

Genome analysis

Hierarchical Sets: Analyzing Pangenome Structure through Scalable Set Visualizations

Thomas Lin Pedersen^{1,2,*}

¹Center for Biological Sequence Analysis, Department of Systems Biology, The Technical University of Denmark, Building 208, DK-2800 Lyngby, Denmark,

²Assays, Culture and Enzymes Division, Chr. Hansen A/S, DK-2970 Hørsholm, Denmark

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: The increase in available microbial genome sequences has resulted in an increase in the size of the pangenomes being analyzed. Current pangenome visualizations are not intended for the pangenome sizes possible today and new approaches are necessary in order to convert the increase in available information to increase in knowledge. As the pangenome data structure is essentially a collection of sets we explore the potential for scalable set visualization as a tool for pangenome analysis.

Results: We present a new hierarchical clustering algorithm based on set arithmetics that optimizes the intersection sizes along the branches. The intersection and union sizes along the hierarchy are visualized using a composite dendrogram and icicle plot, which, in pangenome context, shows the evolution of pangenome and core size along the evolutionary hierarchy. Outlying elements, i.e. elements whose presence pattern do not correspond with the hierarchy, can be visualized using hierarchical edge bundles. When applied to pangenome data this plot shows putative horizontal gene transfers between the genomes and can highlight relationships between genomes that is not represented by the hierarchy. We illustrate the utility of hierarchical sets by applying it to a pangenome based on 46 *Streptococcus* genomes and find it provides a powerful addition to pangenome analysis.

Availability: The described clustering algorithm and visualizations are implemented in the hierarchicalSets R package available from CRAN (<https://cran.r-project.org/web/packages/hierarchicalSets>)

Contact: Thomas Lin Pedersen (thomasp85@gmail.com)

Supplementary information Supplementary data are available at Bioinformatics online.

1 Introduction

Pangenome analysis is concerned with the investigation of multiple bacterial genomes whose genes have been grouped according to similarity. A pangenome is thus defined as a set of gene groups containing members from one or more of genomes. Figure 1 shows the general structure of a pangenome as visualized by a presence/absence matrix. Gene groups are often classified by their ubiquity in the genomes making up the pangenome. *Core* gene groups are present in all genomes, *accessory* gene groups are present in more than one, but not all genomes and *singleton* gene groups are only present in one genome. This classification of gene groups gives a broad overview of the heterogeneity of the pangenome through the number of core gene groups and total gene groups, but is also used to

pinpoint the nature of the genes within each group. Core genes are likely genes that define the unique traits of the genomes under investigation, while accessory genes are disposable genes that define more specialized behavior. Singleton genes can be strain specific genes, pseudogenes, or annotation errors. As is evident from figure 1, there are clear overlaps between the nomenclature associated with pangenome data and that of set algebra, where genomes can be considered sets and gene groups elements in these sets. Furthermore, intersection and core size as well as pangenome and union size are equivalent.

The first published pangenome covered eight strains of *Streptococcus agalactiae* (Tettelin *et al.*, 2005), reflecting the number of available genome sequences for that species at the time. The number of genomes included in pangenome analyses has since increased along with the increased availability of sequenced bacterial genomes and now contains 100s or

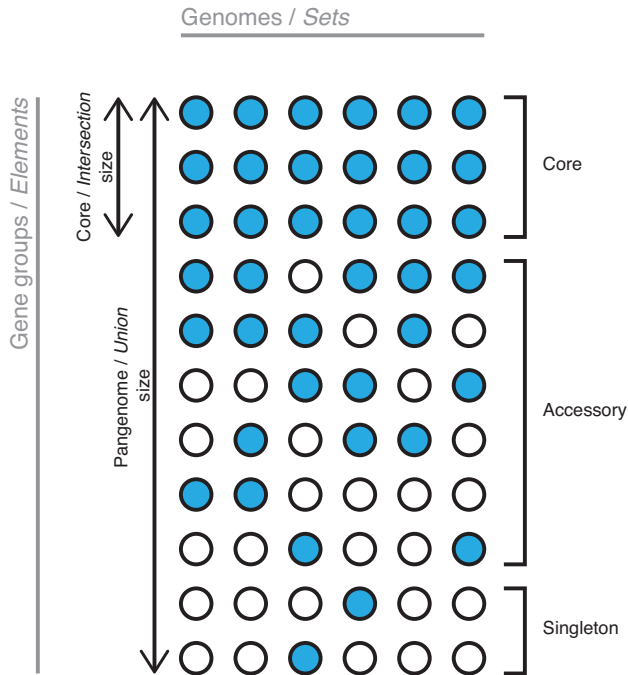


Fig. 1. Overview of the nature of pangenome data and the nomenclature associated with it. Equivalent set algebra terms are shown in *italic*. Columns define genomes and rows gene groups. A filled circle indicates the presence of a member of the respective gene group in the genome while an empty circle indicates absence.

