westh H, Zemlicková H, Coombs
419       G, Kearns AM, Hill RL, Edgeworth J, Gould I, Gant V, Cooke J, Edwards GF,
420       McAdam PR, Templeton KE, McCann A, Zhou Z, Castillo-Ramírez S, Feil EJ,
421       Hudson LO, Enright MC, Balloux F, Aanensen DM, Spratt BG, Fitzgerald JR,
422       Parkhill J, Achtman M, Bentley SD, Nübel U.. 2013. A genomic portrait of

423        the emergence, evolution, and global spread of a methicillin-resistant

424        Staphylococcus aureus pandemic. Genome research; 23(4):653–664.

425   [7]   Moura A, Criscuolo A, Pouseele H, Maury MM, Leclercq A, Tarr C, Bjrkman

426        JT, Dallman T, Reimer A, Enouf V, Larsonneur E, Carleton H, Bracq-Dieye H,

427        Katz LS, Jones L, Touchon M, Tourdjman M, Walker M, Stroika S, Cantinelli

428        T, Chenal-Francisque V, Kucerova Z, Rocha EPC, Nadon C, Grant K, Nielsen

429        EM, Pot B, Gerner-Smidt P, Lecuit M, Brisse S. 2016. Whole genomebased

430        population biology and epidemiological surveillance of Listeria

431        monocytogenes. Nat. Microbiol. 2:16185. http://www.nature.

432        com/articles/nmicrobiol2016185.

433   [8]   Mwangi MM, Wu SW, Zhou Y, Sieradzki K, Lencastre Hd, Richardson P,

434        Bruce D, Rubin E, Myers E, Siggia ED, Tomasz A. 2007. Tracking the in vivo

435        evolution of multidrug resistance in Staphylococcus aureus by whole-

436        genome sequencing. Proc. Natl. Acad. Sci.; 104(22):9451–9456.

437        http://www.pnas.org/content/104/22/9451.

438   [9]   Walker TM, Merker M, Knoblauch AM, Helbling P, Schoch OD, Werf MJvd,

439        Kranzer K, Fiebig L, Kroger S, Haas W, Hoffmann H, Indra A, Egli A, Cirillo

440        DM, Robert J, Rogers TR, Groenheit R, Mengshoel AT, Mathys V, Haanpera

441        M, Soolingen Dv, Niemann S, Bottger EC, Keller PM, Avsar K, Bauer C,

442        Bernasconi E, Borroni E, Brusin S, Devis MC, Crook DW, Dedicoat M,

443        Fitzgibbon M, Gagneux S, Geiger F, Guthmann JP, Hendrickx D, Hoffmann-

444        Thiel S, Ingen Jv, Jackson S, Jaton K, Lange C, Stalder JM, O'Donnell J, Opota

445        O, Peto TEA, Preiswerk B, Roycroft E, Sato M, Schacher R, Schulthess B,

19

446    Smith EG, Soini H, Sougakoff W, Tagliani E, Utpatel C, Veziris N, Wagner-
447    Wiening C, Witschi M. 2018. A cluster of multidrug-resistant
448    Mycobacterium tuberculosis among patients arriving in Europe from the
449    Horn of Africa: a molecular epidemiological study. The Lancet Infect. Dis.;.

450  [10]  Carleton HA, Gerner-Smidt P. 2016. Whole-Genome Sequencing Is Taking
451    over Foodborne Disease Surveillance. Microbe Mag;11(7):311–317.
452    http://www.asmscience.org/content/journal/microbe/10.1128/microbe.
453    11.311.1.

454  [11]  Stevens EL, Timme R, Brown EW, Allard MW, Strain E, Bunning K, Musser
455    S. 2017. The Public Health Impact of a Publically Available, Environmental
456    Database of Microbial Genomes. Front. Microbiol.; 8.
457    https://www.frontiersin.org/articles/10.3389/ fmicb.2017.00808/full.

458  [12]  Beerenwinkel N, Daumer M, Oette M, Korn K, Hoffmann D, Kaiser R,¨
459    Lengauer T, Selbig J, Walter H. 2003. Geno2pheno: estimating phenotypic
460    drug resistance from HIV-1 genotypes. Nucleic Acids Res; 31(13):3850–
461    3855. https://academic.oup.com/nar/article/ 31/13/3850/2904197.

