

The Polygraph: A Data Structure for Genome Alignment and Variation Detection

M. Stanley Fujimoto, Cole Lyman and Mark Clement
Department of Computer Science, Brigham Young University
Provo, UT, 84606, USA
sfujimoto@gmail.com

Abstract

Comparing whole genomes and finding variation is an important and difficult bioinformatic task. We present the Polygraph, a data structure for reference-free, multiple whole genome alignment that can be used to identify genomic structural variation. This data structure is built from assembled genomes and preserves the genomic structure from the assembly. It avoids the “hairball” graph structure that can occur in other graph methods such as de Bruijn graphs. The Polygraph can easily be visualized and be used for identification of structural variants. We apply the Polygraph to *Escherichia coli* and *Saccharomyces cerevisiae* for finding Structural Variants.

keywords: genome alignment, comparative genomics, graph, homology, structural variants

1 Introduction

Sequence alignment is one of the most basic tools in bioinformatics. Algorithms for sequence comparison, however, are often limited to short sequences and cannot be applied to whole genome sequences due to computational complexity. Aligning only short sequences captures small, local mutations that occur while leaving large-scale mutations undetected. Complete and accurate whole genome alignment is necessary for understanding evolutionary histories of related organisms.

The genome of an organism can evolve in many ways. Small, local mutations include insertions and deletions (indels) and point substitutions. Large-scale genomic modifications include structural variants (SVs) such as large (> 50 base pair) indels, inversions, duplications and rearrangements such as translocations. As genomes diverge evolutionarily, genomic regions that are ancestrally linked are called homologous. Genome alignment attempts to identify homologous regions amongst a set of genomes.

Previous work in the area of genome alignment has

been limited to pairwise alignment or limited to core-genome identification. Methods such as progressiveMauve and Mugsy rely on all-versus-all progressive alignments when applied to many genomes [2, 1]. Methods such as the Harvest Suite rely on core-genome alignment which is a subset of the genome alignment [13]. Core-genome alignment seeks to find orthologous sequences conserved in all aligned genomes. This process is limiting because an all-or-nothing approach does not allow for relationships that exist between subsets of genomes to appear.

Current algorithms usually align using genome anchoring heuristics based on substring seeds. progressiveMauve and Mugsy are both reference-free genome alignment algorithms that use seed anchors [1, 2]. progressiveMauve relies on local multiple alignments (LMAs) which are maximal unique matches (MUMs) [3] that allow for mismatches and occur in multiple genomes. Mugsy first performs pairwise genome alignment using nucmer [4]. The Harvest Suite’s Parsnp aligns genomes by identifying MUMs using a compressed suffix graph and is designed specifically for microbial genomes. Parsnp does not identify SVs, instead focussing only on identifying core-genome regions. Mugsy and progressiveMauve tend to be conservative in their alignments and miss SVs by preferring a consistent global alignment.

In this work, we present a method for *positional homology multiple genome alignment* [2] that extends our previous work [5]. Genome alignment is made possible by a graph data structure called the Polygraph (PG) which can house multiple genomes and is constructed in a reference-free manner. This data structure contains vertices where homologous regions of genomes are collapsed and edges can show shared recombination events amongst subsets of genomes. Storing multiple genomes in this format facilitates the discovery genomic features useful in comparative genomic analyses.

We demonstrate the efficacy of genome alignments produced by the Polygraph in detecting inversions, translocations and indels. First, we align two yeast (*Saccharomyces cerevisiae*) genomes to verify previ-

ously annotated SVs [10] are identified by the PG. We compare these results to Mugsy, Mauve and the Harvest Suite’s Parsnp module. We then compute the PG for 5 *Escherichia coli* and demonstrate how it can be used to identify conserved regions amongst subsets of genomes. The Polygraph provides a method for storing multiple genomes as a graph that allows for the discovery of structural variants.

2 Methodology

Initially, the Polygraph is a data structure that represents a rough alignment of multiple genomes. It is created by merging the different genomes together on a special set of k-mers (*shared-unique k-mers*). We call this initial alignment rough because only regions we are highly-confident are homologous are merged together. Once the genomes are merged additional context is created that can inform if other regions should be collapsed. The alignment is then refined using the generated context by collapsing additional regions together. Through this process, a graph is formed that has Particulars of the Polygraph are detailed below.

2.1 Graph Properties

A Polygraph $P = (V, E)$ is a directed graph where V is a set of vertices and E is a set of edges. A vertex $v \in V$ represents sequence from one or more genomes, storing chromosome, start and end coordinates. The amount of sequence that v represents can be the same or different for each genome present in the vertex and may be considered to represent a syntenic region. An edge $e \in E$ represent paths that different genomes will take through the graph. Each edge stores an array containing genome identifiers.

2.2 Shared-Unique k-mers

The first step in Polygraph construction is identifying shared-unique k-mers. *Shared-unique k-mers* are k-mers that occur only once within more than one of the individual genomes and are assumed to be homologous. Given a set $G = \{g_i \mid 0 \leq i < n\}$ for n genomes where each g_i is a single genome, a k-mer s is in the set of shared-unique k-mers S if:

$$S = \{s \mid \mathbf{card}(G_s) \geq 2\} \quad (1)$$

where

$$G_s = \{i \mid \mathit{occ}(s, g_i) = 1\} \quad (2)$$

$\mathit{occ}(s, g_i)$ returns the frequency of the k-mer s in genome i and $\mathbf{card}(G_s)$ returns the cardinality or size of G_s .

For example, given two genomes A and B , a shared-unique k-mer is a k-mer that occurs only once in genome

A and once in genome B . If a third genome, C , were present also, a shared-unique k-mer need not occur on all three genomes to be considered shared-unique. K-mers that are not shared-unique are called *common*.