1,000s of genomes (Leekitcharoenphon *et al.*, 2016; Land *et al.*, 2015; Jun *et al.*, 2014; Méric *et al.*, 2013; Snipen and Ussery, 2012; Kaas *et al.*, 2011) resulting in >10,000 gene groups. A main concern when evaluating the result of pangenome analyses is how the pangenome and core size change as genomes are added to the pangenome. Sudden drops in core size or jumps in pangenome size indicate the addition of a genome deviating strongly from the genomes already present in the pangenome. The standard approach to show this evolution in pangenome and core size is through a simple line-plot as shown in figure 2 (Smokvina *et al.*, 2012; De Maayer *et al.*, 2014; Lukjancenko *et al.*, 2010). This approach has considerable drawbacks as the shape of the line is determined by the order in which genomes are added. While it is possible to define a progression of genomes that ensures that similar genomes follow each other, changes between genomes will still be obscured by the level of heterogeneity between the genomes that comes before it. The extreme case is a pangenome without any core gene groups. At some point along the line-plot the core line will drop to zero and any difference between genomes that follows this point will be invisible. The set nature of pangenome data could offer a better way of visualizing the change in core and pangenome size without imposing a specific order to the genomes. Set algebra has been used sparingly in pangenome visualizations. GenoSets (Cain *et al.*, 2012) and PanViz (Pedersen *et al.*, 2016) both apply set arithmetic to create visual queries for gene group subsets. Apart from query construction though, set algebra is largely unexplored when it comes to visualizing the relational structure between genomes. While visualizing relations between large numbers of sets is difficult due to the combinatorial explosion of possible set combinations, different visualization techniques have been developed to show intersection sizes between sets in a scalable manner such as UpSet (Lex *et al.*, 2014) and Radial Sets (Alsallakh *et al.*, 2013). These techniques do not scale to the number of sets that is exposed in contemporary pangenomes though and are thus a poor fit for investigating all but the simplest pangenomes.

Here, we present a new approach to set analysis and visualization called Hierarchical Sets, that works particularly well on large structured collections of sets such as pangenomes. Hierarchical Sets limits the comparisons between sets to branch points of a hierarchical clustering; thus, achieving good scalability at the expense of not showing direct comparisons between very dissimilar sets. While the focus in this paper is on the use of Hierarchical Sets in pangenome visualization, the technique can be applied equally well to other problems involving large numbers of sets.

2 Data

The data set used for the examples is a pangenome based on 46 different *Streptococcus* genomes. Each genome is identified by its NCBI assembly accession number. The pangenome has been created using FindMyFriends (Pedersen, 2015) using default parameters and consists of 37,889 gene groups distributed among 49 core groups, 20,337 accessory groups and 17,503 singleton groups. A standard line-plot representation of the pangenome can be seen in figure 2. A presence/absence matrix of the pangenome is available in the supplementary material.

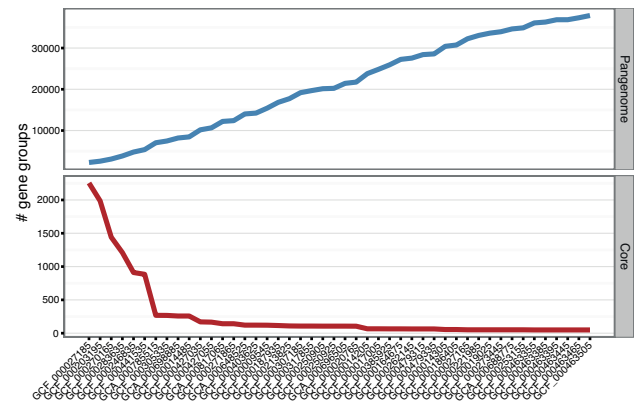


Fig. 2. Line plot showing the evolution of pangenome and core genome size as genomes are added to the pangenome. The ordering along the x-axis follows the ordering given by the Hierarchical Sets clustering described in the Hierarchical clustering approach section and matches that in figure 4.

3 Algorithm

Existing approaches for hierarchical clustering of sets or pangenomes usually follows a conversion of the data into a distance matrix followed by an agglomerative clustering. For pangenomes several distance measures have been used, e.g. binary (Richards *et al.*, 2014), Jaccard (Kuenne *et al.*, 2013) or Manhattan distance (Jacobsen *et al.*, 2011) as well as several clustering algorithms such as average (Karlsson *et al.*, 2011) or single linkage (Tettelin *et al.*, 2005). These approaches have several drawbacks when it comes to interpreting the results in a set algebraic context. The reliance of a conversion to a distance matrix makes the clustering extremely sensitive to the choice of clustering algorithm as the clustering is no longer based on the original data. Furthermore it implies that a distance exists for combinations of sets which might not make sense if two sets are fully independent (no intersecting elements). The result of the former is that the result of standard hierarchical clusterings can be hard to translate back to features of the set data, while the latter results in all sets being merged into a final cluster even though there might not be any similarity between

all sets in the analysis. To address these shortcomings we introduce a new agglomerative hierarchical clustering approach for sets that works directly with the set data itself, by means of a set family homogeneity measure defined below. The clustering happens through the following steps:

1. Let each set in the analysis define their own set family of size 1.
2. For each pair of set families calculate the homogeneity, λ , of the combined set family.
3. Choose the pair that exhibit the highest λ (on ties choose the pair with the smallest union) and let the pair define a new set family.
4. Repeat 2-3 until all available set family pairs have $\lambda = 0$.

Note that this approach specifically terminates the clustering before all sets have been combined to a single cluster if the remaining clusters have no pairwise homogeneity.