462  [13]  Votintseva AA, Bradley P, Pankhurst L, del Ojo Elias C, Loose M, Nilgiriwala
463    K, Chatterjee A, Smith EG, Sanderson N, Walker TM, Peto TEA, Crook DW,
464    Iqbal Z. 2017. Same-day diagnostic and surveillance data for tuberculosis
465    via whole-genome sequencing of direct respiratory samples. J. clinical
466    microbiology; 55(5):1285– 1298.

20

467   [14]  Gardy JL, Loman NJ. 2018. Towards a genomics-informed, real-time, global
468          pathogen surveillance system. Nat. Rev. Genet.; 19(1):9.
469          https://www.nature.com/articles/nrg.2017.88.

470   [15]  Drummond AJ, Suchard MA, Xie D, Rambaut A. 2012. Bayesian
471          Phylogenetics with BEAUti and the BEAST 1.7. Mol Biol Evol; 29(8):1969–
472          1973.

473   [16]  To TH, Jung M, Lycett S, Gascuel O. 2016. Fast Dating Using Least-Squares
474          Criteria and Algorithms. Syst. Biol.; 65(1):82–97.

475   [17]  Volz EM, Frost SDW. 2017. Scalable relaxed clock phylogenetic dating. Virus
476          Evol.; 3(2). https://academic.oup.com/ve/article/ 3/2/vex025/4100592.

477   [18]  Sagulenko P, Puller V, Neher RA. 2018. TreeTime: Maximum-likelihood
478          phylodynamic analysis. Virus Evol; 4(1). https://academic.oup.
479          com/ve/article/4/1/vex042/4794731.

480   [19]  Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko
481          P, Bedford T, Neher RA. 2018. Nextstrain: real-time tracking of pathogen
482          evolution. Bioinformatics.
483          https://academic.oup.com/bioinformatics/advance-
484          article/doi/10.1093/bioinformatics/bty407/5001388.

485   [20]  Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, Fookes
486          M, Falush D, Keane JA, Parkhill J. 2015. Roary: rapid largescale prokaryote
487          pan genome analysis. Bioinformatics. 31(22):3691–3693.
488          http://bioinformatics.oxfordjournals. org/content/31/22/3691.

489 [21] Ding W, Baumdicker F, Neher RA. 2017. panX: pan-genome analysis and
490 exploration. Nucleic Acids Res.;
491 https://academic.oup.com/nar/article/doi/10.1093/nar/gkx977/4564799
492 /

493 [22] Bradley P, Bakker Hd, Rocha E, McVean G, Iqbal Z. 2017. Real-time search of
494 all bacterial and viral genomic data. bioRxiv; p. 234955.
495 https://www.biorxiv.org/content/early/2017/12/15/234955.

496 [23] Argimon S, Abudahab K, Goater RJE, Fedosejev A, Bhai J, Glasner C, Feil EJ,
497 Holden MTG, Yeats CA, Grundmann H, Spratt BG, Aanensen DM. 2016.
498 Microreact: visualizing and sharing data for genomic epidemiology and
499 phylogeography. Microb. Genomics. 2(11).
500 http://mgen.microbiologyresearch.org/content/journal/mgen/10.1099/mg
501 en.0.000093.

502 [24] Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, Harris SR.
503 2018. Phandango: an interactive viewer for bacterial population genomics.
504 Bioinformatics, 34(2):292–293. https://academic.oup.com/
505 bioinformatics/article/34/2/292/4212949.

506 [25] Bielejec F, Baele G, Vrancken B, Suchard MA, Rambaut A, Lemey P. 2016.
507 SpreaD3: interactive visualization of spatiotemporal history and trait
508 evolutionary processes. Mol Biol Evol; 33(8):2167–2169.

509 [26] Libin P, Vanden Eynden E, Incardona F, Nowe A, Bezenchek A, Group ES,
510 Sonnerborg A, Vandamme AM, Theys K, Baele G. 2017. PhyloGeoTool:

511    interactively exploring large phylogenies in an epidemiological context.

