

panX: pan-genome analysis and exploration

Wei Ding¹, Franz Baumdicker² and Richard A. Neher^{1,3,*}

¹Max Planck Institute for Developmental Biology, 72076 Tübingen, Germany, ²Mathematisches Institut, Albert-Ludwigs University of Freiburg, 79104 Freiburg, Germany and ³Biozentrum and SIB Swiss Institute of Bioinformatics, University of Basel, 4056 Basel, Switzerland

Received September 14, 2016; Revised September 28, 2017; Editorial Decision October 06, 2017; Accepted October 10, 2017

ABSTRACT

Horizontal transfer, gene loss, and duplication result in dynamic bacterial genomes shaped by a complex mixture of different modes of evolution. Closely related strains can differ in the presence or absence of many genes, and the total number of distinct genes found in a set of related isolates—the pan-genome—is often many times larger than the genome of individual isolates. We have developed a pipeline that efficiently identifies orthologous gene clusters in the pan-genome. This pipeline is coupled to a powerful yet easy-to-use web-based visualization for interactive exploration of the pan-genome. The visualization consists of connected components that allow rapid filtering and searching of genes and inspection of their evolutionary history. For each gene cluster, panX displays an alignment, a phylogenetic tree, maps mutations within that cluster to the branches of the tree and infers gain and loss of genes on the core-genome phylogeny. PanX is available at pangenome.de. Custom pan-genomes can be visualized either using a web server or by serving panX locally as a browser-based application.

INTRODUCTION

In addition to vertically passing down their genome to offspring, bacteria have the capability to acquire genetic material from the environment via horizontal transfer (1). Genes are transferred among bacteria by a variety of mechanisms including active uptake, mobile genetic elements, and gene transfer by viruses (2). In addition to gene gain, genes are frequently duplicated or lost. The mix of vertical transmission and horizontal transfer complicates the phylogenetic analysis of bacterial genomes and results in patterns of genetic diversity that are difficult to interpret (3).

A common approach when analyzing collections of bacterial genomes is categorizing genes into the *core* or *accessory* genome (4–6). Core genes are shared by all strains in a group of isolates, accessory genes shared by two or more but not all strains, and unique genes are specific to

a single strain. The union of all genes found in a group of strains (e.g. strains from one species) is called the pan-genome, which is typically several times larger than the core genome. The core genome is often used to assess the relatedness among the genomes in the sample and to approximate the species tree, but extensive horizontal transfer has been documented in the core genome as well (7) such that a tree reconstructed from core genome diversity does not necessarily reflect the phylogeny. Different software tools try to infer or remove the impact of recombination on the species level phylogeny (8,9).

By providing a repertoire of functional genes, gene gain from the pan-genome can facilitate the acquisition of new metabolic pathways (10), the adaptation to new habitats, or the emergence of drug resistant variants (11). With the rapidly increasing number of sequenced bacterial genomes, it is now possible to detect associations between metadata such as habitats, phenotypes, clinical manifestations and the presence or absence of particular genes (12–15).

Pan-genome construction from a group of related bacterial genomes typically involves the identification of homologous regions by all-against-all comparisons followed by clustering orthologous genes (4). Several software packages and pipelines have been developed to construct such pan-genomes that differ in the heuristics used to compare strains and generate clusters (16–19).

One fundamental limitation, however, is the difficulty to interrogate, explore, and visualize the pan-genome and the evolutionary relationships between strains. In absence of recombination, the purely vertical evolutionary history of strains would be represented by a single phylogenetic tree, the species tree. With horizontal transfer, the history of different loci in the genome is described by different trees resulting in a phylogenetic forest or network (20,21). While phylogenetic networks can be visualized using consensus representations such as split networks (22), the history and distribution of individual proteins are often critical, for example when searching for associations with phenotypes like drug resistance. Individual clusters of orthologous sequences, however, can again be represented by a tree if genes are short enough that recombination within the gene can be ignored. Some gene trees might be similar to the species tree, while others might deviate dramatically from the species

*To whom correspondence should be addressed. Tel: +41 612075834; Fax: +41 612072078; Email: richard.neher@unibas.ch

tree. The degree of incongruence of the gene tree with the species tree contains important information about the dynamics of gene gain and loss.

Here, we present panX, a web-based environment for microbial pan-genome data visualization and exploration based on an automated pan-genome identification pipeline. The pipeline breaks the genomes of a large number of annotated genomes (e.g. NCBI reference sequences) into genes and then clusters genes into orthologous groups. From these clusters, panX identifies the core genome, builds a strain-level phylogeny using SNPs in the core genome, constructs multiple alignments of sequences in gene clusters, builds trees for individual genes and maps the gene presence/absence pattern onto the core genome tree. The interactive browser-based application then allows the exploration of the above features and provides flexible filter, sort, and search functionalities. This application is available at pangenome.de with a collection of pan-genomes prepared by us, but can also be deployed on other servers with custom pan-genomes, or can be run locally as a browser-based desktop application.

MATERIALS AND METHODS

Identification of orthologous gene clusters

The initial steps in the computational pipeline underlying panX (illustrated in Figure 1) are broadly similar to other tools used to construct pan-genomes (16–19,23). PanX algorithm identifies groups of homologous genes by similarity search using DIAMOND and clustering by MCL. In a second step, panX builds phylogenies of these groups of genes and splits them into approximately orthologous clusters by examining the structure of the trees. PanX thus combines the speed of graph methods to identify groups of similar sequences with tree based methods applied to individual clusters to accurately split homologous sequences into orthologous groups.

Identification of groups of homologous sequences. As input, panX uses annotated genome sequences in GenBank format. To identify homologous proteins, panX performs an all-against-all similarity search using DIAMOND (24) with default *e*-value cut-off of 0.001. The DIAMOND similarity search can be multi-threaded and panX uses 64 CPUs by default if run on a compute cluster. From the diamond output, panX constructs a file with all pairs of genes with significant hits and the corresponding bitscore. Using bitscore instead of *e*-value avoids underflow problems and combines similarity and length of the homologous region (25). The table of similarity scores serves as input for the Markov Clustering Algorithm (MCL) (26,27) to create the clusters of putatively orthologous genes.

Since DIAMOND aligns proteins, ribosomal RNAs (rRNA) and other non-protein coding genes have to be handled separately. PanX extracts annotated rRNAs from annotated genomes (GenBank format) and compares sequences to each other via blastn. The output of blastn is then clustered by MCL in the same way as the protein comparison by DIAMOND. If desired, the DIAMOND similarity search can be replaced completely by other sequence similarity search tools such as blastx or blastn.

Divide-and-conquer strategy for large data sets. The all-against-all similarity search scales quadratically with the number of genomes, making the naive implementation infeasible for thousands of genomes, see Figure 1. However, the majority of these comparisons is redundant and can be avoided by first clustering small batches of genomes and subsequently combining different batches. Specifically, we apply the DIAMOND and MCL steps to subsets of 50 genomes (large enough to benefit from DIAMOND's double indexing strategy, small enough such that the all-against-all comparison is not yet prohibitive) and derive gene clusters of this “sub-pan-genome”. Each gene cluster is then reduced to a representative sequence and the representative sequences of all gene clusters are used as a “pseudo genome” representing the entire batch. The pseudo genomes representing the different batches are then again clustered using the DIAMOND + MCL steps. Eventually, complete clusters are constructed by combining sequences represented by the pseudo genomes. This “divide-and-conquer” strategy can be applied repeatedly for very large pan-genomes and keeps scaling of clustering approximately linear, see Figure 1.

Splitting into orthologous clusters. In our experience, it is advisable to cluster proteins aggressively and split clusters with paralogous sequences in a post-processing step. The groups of paralogs are often readily apparent in a phylogenetic tree. PanX reconstructs trees from sequences in each cluster by first aligning the protein sequences using MAFFT (28). The protein alignment is then used to construct a codon-alignment of the corresponding nucleotide sequences by inserting a gap of length three for every gap in the amino acid alignment. From nucleotide sequence alignment panX then reconstructs a tree using FastTree (29). The runtime of FastTree scales approximately as $n^{3/2}$ with the number of sequences. While superlinear, this scaling still allows the analysis of thousands of genomes. Tree building for the thousands of gene clusters can be trivially parallelized on all available CPUs. For 600 genomes, this step takes about twice as long as the initial clustering (see Figure 1). Once the tree of a cluster is available, panX employs a three-step procedure to decide whether and where a cluster should be split into sub-clusters.