3.1 Set family homogeneity and heterogeneity measure

Similarity between two sets are often measured using Jaccard similarity defined as the size of their intersection divided by the size of their union. The similarity between two sets can also be thought of as the homogeneity of a set family consisting of the two sets. The Jaccard similarity can then be generalized to a measure of set family homogeneity for set families of any size by dividing the total intersection size with the total union size. Formally, for a set family A , the set family homogeneity λ is defined by:

$$\lambda(A) = \frac{|\cap(A)|}{|\cup(A)|}$$

In the case of pangenomes, the data is often incomplete as there is a chance to miss genes during sequencing and annotation. Therefore, core size can be underestimated and it is a custom to loosen the requirement for gene groups to be considered core by requiring the fraction of genomes represented in a core group to be above a fixed threshold (such as 0.95). The set family homogeneity definition can be modified to accommodate this practice by introducing a parameter $t \in [0, 1]$ that defines the ratio threshold for an element to be considered part of the intersection ($t = 1$ will result in the standard intersection definition). The set family homogeneity subject to t can thus be defined as:

$$\lambda(A)_t = \frac{\sum_{i=1}^n \frac{\sum_{j=1}^m A_{i,j}}{m}}{|\cup(A)|}$$

where A is the set family, n is the universe size, m is the number of sets in the family, and t a value between 0 and 1. $A_{i,j}$ is 1 if element i is present in set j and 0 otherwise. Similar to the Jaccard similarity the set family homogeneity is bound between 0 and 1 ($\lambda \in [0, 1]$). Conversely, the set family heterogeneity is defined as:

$$\lambda'(A)_t = \lambda(A)_t^{-1} - 1$$

And it follows that $\lambda' \in [0, \infty]$. This definition makes λ' undefined for set families with $\lambda = 0$, which is sensible as the heterogeneity of a collection with no homogeneity must be undefined.

4 Results

4.1 Visualizing set family heterogeneity

An obvious way to present the result of the clustering is through the use of a dendrogram. By encoding the height of the branch points to λ' , the dendrogram will illustrate how the heterogeneity increases as set families are combined. This dendrogram encoding is particularly good at

identifying clusters of highly homogeneous sets as well as independent clusters (figure 3).

Compared to a standard hierarchical clustering (Jaccard distance + complete linkage) there are a few differences in the clustering itself, but more so in the interpretation of the dendrogram. While λ' can clearly be interpreted as the ratio of intersection to union for the sets contained in each branch point, complete linkage only shows the Jaccard distance between the most distant sets in the branch leaving the relationship of the remainder of the sets obscured. Furthermore, as the Jaccard distance is bound below 1 the branch points stack closer and closer to 1 giving the appearance of increasingly smaller changes in distance as clusters are merged. Even for distance measures without an upper bound, such as Manhattan, the conversion to a distance matrix implicitly creates an upper bound (maximum distance between two sets in the data) and branch points tends to be stacked up against this (alternative clusterings can be seen in figure S1 in the supplementary material).

4.2 Visualizing intersection and union sizes

Often in set analysis there is an interest in the intersection sizes of the different combinations of sets. For a number of sets, n , the number of possible set families are $2^n - 1$, resulting in 7e13 possible set families for the 46 sets used as example in this paper. This combinatorial explosion has made it difficult to visualize intersection sizes for large numbers of sets. The Hierarchical Sets clustering offers a way to decrease the number of set families by only considering set families at branch points. The intersection sizes of each branch point can be visualized while preserving the hierarchical layout by using an inverted icicle plot with bar height encoded to intersection size (figure 4, bottom). The plot can be envisioned as a stack of blocks where the height of the stack denotes the total value and the height of the block denotes the contribution of that single block.

Based on this plot a lot of information can be decoded. The intersection size of the different set families defined by the branch points are shown as the absolute height of the stacks while the drop in intersection size is shown as the height of each block. To improve visual separation of the blocks, their fill color is encoded to the number of the sets represented by the family. This type of plot can show relational structure between the different sets: Dark, narrow bars starting close to the x-axis (GCF_000009545 in figure 4) represent sets having little overlap with the rest of the sets, while light and wide bars represent larger collections of sets showing large overlaps. The near absence of single-width bars (GCF_000463395 and GCF_000463445 in figure 4) indicate near-similar sets and the height of each single-width bar shows the total size of each set.

In the same way as intersection at each branch point can be shown, so can the union. In contrast to the intersection, the union decreases as you approach the leafs of the hierarchy, making a dendrogram a better choice for this (figure 4, top). While the union dendrogram would extend naturally from the top of the bars in the icicle plot, as the union and intersection of a single set are equal, the range of unions often vary substantially from that of intersections. Thus, it is a better choice to plot them in separate plots, but stacked so that they share the x-axis.

4.3 Visualizing deviations from hierarchy

Imposing a hierarchy on a dataset is likely to distort the data as complete adherence to a hierarchical structure is rare. In the case of a hierarchical set analysis, deviations from the hierarchy can be thought of as elements shared by two sets, but not part of their common set family (figure 5). More formally outlying elements \hat{x} can be defined as