512    Bioinformatics. 33(24):3993–3995.

513    [27]  Allende C, Sohn E, Little C. 2015. Treelink: data integration, clustering and

514    visualization of phylogenetic trees. BMC Bioinformatics. 16(1):414.

515    [28]  Bedford T, Suchard MA, Lemey P, Dudas G, Gregory V, Hay AJ, McCauley JW,

516    Russell CA, Smith DJ, Rambaut A. 2014. Integrating influenza antigenic

517    dynamics with molecular evolution. eLife; 3.

518    https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3909918/.

519    [29]  Neher RA, Bedford T, Daniels RS, Russell CA, Shraiman BI. 2016. Prediction,

520    dynamics, and visualization of antigenic phenotypes of seasonal influenza

521    viruses. Proc. Natl. Acad. Sci. 113(12):E1701– 1709.

522    [30]  Crisan A, McKee G, Munzner T, Gardy JL. 2018. Evidence-based design and

523    evaluation of a whole genome sequencing clinical report for the reference

524    microbiology laboratory. PeerJ. 6:e4218. https://peerj.com/ articles/4218

525

23

## Figure captions:

526

527 Figure 1: The number of complete pathogen genomes has increased dramatically
528 over the last few years. More than 4000 complete influenza A (IAV) subtype
529 H3N2 virus genomes have been deposited in GISAID in 2017. The GenomeTrakr
530 network sequenced in excess of 40,000 Salmonella genomes and 25,000 other
531 bacterial genomes (mostly *Listeria*, *E. coli/Shigella*, and *Campylobacter*) in 2017
532 (11).

533 Figure 2: Phylogeographic analysis of Zika virus sequences on nextstrain.org (19).
534 Whole genomes sequences sampled between 2013 and 2017 were processed
535 using the Nextstrain pipeline. Nextstrain reconstructs likely time and place of
536 each internal node of the tree and from this assignments infers possible
537 transmission patterns that are displayed on a map. Molecular analysis of this sort
538 reveals for example multiple introductions of Zika virus into Florida originating
539 most likely from viruses circulating in the Caribbean in 2015-2016.

540 Figure 3: Integration of HI titer data with molecular evolution influenza virus.
541 Each year, influenza laboratories determine thousands of HI titers of test viruses
542 relative to sera raised against several reference viruses (indicated by gray cogs).
543 These data can be integrated with the molecular evolution of the virus and
544 visualized on the phylogeny (here showing inferred titers using a model). The
545 reference virus with respect to which titers are displayed can be chosen by

24

546   clicking on the corresponding symbol in the tree (29). The visualization exposes

547   both raw data (via tool tips for each virus) as well as a model inference that

548   integrates many individual measurements (hi.nextflu.org).

Phylogeny

Country ∧

| | |
|---|---|
| Thailand | Panama |
| Vietnam | Nicaragua |
| Singapore | Honduras |
| French Polynesia | El Salvador |
| American Samoa | Guatemala |
| Fiji | Mexico |
| Tonga | Martinique |
| China | Guadeloupe |
| Taiwan | Saint Barthelemy |
| Japan | USVI |
| Brazil | Puerto Rico |
| Peru | Jamaica |
| Ecuador | Dominican Republic |
| Colombia | Haiti |
| French Guiana | Cuba |
| Suriname | USA |
| Venezuela | Philippines |

reset layout

2012.0    2013.0    2014.0    2015.0    2016.0    2017.0

Transmissions

Play    Reset

Leaflet | © OpenStreetMap contributors

log₂ titer distance from
A/Stockholm/6/2014

| | | | |
|---|---|---|---|
| 0 | | 2 | |
| 0.4 | | 2.4 | |
| 0.8 | | 2.8 | |
| 1.2 | | 3.2 | |
| 1.6 | | 3.6 | |

focal virus

A1b/135K
A1b/135N
A1b
A1a
A1
A2
A3
A4
3c2.A
3c3.A
3c3.B
3c3
3c3

2010  2011  2012  2013  2014  2015  2016  2017  2018