Splitting of distantly related homologs. Since branch lengths reflect evolutionary distances among genes within one cluster, groups of distantly related genes, for example resulting from an ancient duplication, are connected by long branches and can be easily spotted in a gene tree—at least for pan-genome of low or moderate diversity. PanX splits trees into subtrees at branches whose length exceeds an adaptive threshold. This threshold is determined from the average diversity d_c of single copy core genes via

$$b_c = \frac{0.1 + 2d_c}{1 + 2d_c}. \quad (1)$$

This cut-off increases as $0.1 + 2d_c$ for very similar strains and eventually saturates at 1. The genetic diversity of the core genes will be of the same order of magnitude as mutational distance to the most recent common ancestor (MRCA) of the collection of genomes. Branches much

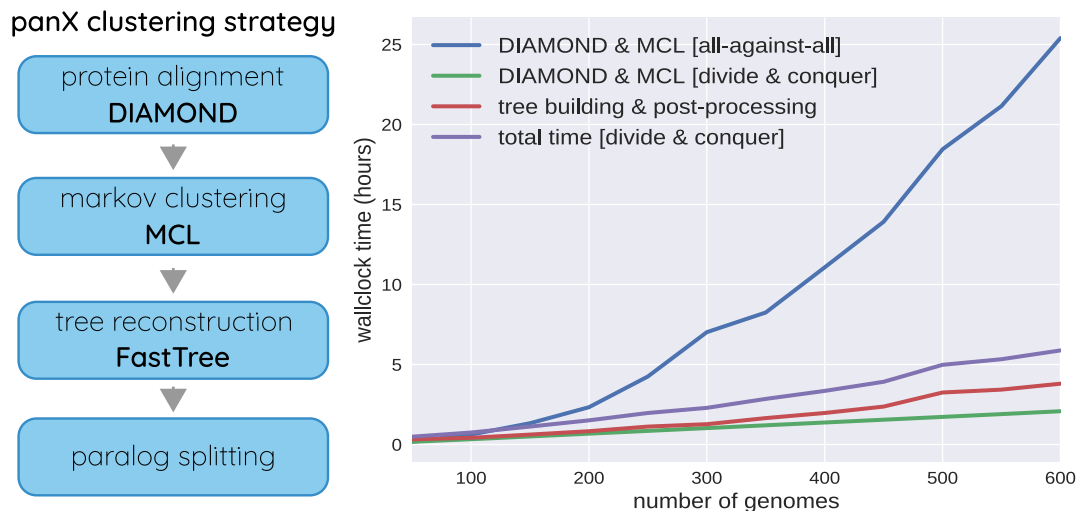


Figure 1. panX analysis pipeline. PanX uses DIAMOND (24) and MCL (26,27) to identify clusters of homologous genes from a collection of annotated genomes. These clusters are then analyzed phylogenetically and split into orthologous groups based on the tree structure. The graph on the right shows the time required to identify orthologous clusters in pan-genomes of different size on a compute node with 64 cores. The naive all-against-all comparison with DIAMOND scales quadratically with the number of genomes (blue line, ‘DIAMOND & MCL [all-against-all]’). The “divide and conquer” strategy where clustering is first applied to batches of sequences and batches are subsequently clustered (see text) reduces this scaling to approximately linear (green line). Tree building and post-processing take about as long as the clustering itself for pan-genomes of 500 genomes.

longer than this diversity will typically correspond to duplications long before the MRCA. Hence, the cluster should be cut along these branches. For very diverse pan-genome with $d_c > 0.25$, however, this simple threshold splitting will result in under-clustering and should be switched off or modified.

The newly formed clusters are then re-aligned, a new tree is built, and further split using above-mentioned method until all long branches have been cut.

Splitting of closely related paralogs. Splitting branches longer than b_c will miss recent gene duplication events. To detect paralogous clusters more sensitively, panX calculates a paralogy score for each branch in each tree. The paralogy score of a branch is the number of strains represented on both sides of the branch in the phylogeny. This score can be calculated in linear time for all branches simultaneously in two tree traversals. Clusters are then split into two sub-clusters if the highest paralogy score ϕ_{max} and the length ℓ of the corresponding branch fulfill the following criteria:

$$\begin{aligned} \phi_{max} &> 0 & (2) \\ \frac{\ell}{b_c} + \frac{\phi_{max}}{1.5 \times \#strains} &> 1.0 \end{aligned}$$

Here, b_c refers to the cut-off defined in Eq. (1). The criterion $\phi_{max} > 0$ prevents splitting irrespective of paralogy. Other criteria could be used but in our experience, this linear discriminator works well for many different applications.

This paralog splitting is iterated until no gene cluster meets the condition (2). Some heavily duplicated genes require more than five rounds of splitting. The parameters of the condition (2) can be set by the user.

Merging of fragmented clusters. A small number of genes are not properly clustered either because no homology was detected initially or the clustering by MCL failed. Such unclustered sequences manifest themselves as many singleton

clusters of identical length. To detect those sequences, panX calculates the average length of sequences in each cluster and searches for peaks in the distribution of gene cluster length. Unclustered genes show as spikes in this empirical gene length distribution, which panX can identify by detecting peaks in this distribution relative to a smoothed background distribution. For each detected peak, all involved genes are gathered in one pre-cluster, their sequences are aligned and the corresponding phylogeny is inferred. These clusters are then processed as described above. Currently, only sequences of identical length are combined into tentative clusters. This condition could be relaxed but this has not been necessary in our experience.

The number of clusters that require post-processing depends on the diversity of the data set. On simulated data, ~40% of initial clusters needed splitting for the least diverse sets, while only a small fraction of clusters required post-processing for the more diverse data sets.

Phylogenetic analysis of gene clusters

To reduce the computational burden of the subsequent visualization of the pan-genome, alignments, trees and other properties of the gene clusters are precomputed. The input for this phylogenetic analysis is either the output of the panX pipeline presented above, or the output of Roary. Other pan-genome tools could be used when a script parsing the clustering output is supplied.

Tree building and ancestral reconstruction. PanX extracts all variable positions from the nucleotide alignments of all single copy core genes (those gene clusters in which all strains are represented exactly once) to construct a core-genome SNP matrix. This SNP matrix is used to build a core genome phylogenetic tree using FastTree (29), which is further refined by RaxML (30) following a similar strat-

egy as implemented in nextflu (31). Due to homologous recombination, this core genome tree may not reflect the true history for each of the genes in the core genomes (7) and branch lengths do not reflect sequence similarity since only variable sites are used (32). Nevertheless, this core genome SNP phylogeny is still a useful approximation of the relationships of the different strains that can be used as a scaffold to investigate the evolution of the mobile genome and the distribution of phenotypes.

Phylogenetic trees for a gene cluster have already been inferred in the cluster post-processing step. PanX uses these trees to infer ancestral sequences of internal nodes using a joint maximum likelihood approach (33) as implemented in treetime github.com/neherlab/treetime. Likely mutations are mapped onto the branches of the tree using this ancestral reconstruction.

Then, we infer the presence or absence of each gene cluster on internal nodes of the core genome SNP tree using an analogous ancestral inference procedure. Individual gain and loss events are associated with branches based on this ancestral reconstruction. The gain and loss rates are optimized such that the likelihood for the observed presence/absence pattern of genes is maximized (34,35). We found that optimal loss rates are always larger than the gain rates but their ratio is variable among species with a median ratio of 22 (inter-quartile range 9–35).

Gene clusters, trees, mutations, and metadata are stored as JSON files for the web visualization.

Associations

If informative numerical meta data are attached to the genomes, panX can quantify the association of genetic signatures with these meta data. PanX considers associations of two types: Either particular variants of a gene are associated with a phenotype, or the presence/absence of a gene can be linked to a phenotype.

To quantify how well a phenotype is associated with particular variants of a gene, panX computes a normalized difference of phenotypes of strains on either side of a branch in the tree as follows

$$b_a = \frac{\mu_1 - \mu_2}{\sqrt{\sigma_1^2 + \sigma_2^2}} \quad (3)$$

where $\mu_{1/2}$ and $\sigma_{1/2}^2$ are the mean and variance of the phenotypes of either side of the branch, respectively. This score is calculated for every branch of the tree and the maximum score is reported.