$$\hat{x} \in \cap(A, B) \wedge \hat{x} \notin \cap(C_{A,B})$$

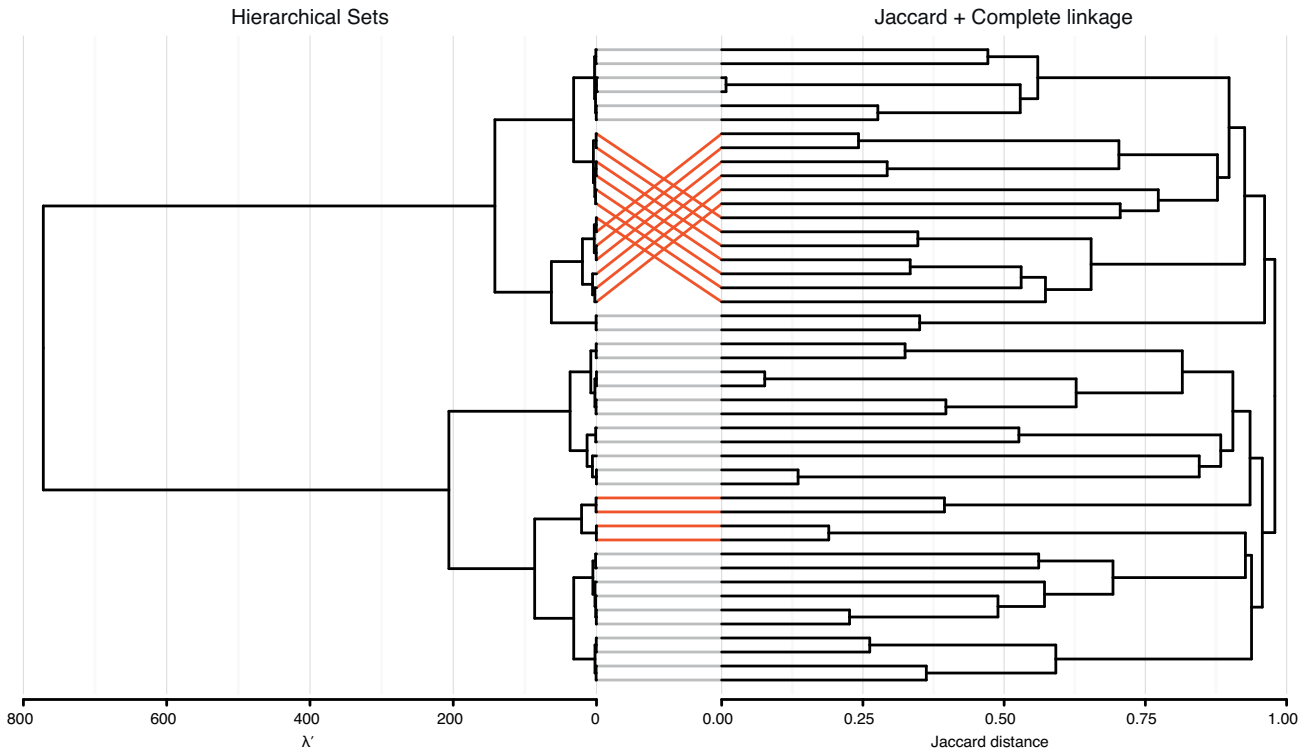


Fig. 3. Comparison of hierarchical set clustering and complete linkage clustering based on Jaccard distance as performed on a pangenome based on 46 *emphStreptococcus* strains. A grey link between leaf nodes indicate agreement between the methods while a red link indicates disagreement (i.e., leafs are part of different sub-clusters).

Where A and B are sets and $C_{A,B}$ the smallest set family containing A and B derived from the hierarchical clustering.

Visualizations of outlying elements can be either set- or element centric, depending on whether the focus is on how pairs of sets deviate from the hierarchy or on the individual elements that make up the deviation. Showing statistics on pairs of sets can be done effectively using a heatmap. By overlaying both hierarchy information and pair information in the same way as done by dendrogramix (Blanch *et al.*, 2015), it is possible to get a matrix plot that both shows the intersection at each branch point, as well as the intersection and union size of each set pair. The contrasts between the branch point intersections and the set pair intersections are thus indicative of the amount of deviation from the hierarchy that each pair of sets exhibit (see figure S2 in supplementary material).

An alternative way to show connections between leafs in a hierarchical clustering is by using hierarchical edge bundling (Holten, 2006). To avoid overplotting, edges can be filtered by weight (number of outlying elements), in order to only show the strongest deviations from the structure (figure 6).

The elements themselves can be investigated as well, based on the outlying elements approach outlined above. Counting the number of times each element appears as outlying will give an indication of each elements propensity to not conform with the hierarchy. As the number of times an element can appear as an outlier is governed by the number of times it appears in a set, these two values can be shown in a scatter plot (see figure S3 in supplementary material) to quickly identify elements exhibiting unexpectedly high deviation.

5 Discussion

We have presented a new approach to hierarchical clustering of set data, a range of scalable visualizations that builds on top of the clustering, and

an outlier definition for elements based on the clustering. Hierarchical Set analysis optimizes intersection size at each branch point, making it easier to reason about the clustering and, as a consequence, the visualizations. Hierarchical Set analysis is particularly well-suited for pangenome analysis as pangenome data often consists of a large number of sets with a clear hierarchical structure due to the evolutionary nature of genomes.

5.1 Pangenome evolution

In the context of pangenomes the intersection is equivalent to the core, while the union equates the pangenome. As such there is strong similarity between figure 2 and figure 4 as they both try to convey the same type of information (i.e., the change in pangenome and core size as additional genomes are added). The main difference is that figure 4 shows the core and pangenome sizes along a hierarchy instead of along a linear progression as in figure 2. The benefit of the hierarchical sets approach is that evolutionary features are not obscured. The line-plot hardly shows any change in core size in the last half of the plot despite the fact that this group of genomes are just as diverse as the first half. Furthermore, figure 4 also conveys the hierarchical structure of the pangenome, information that is very relevant when evaluating core and pangenome sizes of different subsets of the pangenome. Based on figure 3 and 4, it is obvious that the pangenome consists of at least four clusters of genomes sharing very little genomic material with each other (figure 4A1–4). Furthermore, three pairs of genomes shows little similarity to any of the other genomes (figure 4B) and two pairs of genomes are almost completely homogeneous (figure 4C). In addition, the size of each genome is clearly visible as well as pangenome size for relevant subsets of genomes. Compare this with what can be decoded from the current approach to showing evolution in core and pangenome size (figure 2), namely that pangenome continue to increase