To quantify associations of a phenotype with the presence and absence of a gene, panX uses the average phenotype μ_p of strains carrying the gene and the average μ_a of strains without the gene, the overall variance of the phenotype σ^2 , and the number of gain/loss events n to calculate the score

$$p_a = \sqrt{n} \frac{\mu_a - \mu_p}{\sigma} \quad (4)$$

Simple scores like the above won't reliably separate true associations from all false positives since that would require careful population structure correction (13–15). However,

we will discuss below that these scores recapitulate previously described associations very well. In the panX web application, genes can be sorted by association score such that strongly associated genes can be rapidly found. By inspecting the species and gene tree, colored by the corresponding meta-data, the user can get a visual impression of the degree, consistency, and population structure confounding of the association.

Simulation of pan-genomes

To assess the accuracy of clustering methods we simulated 120 pan-genomes with 30 artificial genomes each. The simulation evolves ancestral sequences along coalescent trees generated by the software ms (36) and allows for horizontal transfer as well as gene loss and gain. To get realistic ancestral sequences we used one representative gene from each KEGG ortholog group (37) present in the *Escherichia coli* strain K-12 (NC_000913) as a starting point for the simulation. This yielded 2803 different genes.

Using these as ancestral sequences, we simulated pan-genomes by the following procedure: For each of the 2803 genes we generated correlated trees using the software ms (36) with different rates of horizontal transfer. If the gene transfer rate is zero, all 2803 genes evolve according to the same clonal genealogy of the population, i.e. one common species tree. By contrast, the individual gene trees may differ if some genes are affected by gene transfer. Nonetheless, the gene trees are still strongly dependent on each other due to the common link to the clonal genealogy. To investigate the effect of transfer on the accuracy of reconstruction, we used three different rates of gene conversion for the simulation of gene trees with ms (option `-c` with values 0, 2000 and 4000 with 6000 potential sites for gene conversion).

Among 2803 genes, 2100 were assigned to the most recent common ancestor (MRCA) at the root of the simulated gene tree while the remaining genes are gained at uniformly distributed points at the branches after the MRCA on the gene tree. 300 of the 2100 ancestral genes were assigned to be present at all times, the remaining 2503 genes were lost at rate 2.1 along the branches of the corresponding gene tree as defined in (38,39). After a loss event, the corresponding gene is absent from all individuals descending from the branch of the loss event.

Given the gene trees and the presence absence pattern for each representative K-12 sequence, substitution can occur along the branches of the corresponding gene tree. We used seq-gen (40) to simulate these mutations according to the HKY model (41), setting the base frequencies to empirical *E. coli* base frequencies and the transition-transversion bias to 1.1. We simulated 5 sets of gene trees for no, occasional and frequent gene conversion. For each set we used 8 different substitution rate distributions to simulate pan-genomes: an exponential distribution with mean 0.06, uniform distributions between 0.05 and 0.1 and between 0.1 and 3, and constant substitution rates of 0.3, 0.2, 0.1, 0.05 and 0.01. The substitution rate μ of each gene was drawn from the corresponding distribution. The mean number of substitutions per site between two strains is given by $1 - e^{-\mu T}$, where T is the distance between both strains in the gene tree.

RESULTS AND DISCUSSION

Benchmarking and comparison to other tools

To compare the clustering performance of different methods, one has to know which genes belong to the same cluster of orthologs. However, the *orthoBench* collection of manually curated groups of orthologous proteins (42) has been used to compare very diverged proteins across different domains of life and is far too diverse to benchmark a tool meant for pan-genomes of closely related bacteria. The orthologous groups in inferred pan-genomes depend on the software used for pan-genome inference and the ground truth is unknown. While pan-genomes based on real genomes can be used to compare clustering methods against each other, they are not immediately useful to assess accuracy.

To evaluate the performance of panX in an absolute sense and in comparison to state-of-the-art tools Roary (16), OrthoMCL (43), PanOCT (23) and OrthoFinder (44), we generated simulated pan-genomes for which the ground truth is known. In addition, we investigated the consistency of the orthologous clusters between tools in pan-genomes constructed from real bacterial genome sequences.

OrthoMCL and OrthoFinder were designed to identify orthologous groups across different domains of life, not for bacterial pan-genome inference. Nonetheless, they are often used in this context and we therefore included them here. By contrast, Roary was designed to cluster very large numbers of similar genomes. These tools are therefore expected to work well in different parameter ranges.

As a unique feature among the tools we compared, panX relies on phylogeny based post-processing of the initial MCL clustering. This post-processing step is adaptive in the sense that thresholds are scaled relative to the core genome diversity. As a result, panX works well across a large range of diversities.

Comparison of clustering accuracy on simulated datasets. We constructed pan-genomes by evolving 2803 genes from the *E. coli* genome along gene trees generated by the software ms (36). Gene sequences evolve along the gene trees with a gene specific mutation rate, can be gained or lost, and undergo horizontal transfer, see Materials and Methods for details.

We subjected 40 simulated pan-genomes of size 30 to analysis by Roary, PanOCT, OrthoFinder, OrthoMCL and panX and compared the inferred orthologous clusters to the actual clusters generated by the simulation. For each cluster, there are four possible outcomes: (i) a cluster is correct if it contains all and only genes from one true cluster, (ii) a cluster is incomplete but contains only genes from one true cluster, (iii) a cluster contains all genes from one true cluster but also genes from other clusters, or (iv) a cluster could fail on both counts. Figure 2 shows how different tools perform at different levels of diversity of the pan-genome.

OrthoMCL and OrthoFinder were designed for cross-species comparisons at large evolutionary distances. It is hence not surprising that when analyzing low diversity pan-genomes these two tools merged many clusters that should have been kept separate. This effect is most pronounced for rare clusters, predominantly singletons, that get combined

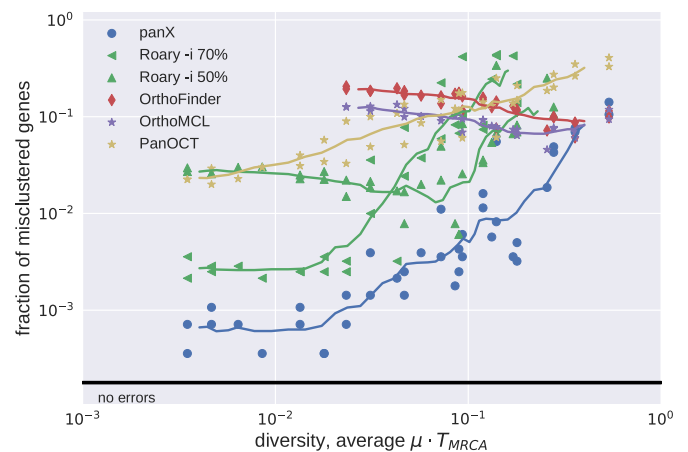


Figure 2. Accuracy of clustering by different tools. The fraction of misclustered genes increased with diversity of the pan-genome. We ran Roary with options *-i* 70 and *-i* 50. At low diversity, panX and Roary (*-i* 70) performed with similar accuracy and mis-clustered about 1 in 1000 genes. At high diversities, all tools showed similar accuracy and mis-clustered 1 in 10 genes. Results for tools designed for high diversity data sets (OrthoMCL and OrthoFinder) are only shown for diversities above 0.02. Similarly, results for Roary are suppressed at high diversity to improve clarity of the graph.

with other clusters, see Figure 3. Core genes and other common gene clusters are typically correctly reconstructed. At very large diversities, we found that OrthoMCL and OrthoFinder had an accuracy similar to that of panX. Roary and panX showed similar behavior across a wide range of diversities from below 1% to 30%, with panX typically making a factor of two fewer mistakes. However, we were unable to find a parameter set for Roary that worked well across the entire range of diversities. Using an identity cutoff of 70% (*-i* 70) worked best at low diversity, while lower cutoffs were required at high diversity. PanOCT didn't perform very well on our simulated data predominantly because it split too many clusters, see Figure 3.

The different types of errors (erroneous merging/splitting) are shown separately in Supplementary Figure S1. We repeated the analysis for different gene conversion rates and found mainly comparable results for no, occasional and frequent gene conversions (Supplementary Figures S2 and S3).

To test the divide-and-conquer strategy, we generated an additional set of simulated pan-genomes comprising 500 strains (exponentially distributed substitution rates with mean 0.06 per coalescent time scale). We analyzed this pan-genome with panX using divide-and-conquer, panX without divide-and-conquer, and Roary. PanX correctly identified 2780 out of 2803 gene clusters (23 errors), while panX made four additional errors when switching the divide-and-conquer strategy off. For comparison, Roary mis-clustered 75 genes when running on this set of genomes with parameter *-i* 50 (146 false gene clusters with parameter *-i* 70). We conclude that divide-and-conquer does not reduce the clustering accuracy significantly. Additional tests on a large *S. pneumoniae* data set are presented in Supplementary Figure S4.

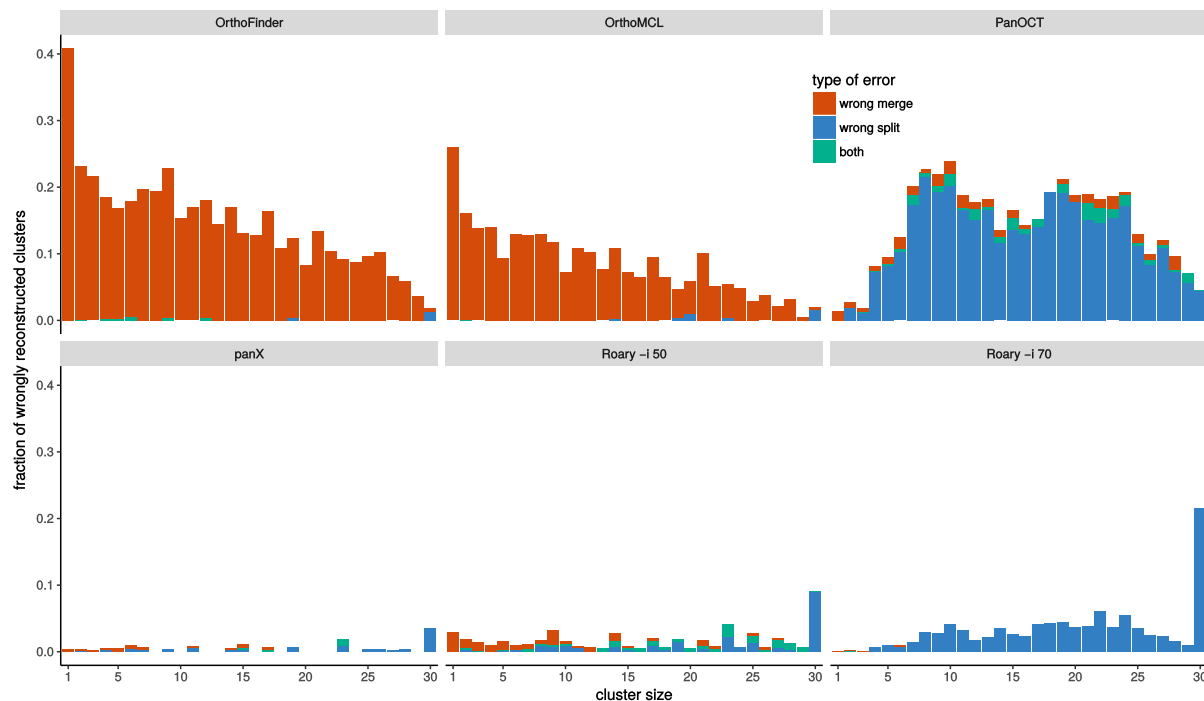


Figure 3. Type of mis-clustering by tool and gene frequency. The histograms show the fraction of wrongly merged (red) and wrongly split (blue) clusters by gene frequency and clustering tool across 5 simulated datasets with exponentially distributed substitution rates with mean rate $\mu = 1/15$.

Due to the gene-centric alignment and clustering strategy, panX is insensitive to incomplete and fragmented genomes, as is common when genomes are assembled from short reads. We tested this insensitivity explicitly by removing 10, 20 or 30% of genes from 10, 20 or 30% of all strains, respectively, as well as removing 10% of genes from all strains. Supplementary Figure S5 demonstrates that incomplete genomes have negligible effect on panX's accuracy when tested on our simulation data.

Real pan-genomes. While the ground truth of simulated pan-genomes is known, real pan-genomes lack an obvious point of comparison. Nonetheless, a comparison of the orthologous clusters identified by different approaches can highlight the similarity and differences between clustering strategies.

We used *S. pneumoniae* and *Prochlorococcus* pan-genomes to compare the results of the panX ortholog identification pipeline to that of Roary, OrthoMCL, PanOCT and OrthoFinder. We computed the size distribution of clusters, the size of the core genomes, and the total number of clusters from the pan-genomes estimated by different tools, see Figure 4. For the low diversity collection of 33 *S. pneumoniae* genomes from the RefSeq database, the cluster size distributions inferred by the different tools are very similar (Figure 4A) and the number of inferred core genes differs by less than 10%. Between 78 and 86% of clusters inferred by panX are found by other tools, see Figure 4B. The greatest overlap is with Roary when run with identity threshold $-i\ 90$.

More variation is observed in inferred pan-genomes of 40 *Prochlorococcus* strains, see Figure 4C and D. Roary (16) separates nearly all *Prochlorococcus* genes and identi-

fies only 10 core genes when using standard parameters. After lowering the minimum percentage identity for blastp in Roary to 30% ($-i\ 30$), Roary identified 1111 core genes vs. 1214 identified by panX. While Roary warns that it has not been designed to support such diverse datasets, 60% of the resulting clusters agree with those identified by panX. PanX and Roary identified 5407 and 6981 clusters of orthologous genes, respectively—not too far from the estimated pan-genome size of >8500 genes present in more than one percent of the population (45). The results for OrthoFinder and OrthoMCL are comparable to those of Roary and panX. In contrast, the tool PanOCT separates many more genes than panX or Roary with parameter $-i\ 30$. PanOCT is designed for closely related prokaryotic strains and therefore splits the diverse *Prochlorococcus* genomes into 16 820 clusters with mainly 1–3 genes. Only 109 core gene clusters are identified. The tools OrthoMCL and OrthoFinder, designed for more diverse data sets, generate cluster size distributions similar to those by panX and between 71 and 87% of clusters found by one tool are also found by another.

Testing this collection of pan-genome tools on larger data sets proved prohibitive since only Roary and panX scale well with number of genomes.

Web application for pan-genome exploration

To explore the pan-genome constructed by the pipeline described above, we developed a browser based visualization application. The layout of the application is that of a large dashboard (see Figure 5), on which multiple aspects of the pan-genome can be interrogated simultaneously.

At the top, three graphs provide basic statistics on the abundance and length distribution of all genes. A search-