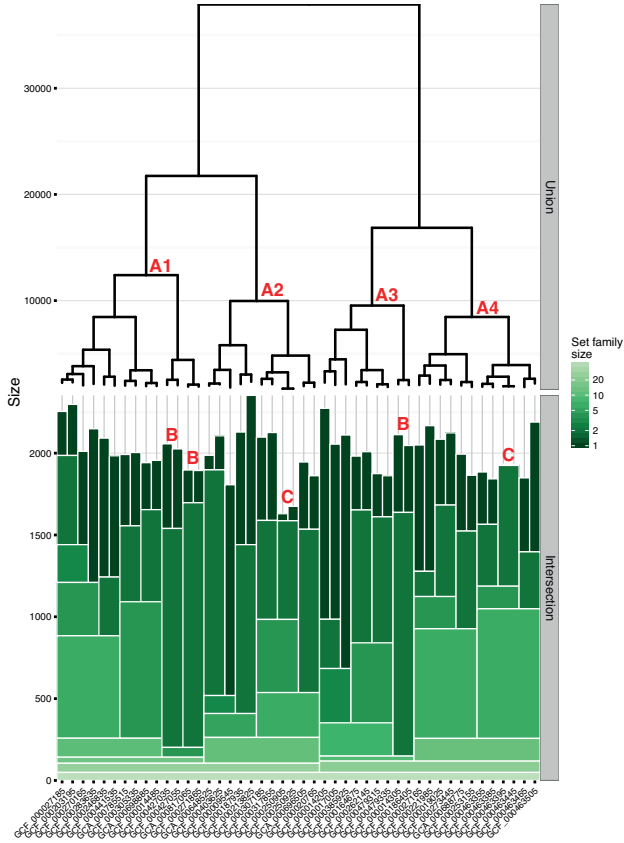


Fig. 4. Intersection and union sizes at the branch points in a hierarchical set clustering, visualized as an icicle plot for the intersections and a dendrogram for the unions. The intersection size of each set family is encoded to the height of the bar and the size of the set family are encoded to the color of the bar. The area of each rectangle is thus proportional to the number of sets it represents and the increase in intersection size relative to the next branch point.

linearly as genomes are added and that the core quickly drops to a very low size.

5.2 Deviations from the hierarchy

There is a clear similarity between the Hierarchical Sets based heatmap visualization (figure S1) and the BLAST matrices often used to show similarities between genomes in a pangenome, e.g., figure 3 in (Lukjancenko *et al.*, 2012). The Hierarchical Sets heatmap provides additional information though, allowing for both an assessment of the pairwise similarities as well as deviation between the pairwise similarity and the similarity defined by their common ancestor. The deviation, defined as outlying elements in the context of Hierarchical Sets, has a clear analogy in gene deletion and horizontal gene transfer events. Such events results in distributions of gene groups not governed by the evolutionary hierarchy of the genomes itself but more related to shared environment. These events can be of just as much interest as the hierarchical structure itself. Detecting structure in where these events occur, in relation to the evolutionary hierarchy, can help researchers detect strong cross-talk between evolutionary unrelated organisms. In contrast to the heatmap approach used in figure S1, hierarchical edge bundles puts focus on larger structures in the deviation, while obscuring the single pairwise values due to overplotting e.g., stronger cross-talk between organisms in group A3-4 (figure 6A) compared to group A1-2 (figure 6B). This also means that odd, but weaker structures, such as the edge crossing the root node

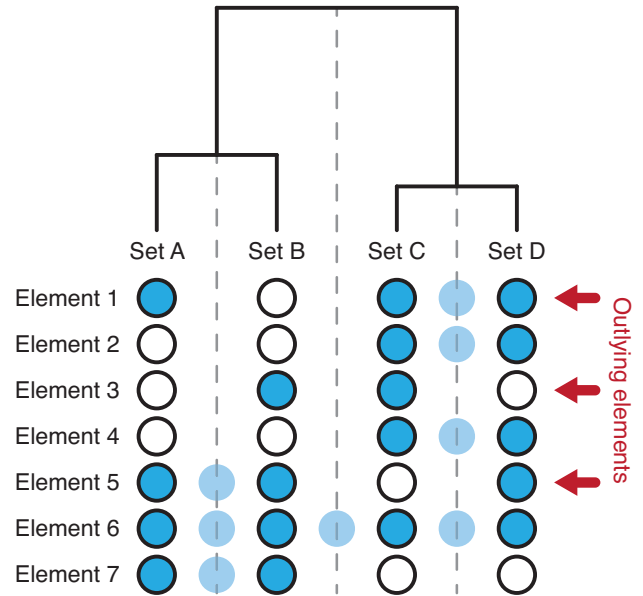


Fig. 5. Definition of outlying elements: Set A-D are sets defined by the presence of element 1-7. Blue filled circles indicate presence while empty circles indicate absence. The shaded circles on the dashed lines shows the set family intersection of the families defined by the clustering. The red arrows shows outlying elements, i.e., elements that are shared by two sets but not shared by their common set family.