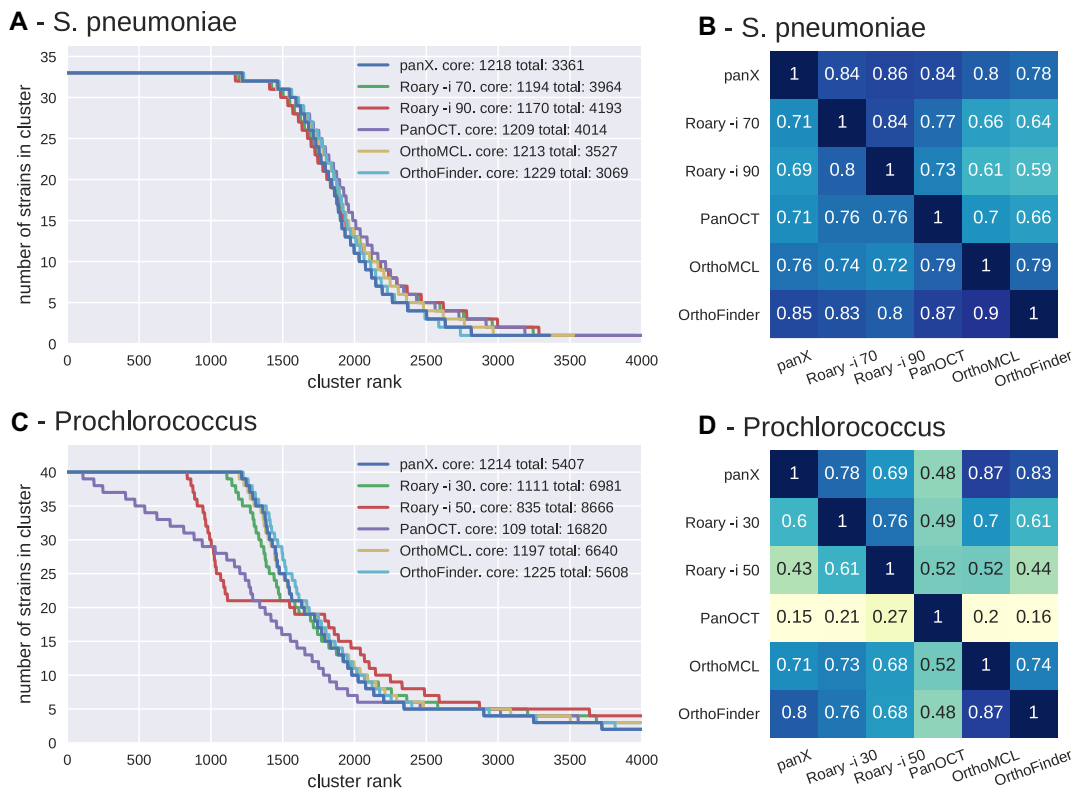


Figure 4. Pan-genome statistics. Panels A and C show the distribution of the number of strains represented in pan-genomes of 33 *S. pneumoniae* and 40 *Prochlorococcus* strains constructed by panX, Roary, OrthoFinder, OrthoMCL and PanOCT (the last two tools are only available for the smaller *Prochlorococcus* data set). To obtain these graphs, clusters are sorted by descending number of strains represented in the cluster. This number is then plotted against the rank of the sorted clusters. The point where the lines drop below the number of strains marks the size of the core genome. PanX, OrthoFinder and OrthoMCL largely agree on the cluster size distribution, the number of core genes and the total size of the pan-genome (with ~10% variation). Roary agrees with the latter tools if the identity cut-off is chosen appropriately, while PanOCT estimates a very small core genome and an extremely large number gene clusters. Panels B and D show the degree to which different pan-genome tools agree with each other. Each row shows the fraction of clusters identified by one tool, which exactly match the clusters identified by another tool. Analogous results for simulated data are given in Supplementary Figure S6.

able table contains summary statistics and annotations for all gene clusters. Below the table, the core genome SNP tree is shown, along with a phylogenetic tree of the currently selected gene cluster. A second searchable table below the trees allows rapid access to meta information available for different strains and can be used to select strains in the tree. An alignment viewer shows the nucleotide or amino acid alignment of the gene cluster selected in the table. The layout of the difference components can be easily rearranged. In our standard layout, the alignment viewer is next to the metadata table. If, however, many data points are to be displayed in the gene cluster table, the table extends across the entire row and the alignment viewer is placed at the very bottom of the page.

The hallmarks of the panX web-application are the interconnected components that illustrate different properties of the gene clusters. The pan-genome statistic charts at the top allow rapid sub-setting of gene clusters by gene length and abundance. The left chart shows an inverse cumulative distribution of clusters sizes, i.e. clusters are sorted by decreasing number of strains represented in the cluster, such that all core genes present in all strains are shown on the left. The size of the core genes is then simply the length of the plateau of the curve to the first drop. The core genome is

followed by gradual decline in gene number from common to rare accessory genes. Lastly, a long tail contains strain-specific singletons. Subsets of genes can be easily defined by selecting a range of the graph with the mouse. Similarly, the center chart shows the distribution of gene length.

The pie chart on the right shows the proportion of core and accessory genome, each of which can be selected by clicking on the sectors in the chart. To allow for soft and strict definitions of the core genome, the cut-off delineating core and accessory genome can be adjusted with a slider.

Rapid and searchable access to alignment and gene trees. The table of all gene clusters is dynamically restricted to the range of gene abundances and gene lengths selected above. The table can be searched by gene name and annotation or sorted by gene count, diversity etc. Annotations of all input sequences (also discordant annotations of genes belonging to the same gene cluster) are accessible by expanding the annotation field. Similarly, the column *duplicated* specifies whether the gene cluster contains more than one gene per strain. The list of strains in which genes are duplicated and copy number of this gene can be accessed by expanding the row. Each row contains triggers to show the corresponding nucleotide or amino acid sequence alignment in the align-



Figure 5. Interconnected components of the panX web application. The top panels provide a statistical characterization of the pan-genome and allow filtering of gene clusters by abundance and gene length. The gene cluster table below is searchable and sortable and allows the user to select individual gene clusters for closer inspection. Upon selection in the table, the alignment of gene cluster is loaded into the viewer on the center right, the gene tree is loaded into the tree viewer at the bottom right, and presence/absence patterns of this gene cluster are mapped onto the core genome tree at the bottom left. The example shows the gene coding for the penicillin binding protein Pbp2x and the color indicates the susceptibility to benzylpenicillin.

ment viewer (MSA) from BioJs (46). In order to highlight difference among sequences, only consensus sequence and variable sites are shown by default, while the corresponding original alignment can be downloaded. This trigger also updates the phylogenetic tree viewers. Searching *mcr-1*, for example, immediately highlights the 11 *E. coli* genomes in the RefSeq database that have an annotated mobile colistin resistance gene.

Interactive core genome tree and gene tree viewers. To facilitate the comparison between the core genome SNP tree and the gene tree, the two trees have connected interactive elements. When placing the mouse on a leaf node in one tree, the corresponding nodes are highlighted in both trees. Similarly, if the mouse is placed over an internal node, all nodes in the corresponding clades are highlighted with different colors for each strain. This gives a rapid impression of whether the core genome tree and the gene tree are compatible and whether the gene is duplicated in some of the strains, see Figure 6.

The most likely gene loss and gain events inferred by the ancestral reconstruction algorithm are indicated on the tree by dashed or thick lines, respectively. Mutations in the amino acid or nucleotide sequence of the gene are mapped onto the gene tree and can be inspected using the tooltips associated with branches in the tree.

In addition to mutation and gain/loss events, the tree can be colored with metadata associated with different strains. Such metadata would typically include collection dates, sampling location, host species or resistance phenotypes.

Pan-Genomes of common bacterial groups

We ran the panX on collections of genomes of bacterial groups for which more than 10 genomes were available in RefSeq, resulting in a total of 94 pan-genomes including many human pathogens. Statistics of a subset are shown in Table 1, the corresponding data for all species are given as supplementary material. The majority of these species exhibit low diversity in their core genomes with typically just a few percent nucleotide differences, sometimes even <0.001 . The median number of core genes is 1800 while the median pan-genome size is ~ 5000 genes. Diverse species tend to have smaller core-genomes and larger pan-genomes, as expected.

Pan-Genomes of diverse collections of genomes

Most of these collections are closely related genomes, but we also included a diverse group of genomes of *Prochlorococcus* and the pan-genomes of the bacterial orders Pseudomonadales, Enterobacteriales and Vibrionales.

Prochlorococcus is a marine cyanobacterium that is responsible for a significant fraction of the marine pri-

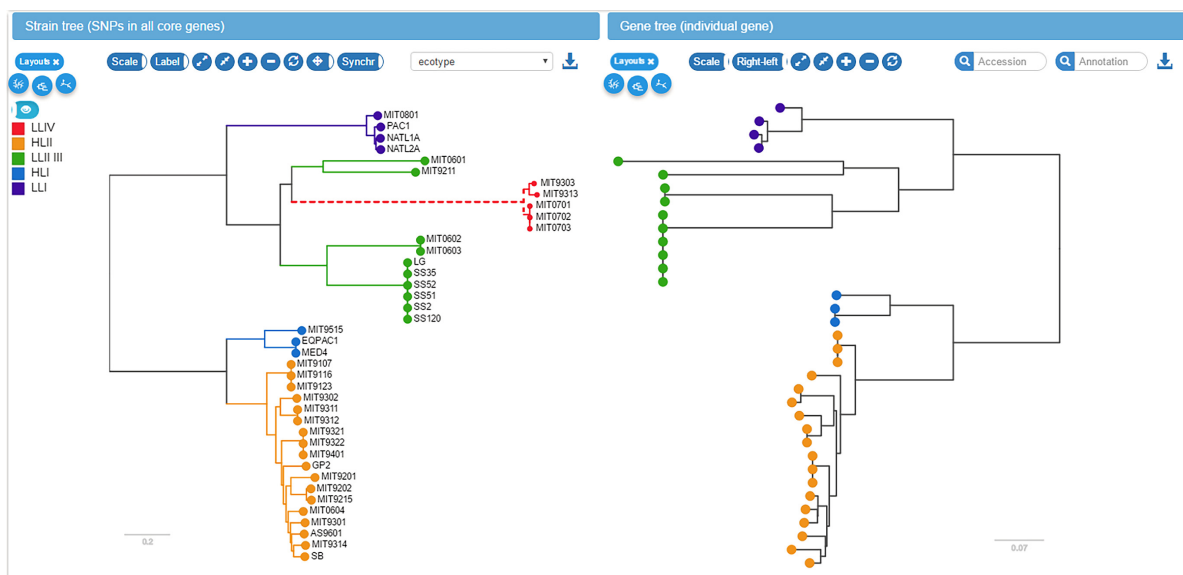


Figure 6. Linked core genome and gene trees. The core genome tree shows the strains in which the current gene is present or absent. Placing the mouse over an internal node in one of the trees (upper clade of the gene tree on the right in this example) highlights all strains in the corresponding clade in both trees. This gives the user a rapid impression of phylogenetic incongruence and likely gene gain and loss events.

Table 1. Summary statistics of pan-genomes available at pangenome.de

Species	Genomes	Core genes	All genes	Singletons	Diversity*
<i>Acinetobacter baumannii</i>	71	1701	8334	1558	0.010
<i>Bacillus anthracis</i>	43	4156	5980	62	1.0e-04
<i>Bacillus cereus</i>	36	2979	13364	3486	0.048
<i>Bordetella pertussis</i>	291	2437	3743	158	4.1e-06
<i>Burkholderia pseudomallei</i>	59	4098	11580	1966	0.003
<i>Campylobacter jejuni</i>	113	935	3166	526	0.014
<i>Chlamydia trachomatis</i>	68	809	978	12	0.005
<i>Clostridium botulinum</i>	23	795	9083	2294	0.147
<i>Corynebacterium pseudotuberculosis</i>	59	1133	2316	65	0.005
<i>Enterobacter cloacae</i>	22	2971	10783	3211	0.087
<i>Escherichia coli</i>	307	778	23107	6339	0.015
<i>Francisella tularensis</i>	35	838	2339	302	0.007
<i>Helicobacter pylori</i>	85	694	2371	328	0.042
<i>Klebsiella pneumoniae</i>	109	2545	15978	4004	0.007
<i>Listeria monocytogenes</i>	95	1907	4947	485	0.031
<i>Mycobacterium tuberculosis</i>	51	2665	4350	93	2.0e-04
<i>Neisseria meningitidis</i>	78	1071	3375	426	0.015
<i>Prochlorococcus marinus</i>	40	1047	5407	1262	0.291
<i>Pseudomonas aeruginosa</i>	70	3264	12768	3195	0.006
<i>Salmonella enterica</i>	260	1327	15521	3996	0.009
<i>Staphylococcus aureus</i>	146	1229	5206	731	0.008
<i>Streptococcus pneumoniae</i>	33	1188	3361	540	0.010
<i>Streptococcus pyogenes</i>	50	970	2856	341	0.008
<i>Vibrio cholerae</i>	28	2412	5156	771	0.005
<i>Xanthomonas citri</i>	26	3385	5261	291	0.001
<i>Yersinia pestis</i>	33	2557	4587	172	1.0e-04
Pseudomonadales	119	966	42520	20577	0.194
Enterobacteriales	33	1998	16413	6988	0.112
Vibrionales	66	716	30461	15643	0.193

*Average number of pairwise differences per nucleotide in core gene alignments.

mary production and serves as a model system in marine microbial ecology (47). While we relied on annotations available in NCBI for most species, we re-annotated the genomes of 40 *Prochlorococcus* sequences (48) using Prokka (49). The annotation was derived from a custom database based on the 12 annotated *Prochlorococcus* strains CCMP1375, MED4, MIT9313, NATL2A,

MIT9312, AS9601, MIT9515, NATL1A, MIT9303, MIT9301, MIT9215 and MIT9211. *Prochlorococcus* is a much more diverse population than the other species we investigated, see Table 1, which makes it a challenging case for pan-genome analysis.

While the 16S rRNA sequences of all 40 *Prochlorococcus* strains do not differ by more than 3%, *Prochlorococcus*

can be divided into ecotypes that are remarkably different in genome size and GC content (47). These ecotypes correspond to high and low light intensity adapted *Prochlorococcus* populations and can be visualized along the species and gene trees in panX. While the genomes of *Prochlorococcus* have likely been streamlined by strong selective forces to lose genes (50), gene gains and duplications have frequently occurred in all *Prochlorococcus* lineages. Two well-known examples for *Prochlorococcus* are the gain of nitrate assimilation genes *nirA* (51) and gene duplication and phage mediated gene transfer of the photosynthesis gene *psbA* (52). PanX highlights ancestral gene gain or loss events on the species tree. The gain/loss history of each gene cluster can be investigated by selecting the gene in the gene cluster table choosing 'Gene presence/absence' as metadata coloring.

In addition to the diverse *Prochlorococcus* genomes, we analyzed collections of genomes that encompass the entire bacterial orders *Pseudomonadales*, *Enterobacteriales*, and *Vibrionales*. For each of these orders, we collected at most 10 genomes from each species (based on the species designation in the RefSeq files) to avoid over-representation of human pathogens. Running panX on these orders resulted in ~1000 core genes for *Pseudomonadales* and *Vibrionales* and about 2000 core genes in case of the smaller collection of *Enterobacteriales*. Core genes of *Pseudomonadales* and *Vibrionales* were typically 20% diverged from each other, while *Enterobacteriales* core genome was less diverse at 11%. The core genome SNP tree of the *Vibrionales* clearly separates the genomes by species, while the core genome trees of the other orders show considerable mixing of species.

Large pan-genome of *Streptococcus pneumoniae*

The utility of the interactive web application is most evident for collections of genomes with rich meta data. One such collection is the *S. pneumoniae* data set generated by Croucher *et al.* (53). This data set consists of 616 whole genome sequences and rich meta data including antibiotic susceptibility and host characteristics. The pan-genome inferred by panX consists of a total of 4241 gene clusters with 34% core genes, consistent with previous analysis (54).

For *S. pneumoniae*, we calculated branch and presence/absence association scores for every gene cluster. All scores are included in the gene cluster table, which can be sorted by each score. The branch association score quantifies the degree to which a branch in the gene tree separates isolates with low and high phenotypes. We compared the genes which panX ranks as highly associated with benzylpenicillin resistance with beta-lactam resistance associations reported by Chewapreecha *et al.* (14). Their analysis revealed >2000 associations, many of which were deemed false positives by the authors. The associations discussed in the paper as biologically plausible are SNPs in *pbp2x*, *pbp1a*, *pbp2a*, *mraY*, *mraF*, *ftsL*, *gpsB*, *recU*, *clpL*, *clpX*, *dhfR*. Five of them (*pbp2x*, *pbp1a*, *mraY*, *gpsB*, *recU*) are among the six most associated genes according to the panX branch association score. The coloring of the gene tree by the benzylpenicillin susceptibility confirms that the resistant and susceptible isolates form two distinct clades separated by a large number of amino acid substitutions,

see Figure 5. While resistant strains are scattered across the species tree, they form a single clade in the tree of *pbp2x*.

Mutations in the dihydrofolate reductase gene (*dhfR/dyr*), along with dihydropteroate synthase (*folP*), are known to confer trimethoprim resistance rather than beta-lactam resistance and indeed these two are the two genes most strongly associated with trimethoprim resistance according to the panX score.

Similarly, the gene cluster table can be sorted by the degree to which presence/absence of a gene is associated with low/high phenotypes. Susceptibilities to benzylpenicillin, trimethoprim, erythromycin, and ceftioxone are strongly associated with the presence of *mefE* and *mel*, as expected for an efflux pump. After a cluster of seven genes with identical association patterns as *mefE* and *mel* (55), the next strongest association with erythromycin resistance is *ermB*. Again, acquisition of *ermB* is well known to confer macrolide resistance (56).