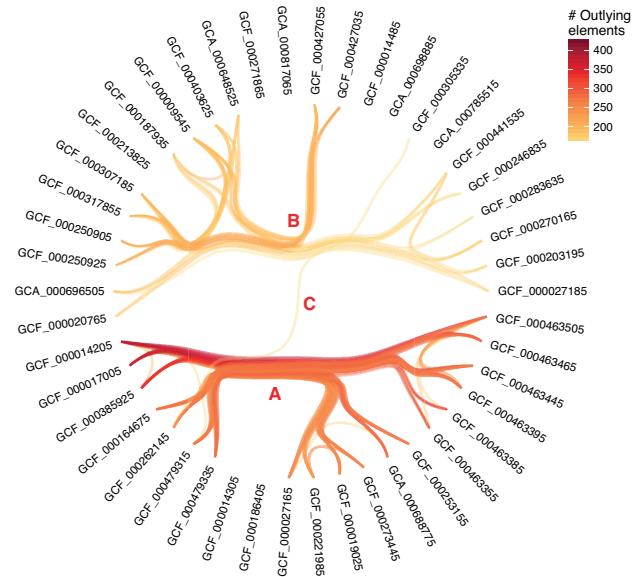


Fig. 6. Hierarchical edge bundling showing the 15% strongest deviations from the hierarchy defined by a hierarchical set analysis (measured in number of outlying elements). Color is mapped to number of outlying elements.

of the hierarchy (figure 6C), can be recognized more easily. Another benefit, in the context of pangenomes and evolutionary analysis, is that origin branch for the deviation is clearly visible, giving a clear hint on the evolutionary distance of the mutation. Furthermore, the likening to a dendrogram structure will feel familiar to many biologists.

5.3 Deviating gene groups

Looking into the diverging elements themselves and the number of times elements appear as outliers can guide researchers looking into mobile

elements. The elements appearing as outliers constitute rows of the presence/absence matrix not conforming to the hierarchical structure. Extracting these rows and performing a second Hierarchical Sets analysis based on them will reveal the second most dominant structure in the dataset. Conceptually, this is equivalent to a principal component analysis (PCA) where components gradually diminish in explanatory power as they focus on structures not captured by the components before them. In evolutionary context the main hierarchy revealed by Hierarchical Sets analysis is likely related (but not necessary identical) to the evolutionary tree of the genomes under investigation, while a secondary hierarchy based on outlying elements would reveal structures pertaining to increased strain interactions such as ecological niches. It is possible to continue creating sub-hierarchies based on outlying elements, but as with PCA the likelihood of beginning to model noise will increase with each step.

5.4 Streptococcus through hierarchical sets

Based on the figures presented here several insights can be gained into the pangenome of the *Streptococcus* genus. The clustering generated based on hierarchical sets provide a well defined clustering of the viridans species in line with the groups defined by (Facklam, 2002). Hemolytic (pyogenic) species are, with the exception of *St. agalactiae* all present in cluster A2. Cluster A1 contains the viridans groups bovis, salivarius, and mutans, as well as the pyogenic *St. agalactiae*. In addition to the pyogenic species, cluster A2 contains the species *St. uberis* and *St. parauberis* that are not part of any grouping scheme (Facklam, 2002). Cluster A3 contains the viridans group sanguini as well as *St. suis*. *St. oligofermentans* is placed in the middle of the sanguinis group based on the hierarchical set clustering. Upon its discovery it was reported to belong to the mitis group (Tong *et al.*, 2003), but the placement based on hierarchical sets is supported by finding in Maruyama *et al.* (2016), indicating an initial misclassification. Also, *St. sp. I-P16* and *St. sp. I-G2* appears to be part of the sanguinis group as well. Cluster A4 contains the viridans groups anginosus and mitis. Three species (*St. agalactiae*, *St. mutans*, and *St. suis*) are shown to be very distant to all other species in the pangenome (figure 4B). All of the above findings are in line with recent phylogenetic studies based on other clustering approaches (Thompson *et al.*, 2013; Richards *et al.*, 2014; Maruyama *et al.*, 2016). There are small differences in how the major clusters relate to each others compared to other sequence based clusterings (see e.g. Thompson *et al.* (2013)), but agglomerative clustering often loose precision at the top of the hierarchy due to the inadequacy of a single distance measure to describe increasingly heterogeneous clusters. Whether this explains the differences in clustering is difficult to assess. It could also be due to pangenome and single sequence based phylogenies describing different relations. The pangenome phylogeny presented by Richards *et al.* (2014) is in general more in agreement with the clustering presented here so it is likely that the latter explanation holds true.

The addition of the outlying element analysis to pangenome investigations allows a second layer to be added to the phylogeny. It is obvious that cluster A3 and A4 exhibits a very strong connection with each other, not captured by their core (figure 6A). The sanguini group making up most of cluster A3 have often been clustered within the mitis group (Doern and Burnham, 2010; Facklam, 2002; Thompson *et al.*, 2013) based on 16S rRNA comparison. The hierarchical sets analysis shows that these two groups have distinct cores but a large overlap in their accessory genome, indicating a high degree of cross-talk after they delineated. The strong connection between the sanguini and anginosus group reflects their proximity in the phylogeny reported by Maruyama *et al.* (2016). *St. suis* again demonstrates its distance to the other members of cluster A3 and A4 by not exhibiting any strong connection to other members of the clusters. In general, the species represented in cluster A3 and A4 are co-inhabitants of the respiratory tract, which could explain

the strong deviations from the pure hierarchical phylogeny, as they have ample opportunity for horizontal gene transfer. The single deviation visible across the top node in the hierarchy (figure 6C) denotes a week, but present tie between *St. salivarius* and *St. parasanguinis*. Both of these species colonizes the oral cavity, again indicating the possible existence of horizontal gene transfer between the two species. The link is only between one of the two *St. salivarius* strains represented in the pangenome, indicating that genomic material has been shared after the salivarius species became defined. Within cluster A1 and A2, the most dominant link is between *St. agalactiae* and the pyogenic species in cluster A2 (figure 6B), reflecting the link between the pyogenic species. This could indicate that the placement of *St. agalactiae* by hierarchical sets is wrong, but it is clear that the other pyogenic species share a sizable core that is not present in *St. agalactiae*. Other phylogenies also place *St. agalactiae* apart from the remaining pyogenic species (Thompson *et al.*, 2013; Maruyama *et al.*, 2016) and the pathogenicity (sepsis in newborns) of the species is also unique among the pyogenic streptococci.