AVAILABILITY

The computational pipeline to identify the pan-genome consists of a collection of python scripts and a master script that runs desired analysis steps in series. The visualization is built on node.js server and makes extensive use of BioJS (46), D3.js (57), dc.js (58), and other javascript libraries. The analysis pipeline and the code for the web application is made available under the GPL3 license on [github](#) as repositories *pan-genome-analysis* and *pan-genome-visualization*.

The web application can either be hosted on a web server or can be used locally to inspect and explore pan-genomes produced by the panX pipeline. We computed a large number of pan-genomes and made those available at [pangenome.de](#). The website currently hosts 93 bacterial species including those listed in Table 1. Several downloading options are available: core gene alignments and all gene alignments can be downloaded via the down-arrow button next to the gene cluster table. Alignments and gene trees for individual clusters and the strain tree can be downloaded via buttons next to the alignment viewer and tree viewers, respectively. All buttons are associated with tooltips explaining the action of the buttons.

CONCLUSIONS

Being able to visualize and explore high dimensional data is often the key to developing insight into the mechanisms driving complex dynamics. PanX is meant to enable such exploration of large sets of bacterial genomes, which are characterized by the evolution of individual genes as well as the gain and loss of genes. The design of panX focused on combined breadth and depth: besides summary statistics and species trees, panX allows selecting interesting sets of genes or searching for individual genes. Alignments and phylogenetic trees of genes can then be analyzed in detail with individual mutations and gain/loss events mapped to the gene tree and the core tree, respectively. The evolutionary patterns of genes can then be compared to meta-information such as resistance phenotypes associated with the individual strains.

By integrating meta-information with the molecular evolution of genes and genomes in one visualization, panX can

assist investigations of the dynamics of pan-genomes and adaptation of bacteria to new habitats and environmental challenges. Horizontal transfer is pivotal for many aspects of bacterial adaptation (11), but at the same time, it remains much more difficult to analyze than evolution by vertical descent (3). The ability to interactively explore such pan-genomes might help to grasp the complexity of this dynamics.

On the other hand, a web-based tool that can be readily kept up-to-date by addition of newly sequenced isolates would be useful in pathogen surveillance. When paired with meta-information such as resistance, pathogenicity, sampling date, location and comorbidities, panX can help to study adaptation, spread, and transmission chains of pathogens. Similar approaches have proved useful at tracking spread and evolution of seasonal influenza virus or Ebola virus during the recent outbreak in West Africa (31,59). Currently, phenotype data are available for a minority of the whole genomes sequences and data sets like the *S. pneumoniae* by Croucher *et al.* (53) are an exception. With increasing availability and timely publication of such data from routine surveillance, panX or derivatives could be used to track food-borne outbreaks, monitor the global spread of drug resistance bacteria (60), or assist infection control in individual hospitals. The time required to build a pan-genome of 1000 strains is less than a day on a 64 core node such that frequent updates of such a tracking tool are possible.

SUPPLEMENTARY DATA

Supplementary data are available at NAR Online.

ACKNOWLEDGEMENTS

We gratefully acknowledge stimulating discussions with Matthias Willmann, Erik van Nimwegen, and Talia Karasov as well as advice on dc.js from Gordon Woodhull.

FUNDING

Max Planck Society (to W.D. and R.A.N.); University of Basel (to R.A.N.); DFG through the priority program [SPP1590 to F.B.]. Funding for open access charge: University of Basel.

Conflict of interest statement. None declared.

REFERENCES

1. Soucy,S.M., Huang,J. and Gogarten,J.P. (2015) Horizontal gene transfer: building the web of life. *Nat. Rev. Genet.*, **16**, 472–482.
2. Thomas,C.M. and Nielsen,K.M. (2005) Mechanisms of, and barriers to, horizontal gene transfer between bacteria. *Nat. Rev. Micro.*, **3**, 711–721.
3. Puigbò,P., Lobkovsky,A.E., Kristensen,D.M., Wolf,Y.I. and Koonin,E.V. (2014) Genomes in turmoil: quantification of genome dynamics in prokaryote supergenomes. *BMC Biol.*, **12**, 66.
4. Vernikos,G., Medini,D., Riley,D.R. and Tettelin,H. (2015) Ten years of pan-genome analyses. *Curr. Opin. Microbiol.*, **23**, 148–154.
5. Lapiere,P. and Gogarten,J.P. (2009) Estimating the size of the bacterial pan-genome. *Trends Genet.*, **25**, 107–110.
6. Tettelin,H., Riley,D., Cattuto,C. and Medini,D. (2008) Comparative genomics: the bacterial pan-genome. *Curr. Opin. Microbiol.*, **11**, 472–477.
7. Everitt,R.G., Didelot,X., Batty,E.M., Miller,R.R., Knox,K., Young,B.C., Bowden,R., Auton,A., Votintseva,A., Larner-Svensson,H. *et al.* (2014) Mobile elements drive recombination hotspots in the core genome of *Staphylococcus aureus*. *Nat. Commun.*, **5**, 3956.
8. Croucher,N.J., Page,A.J., Connor,T.R., Delaney,A.J., Keane,J.A., Bentley,S.D., Parkhill,J. and Harris,S.R. (2015) Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.*, **43**, e15.
9. Didelot,X. and Wilson,D.J. (2015) ClonalFrameML: efficient inference of recombination in whole bacterial genomes. *PLOS Comput. Biol.*, **11**, e1004041.
10. Fournier,G.P. and Gogarten,J.P. (2008) Evolution of acetoclastic methanogenesis in methanosarcina via horizontal gene transfer from cellulolytic clostridia. *J. Bacteriol.*, **190**, 1124–1127.
11. Ochman,H., Lawrence,J.G. and Groisman,E.A. (2000) Lateral gene transfer and the nature of bacterial innovation. *Nature*, **405**, 299–304.
12. Méric,G., Yahara,K., Mageiros,L., Pascoe,B., Maiden,M.C., Jolley,K.A. and Sheppard,S.K. (2014) A reference pan-genome approach to comparative bacterial genomics: identification of novel epidemiological markers in pathogenic *Campylobacter*. *PLoS One*, **9**, e29798.
13. Earle,S.G., Wu,C.-H., Charlesworth,J., Stoesser,N., Gordon,N.C., Walker,T.M., Spencer,C. C.A., Iqbal,Z., Clifton,D.A., Hopkins,K.L. *et al.* (2016) Identifying lineage effects when controlling for population structure improves power in bacterial association studies. *Nat. Microbiol.*, **1**, 16041.
14. Chewapreecha,C., Martinen,P., Croucher,N.J., Salter,S.J., Harris,S.R., Mather,A.E., Hanage,W.P., Goldblatt,D., Nosten,F.H., Turner,C. *et al.* (2014) Comprehensive identification of single nucleotide polymorphisms associated with beta-lactam resistance within pneumococcal mosaic genes. *PLoS Genet.*, **10**, e1004547.
15. Lees,J.A., Croucher,N.J., Goldblatt,D., Nosten,F., Parkhill,J., Turner,C., Turner,P. and Bentley,S.D. (2017) Genome-wide identification of lineage and locus specific variation associated with pneumococcal carriage duration. *eLife*, **6**, e26255.
16. Page,A.J., Cummins,C.A., Hunt,M., Wong,V.K., Reuter,S., Holden,M.T., Fookes,M., Falush,D., Keane,J.A. and Parkhill,J. (2015) Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics*, **31**, 3691–3693.
17. Zhao,Y., Jia,X., Yang,J., Ling,Y., Zhang,Z., Yu,J., Wu,J. and Xiao,J. (2014) PanGP: a tool for quickly analyzing bacterial pan-genome profile. *Bioinformatics*, **30**, 1297–1299.
18. Zhao,Y., Wu,J., Yang,J., Sun,S., Xiao,J. and Yu,J. (2012) PGAP: pan-genomes analysis pipeline. *Bioinformatics*, **28**, 416–418.
19. Laing,C., Buchanan,C., Taboada,E.N., Zhang,Y., Kropinski,A., Villegas,A., Thomas,J.E. and Gannon,V.P. (2010) Pan-genome sequence analysis using Panseq: an online tool for the rapid analysis of core and accessory genomic regions. *BMC Bioinformatics*, **11**, 1.
20. Koonin,E., Wolf,Y. and Puigbò,P. (2009) The phylogenetic forest and the quest for the elusive tree of life. *Cold Spring Harb. Symp. Quant. Biol.*, **74**, 205–213.
21. Huson,D.H. and Bryant,D. (2006) Application of phylogenetic networks in evolutionary studies. *Mol. Biol. Evol.*, **23**, 254–267.
22. Huson,D.H. (1998) SplitsTree: analyzing and visualizing evolutionary data. *Bioinformatics*, **14**, 68–73.
23. Fouts,D.E., Brinkac,L., Beck,E., Inman,J. and Sutton,G. (2012) PanOCT: automated clustering of orthologs using conserved gene neighborhood for pan-genomic analysis of bacterial strains and closely related species. *Nucleic Acids Res.*, **40**, e172.
24. Buchfink,B., Xie,C. and Huson,D.H. (2015) Fast and sensitive protein alignment using DIAMOND. *Nat. Methods*, **12**, 59–60.
25. Gibbons,T.R., Mount,S.M., Cooper,E.D. and Delwiche,C.F. (2015) Evaluation of BLAST-based edge-weighting metrics used for homology inference with the Markov Clustering algorithm. *BMC Bioinformatics*, **16**, 218.
26. Enright,A.J., Van Dongen,S. and Ouzounis,C.A. (2002) An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res.*, **30**, 1575–1584.
27. Dongen,S.V. (2000) *Graph Clustering by Flow Simulation*. University of Utrecht.

28. Katoh, K., Misawa, K., Kuma, K.-i. and Miyata, T. (2002) MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.*, **30**, 3059–3066.
29. Price, M.N., Dehal, P.S. and Arkin, A.P. (2010) FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS ONE*, **5**, e9490.
30. Stamatakis, A. (2014) RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*, **30**, 1312–1313.
31. Neher, R.A. and Bedford, T. (2015) nextflu: real-time tracking of seasonal influenza virus evolution in humans. *Bioinformatics*, **31**, 3546–3548.
32. Bertels, F., Silander, O.K., Pachkov, M., Rainey, P.B. and Nimwegen, E.v. (2014) Automated reconstruction of whole-genome phylogenies from short-sequence reads. *Mol. Biol. Evol.*, **31**, 1077–1088.
33. Felsenstein, J. (2004) *Inferring Phylogenies*. Sinauer.
34. Felsenstein, J. (1992) Phylogenies from restriction sites: a maximum-likelihood approach. *Evolution*, **46**, 159–173.
35. Zamani-Dahaj, S.A., Okasha, M., Kosakowski, J. and Higgs, P.G. (2016) Estimating the frequency of horizontal gene transfer using phylogenetic models of gene gain and loss. *Mol. Biol. Evol.*, **33**, 1843–1857.
36. Hudson, R. (2002) Ms a program for generating samples under neutral models. *Bioinformatics*, 337–338.
37. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. and Tanabe, M. (2016) KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res.*, **44**, D457–D462.
38. Huson, D.H. and Steel, M. (2004) Phylogenetic trees based on gene content. *Bioinformatics*, **20**, 2044–2049.
39. Baumdicker, F., Hess, W.R. and Pfaffelhuber, P. (2010) The diversity of a distributed genome in bacterial populations. *Ann. Appl. Probab.*, **20**, 1567–1606.
40. Rambaut, A. and Grass, N.C. (1997) Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Bioinformatics*, **13**, 235–238.
41. Hasegawa, M., Kishino, H. and Yano, T.-A. (1985) Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. *J. Mol. Evol.*, **22**, 160–174.
42. Trachana, K., Larsson, T.A., Powell, S., Chen, W.-H., Doerks, T., Muller, J. and Bork, P. (2011) Orthology prediction methods: a quality assessment using curated protein families. *Bioessays*, **33**, 769–780.
43. Li, L., Stoekert, C.J. and Roos, D.S. (2003) OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Res.*, **13**, 2178–2189.
44. Emms, D.M. and Kelly, S. (2015) OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.*, **16**, 157.
45. Baumdicker, F., Hess, W.R. and Pfaffelhuber, P. (2012) The infinitely many genes model for the distributed genome of bacteria. *Genome Biol. Evol.*, **4**, 443–456.
46. Gómez, J., García, L.J., Salazar, G.A., Villaveces, J., Gore, S., García, A., Martín, M.J., Launay, G., Alcántara, R., Ayllón, N.D.T. *et al.* (2013) BioJS: an open source JavaScript framework for biological data visualization. *Bioinformatics*, **29**, 1103–1104.
47. Biller, S.J., Berube, P.M., Lindell, D. and Chisholm, S.W. (2014) Prochlorococcus: the structure and function of collective diversity. *Nat. Rev. Microbiol.*, **13**, 13–27.
48. Biller, S.J., Berube, P.M., Berta-Thompson, J.W., Kelly, L., Roggensack, S.E., Awad, L., Roache-Johnson, K.H., Ding, H., Giovannoni, S.J., Rocap, G. *et al.* (2014) Genomes of diverse isolates of the marine cyanobacterium Prochlorococcus. *Scientific Data*, **1**, 140034.
49. Seemann, T. (2014) Prokka: Rapid prokaryotic genome annotation. *Bioinformatics*, **30**, 2068–2069.
50. Sun, Z., Blanchard, J.L., Bleasby, A., Chenna, R. and McGettigan, P. (2014) Strong genome-wide selection early in the evolution of prochlorococcus resulted in a reduced genome through the loss of a large number of small effect genes. *PLoS ONE*, **9**, e88837.
51. Berube, P.M., Biller, S.J., Kent, A.G., Berta-Thompson, J.W., Roggensack, S.E., Roache-Johnson, K.H., Ackerman, M., Moore, L.R., Meisel, J.D., Sher, D. *et al.* (2014) Physiology and evolution of nitrate acquisition in Prochlorococcus. *ISME J.*, **9**, 1195–1207.
52. Lindell, D., Sullivan, M.B., Johnson, Z.I., Tolonen, A.C., Rohwer, F. and Chisholm, S.W. (2004) Transfer of photosynthesis genes to and from Prochlorococcus viruses. *Proc. Natl. Acad. Sci. U.S.A.*, **101**, 11013–11018.
53. Croucher, N.J., Finkelstein, J.A., Pelton, S.I., Parkhill, J., Bentley, S.D., Lipsitch, M. and Hanage, W.P. (2015) Population genomic datasets describing the post-vaccine evolutionary epidemiology of *Streptococcus pneumoniae*. *Sci Data*, **2**, 150058.
54. Marttinen, P., Croucher, N.J., Gutmann, M.U., Corander, J. and Hanage, W.P. (2015) Recombination produces coherent bacterial species clusters in both core and accessory genomes. *Microb. Genomics*, **1**, e000038.
55. Del Grosso, M., Iannelli, F., Messina, C., Santagati, M., Petrosillo, N., Stefani, S., Pozzi, G. and Pantosti, A. (2002) Macrolide efflux genes *mef(A)* and *mef(E)* are carried by different genetic elements in *Streptococcus pneumoniae*. *J. Clin. Microbiol.*, **40**, 774–778.
56. Leclercq, R. and Courvalin, P. (1991) Bacterial resistance to macrolide, lincosamide, and streptogramin antibiotics by target modification. *Antimicrob. Agents Chemother.*, **35**, 1267–1272.
57. Bostock, M. (2016) D3: Data-Driven Documents.
58. Woodhull, G. and Zhu, N. *et al.* (2016) dc.js - Dimensional Charting Javascript Library.
59. Gardy, J., Loman, N.J. and Rambaut, A. (2015) Real-time digital pathogen surveillance—the time is now. *Genome Biol.*, **16**, 155.
60. Holden, M. T.G., Hsu, L.-Y., Kurt, K., Weinert, L.A., Mather, A.E., Harris, S.R., Strommenger, B., Layer, F., Witte, W., de Lencastre, H. *et al.* (2013) A genomic portrait of the emergence, evolution, and global spread of a methicillin-resistant *Staphylococcus aureus* pandemic. *Genome Res.*, **23**, 653–664.