6 Implementation and Availability

The described clustering algorithm as well as the different visualizations are implemented in the hierarchicalSets R package and available for free on all major platforms through CRAN (R Core Team, 2016). hierarchicalSets takes as input either a presence-absence matrix with sets as columns and elements as rows, or a list of sets defined by their elements. For use in pangenome analysis, hierarchicalSets can work directly with the data structures defined in the FindMyFriends package (Pedersen, 2015). hierarchicalSets uses common, memory efficient R data-structures and the clustering algorithm is written in C++ for speed, ensuring that the package scales well to thousands of sets with millions of elements.

7 Conclusion

Pangenome analyses continue to increase in scope, and visualization approaches that gracefully handle this increased complexity are paramount to extract knowledge from the results. Recent advances in pangenome analysis algorithms have facilitated the creation of pangenomes spanning thousands of genomes, covering the full bacterial domain, and current visualization techniques do not adequately support such large and heterogeneous pangenomes. Identifying the overlap between common set arithmetics and pangenome summaries, we have explored different approaches to scalable set visualization that can address the challenges posed by large pangenome datasets. We present a new range of set visualization approaches well-suited to large collections of structured sets, such as genomes in a pangenome. All presented visualizations are centered around a new hierarchical clustering technique, called Hierarchical Sets, that optimizes the intersection size along the branch points. Based on this clustering it is possible to create scalable visualizations of intersection and union sizes (core and pangenome size), as well as visualizing elements (gene groups) that deviate from the overall structure of the data. We show the utility of hierarchical sets in pangenome analysis by applying it to a genus-level pangenome based on 46 *Streptococcus* genomes, where the analysis both reveal support for the current viridans grouping based on shared core genome, as well as indicating several interesting cases of horizontal gene transfer between species, related to shared ecologic niche. The visualizations presented here do not rely on interactions in order to communicate their message, making them easy to incorporate into composite visualization frameworks or directly augment with interactivity. While Hierarchical Sets has been developed for the purpose of visualizing pangenome data, the approach is agnostic to the underlying data type,

and it could equally well be applied to other large-scale set visualization problems, especially set data with a clear hierarchical interpretation.

8 Acknowledgments

The authors thanks Kasper Dinkla and Hendrik Strobelt, Harvard, for fruitful discussion, suggestions, and help with preparing the manuscript, as well as Jan Egil Afset, Norwegian University of Science and Technology, for input related to the classification of *Streptococcus* and Maria Månsson, Chr. Hansen A/S, for help with the manuscript.

Funding: This work was supported by The Danish Agency for Science, Technology and Innovation.

Conflict of Interest: None declared.

References

- Alsallakh, B., Aigner, W., Miksch, S., and Hauser, H. (2013). Radial Sets: Interactive Visual Analysis of Large Overlapping Sets. *IEEE transactions on visualization and computer graphics*, **19**(12), 2496–2505.
- Blanch, R., Dautriche, R., and Bisson, G. (2015). Dendrogramix: A hybrid tree-matrix visualization technique to support interactive exploration of dendrograms. *IEEE transactions on visualization and computer graphics*, pages 31–38.
- Cain, A. A., Kosara, R., and Gibas, C. J. (2012). GenoSets: Visual Analytic Methods for Comparative Genomics. *PLoS ONE*, **7**(10), e46401.
- De Maayer, P., Chan, W. Y., Rubagotti, E., Venter, S. N., Toth, I. K., Birch, P. R. J., and Coutinho, T. A. (2014). Analysis of the *Pantoea ananatis* pan-genome reveals factors underlying its ability to colonize and interact with plant, insect and vertebrate hosts. *BMC Genomics*, **15**(1), 404.
- Doern, C. D. and Burnham, C.-A. D. (2010). It's not easy being green: the viridans group streptococci, with a focus on pediatric clinical manifestations. *Journal of Clinical Microbiology*, **48**(11), 3829–3835.
- Facklam, R. (2002). What Happened to the Streptococci: Overview of Taxonomic and Nomenclature Changes. *Clinical Microbiology Reviews*, **15**(4), 613–630.
- Holten, D. (2006). Hierarchical edge bundles: visualization of adjacency relations in hierarchical data. *IEEE transactions on visualization and computer graphics*, **12**(5), 741–748.
- Jacobsen, A., Hendriksen, R. S., Aarestrup, F. M., Ussery, D. W., and Friis, C. (2011). The *Salmonella enterica* pan-genome. *Microbial Ecology*, **62**(3), 487–504.
- Jun, S.-R., Wassenaar, T. M., Nookaew, I., Hauser, L., Wanchai, V., Land, M., Timm, C. M., Lu, T.-Y. S., Schadt, C. W., Doktycz, M. J., Pelletier, D. A., and Ussery, D. W. (2014). Diversity of *Pseudomonas* Genomes, Including *Populus*-Associated Isolates, as Revealed by Comparative Genome Analysis. *Applied and Environmental Microbiology*, **82**(1), 375–383.
- Kaas, R. S., Friis, C., Ussery, D. W., and Aarestrup, F. M. (2011). Estimating variation within the genes and inferring the phylogeny of 186 sequenced diverse *Escherichia coli* genomes. *BMC Genomics*, **13**, 577–577.
- Karlsson, F. H., Ussery, D. W., Nielsen, J., and Nookaew, I. (2011). A closer look at bacteroides: phylogenetic relationship and genomic implications of a life in the human gut. *Microbial Ecology*, **61**(3), 473–485.
- Kuenne, C., Billion, A., Mraheil, M. A., Strittmatter, A., Daniel, R., Goesmann, A., Barbuddhe, S., Hain, T., and Chakraborty, T. (2013). Reassessment of the *Listeria monocytogenes* pan-genome reveals dynamic integration hotspots and mobile genetic elements as major components of the accessory genome. *BMC Genomics*, **14**(1), 47.
- Land, M., Hauser, L., Jun, S.-R., Nookaew, I., Leuze, M. R., Ahn, T.-H., Karpinet, T., Lund, O., Kora, G., Wassenaar, T., Poudel, S., and Ussery, D. W. (2015). Insights from 20 years of bacterial genome sequencing. *Functional & Integrative Genomics*, **15**(2), 141–161.
- Leekitcharoenphon, P., Hendriksen, R. S., Le Hello, S., Weill, F.-X., Baggesen, D. L., Jun, S.-R., Ussery, D. W., Lund, O., Crook, D. W., Wilson, D. J., and Aarestrup, F. M. (2016). Global Genomic Epidemiology of *Salmonella enterica* Serovar Typhimurium DT104. *Applied and Environmental Microbiology*, **82**(8), 2516–2526.
- Lex, A., Gehlenborg, N., Strobelt, H., Vuillemot, R., and Pfister, H. (2014). UpSet: Visualization of Intersecting Sets. *IEEE transactions on visualization and computer graphics*, **20**(12), 1983–1992.
- Lukjancenko, O., Wassenaar, T. M., and Ussery, D. W. (2010). Comparison of 61 sequenced *Escherichia coli* genomes. *Microbial Ecology*, **60**(4), 708–720.
- Lukjancenko, O., Ussery, D. W., and Wassenaar, T. M. (2012). Comparative genomics of bifidobacterium, lactobacillus and related probiotic genera. *Microbial Ecology*, **63**(3), 651–673.
- Maruyama, F., Watanabe, T., and Nakagawa, I. (2016). *Streptococcus pyogenes* Genomics. In *Streptococcus pyogenes: Basic Biology to Clinical Manifestations*. University of Oklahoma Health Sciences Center, Oklahoma City (OK).
- Méric, G., Yahara, K., Mageiros, L., Pascoe, B., Maiden, M. C. J., Jolley, K. A., and Sheppard, S. K. (2013). A reference pan-genome approach to comparative bacterial genomics: identification of novel epidemiological markers in pathogenic *Campylobacter*. *PLoS ONE*, **9**(3), e92798–e92798.
- Pedersen, T. L. (2015). *FindMyFriends - Fast alignment-free pangenome creation and exploration*, 1.0.2 edition.
- R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 3.2.4 edition.
- Richards, V. P., Palmer, S. R., Bitar, P. D. P., Qin, X., Weinstock, G. M., Highlander, S. K., Town, C. D., Burne, R. A., and Stanhope, M. J. (2014). Phylogenomics and the dynamic genome evolution of the genus *Streptococcus*. *Genome Biology and Evolution*, **6**(4), 741–753.
- Smokvina, T., Wels, M., Polka, J., Chervaux, C., Brisse, S., Boekhorst, J., van Hylckama Vlieg, J. E. T., and Siezen, R. J. (2012). *Lactobacillus paracasei* comparative genomics: towards species pan-genome definition and exploitation of diversity. *PLoS ONE*, **8**(7), e68731–e68731.
- Snipen, L. G. and Ussery, D. W. (2012). A domain sequence approach to pangenomics: applications to *Escherichia coli*. *F1000Research*, **1**, 19.
- Tettelin, H., Maignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., Angiuoli, S. V., Crabtree, J., Jones, A. L., Durkin, A. S., Deboy, R. T., Davidsen, T. M., Mora, M., Scarselli, M., Margarit y Ros, I., Peterson, J. D., Hauser, C. R., Sundaram, J. P., Nelson, W. C., Madupu, R., Brinkac, L., Dodson, R. J., Rosovitz, M. J., Sullivan, S. A., Daugherty, S. C., Haft, D. H., Selengut, J., Gwinn, M. L., Zhou, L., Zafar, N., Khouri, H., Radune, D., Dimitrov, G., Watkins, K., O'Connor, K. J. B., Smith, S., Utterback, T. R., White, O., Rubens, C. E., Grandi, G., Madoff, L. C., Kasper, D. L., Telford, J. L., Wessels, M. R., Rappuoli, R., and Fraser, C. M. (2005). Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial "pan-genome". *Proceedings of the National Academy of Sciences of the United States of America*, **102**(39), 13950–13955.
- Thompson, C. C., Emmel, V. E., Fonseca, E. L., Marin, M. A., and Vicente, A. C. P. (2013). Streptococcal taxonomy based on genome sequence analyses. *F1000Research*, **2**, 67.
- Tong, H., Gao, X., and Dong, X. (2003). *Streptococcus oligofermentans* sp. nov., a novel oral isolate from caries-free humans. *International journal of systematic and evolutionary microbiology*, **53**(Pt 4), 1101–1104.