Shared-unique k-mers are similar to the maximal unique matches (MUMs) [3] but are not constrained by having to appear in all species, a shared-unique k-mer may exist in any subset of species. This is powerful because instead of all-or-nothing relationships amongst genomes any sub-grouping is permissible.

Shared-unique k-mers are identified using k-mer counting canonical kmers with Jellyfish [9]. They are then stored in a modified Bloom Filter Trie [6].

Genomes are then collapsed together using the shared-unique k-mers as anchor points. The graph is simplified by merging non-branching paths together to form unitigs.

2.3 Bubble Removal

After the initial graph is formed, it can be refined by collapsing bubbles. *Bubbles* in the Polygraph represent regions in genomes where polymorphisms such as single nucleotide variants (SNVs) and insertions and deletions (indels) have occurred. A bubble b in the Polygraph consists of a start vertex $start$, and end vertex end and set of middle nodes M . The set of bubbles B is defined as:

$$B = \{b \mid \mathit{end} \in \mathit{grandChildren}(start), \\ M = \mathit{children}(start) \cap \mathit{ancestors}(end)\} \quad (3)$$

where $\mathit{grandChildren}(v)$ returns the set of grandchildren vertices of v , $\mathit{children}(v)$ returns the set of children vertices of v and $\mathit{ancestors}(v)$ returns the ancestor vertices of v .

To collapse a bubble b , all sequence from vertices in M are absorbed into $start$. All vertices in M are removed from a graph $start$ and end are connected by a new edge. The vertex $start$ may now contain sequences of heterogeneous lengths. After bubbles are collapsed, unitigging is performed to compress the graph.

2.4 Removing Weak Vertices

We can further constrain the graph by applying a minimum support requirement for vertices where genomes have merged. Support for a vertex is calculated by the amount of sequence in a vertex is shared-unique. When a PG is initially formed, the support for each merged vertex is the length of sequence that a vertex represents. A vertex will represent both non-shared-unique and shared-unique sequence after bubbles are collapsed. To determine support, shared-unique sequence lengths are tracked and stored within

a vertex that contains a collapsed bubble. Removing a weak vertex v requires that $\forall g_i \in \text{genomesPresent}(v)$ a vertex is created with edges connected to the appropriate ancestor and child vertices.

2.5 Reflowing

We further refine the sequence represented by vertices through by *reflowing* the Polygraph. Unmerged vertices may exist that represent sequence homologous sequence but was not put into a merged node because the homologous sequence did not consist of shared-unique k-mers. To remedy this, we perform coarse- and fine-grained graph reflowing by using the PG construction and refinement method in a recursive manner.

Coarse-grained reflowing involves separating the graph into separate connected components and re-running the Polygraph construction algorithm on only the sequence represented in a component. This refines the graph by allowing more k-mers to be identified as shared-unique and for additional regions of the genomes to collapse.

Fine-grained reflowing involves re-running the PG algorithm subgraphs. Regions $R = \{m, N\}$ where the graph can be reflowed are identified by a merged vertex m and a set of neighboring nodes N where

$$N = \{v \mid \text{unmerged}(v), \\ (v \in \text{children}(m) \vee m \in \text{ancestors}(v))\} \quad (4)$$

and where $N = \text{children}(m) \cup \text{ancestors}(m)$ is true. Once a valid region R is identified, subsequences are gathered from genomes and sent through the PG algorithm. The newly constructed graph is then inserted into the graph where R resides, all vertices in R are then removed.

3 Results

The Polygraph, Mugsy and progressiveMauve were tested on three data sets. First, two yeast *Saccharomyces cerevisiae* strains: EC1118 Genoscope 2009 and the reference genome S288C were used to see if annotated SVs could be identified. Next, we applied the polygraph to five *Escherichia coli* genomes and visualized the alignment to demonstrate the PG aligning multiple genomes.

In all cases, we formed a Polygraph for genomes using $k = 90$ and the minimum unitig support for weak vertex removal was set to 540 base pairs (bps). For the yeast data set, PG construction took 26m30s on an Intel Xeon E5-2650v4 @2.20GHz. Mugsy’s and progressiveMauve’s runtimes were fast at 1m02s and 2m05s, respectively, but both failed to identify verified

SVs that the PG found. We examine three notable structural variations discovered by Novo et al. in chromosomes VI, XIV and XV [10].

3.1 Yeast Structural Variants

Novo et al. have documented several structural variations that occur between EC1118 and the reference [10]. They make special note of three large-scale rearrangements that occur in chromosomes VI, XIV and XV. Genomes were downloaded from yeastgenome.org. We apply the Polygraph, Mugsy and progressiveMauve to these genomes to identify structural variants. We also attempted to use Parsnp even though it is designed specifically for microbial genomes but were not able to produce comparable results to the other algorithms when applied to a eukaryotic genome.

3.1.1 Chromosome VI

Novo et al. identified three SVs in EC1118 chromosome VI. First, a 38 kilobase (kb) novel insertion in the left arm telomere. Second, a 12kb translocation from chromosome VIII situated between the 38kb novel insertion and the left telomere. Lastly, a 23kb deletion in the left arm with 5kb of the deletion translocated to chromosome X.

Using the Polygraph we were able to successfully identify the 38kb insertion and 12kb translocation. Specifically, we found the 38kb insertion to be 38,836bps and located at EC1118:VI (FN393068.1) 0–38,836. The 12kb translocation was a bit shorter at 11,046bps originating from Ref:VIII 53,9634–55,6754 and inserted into EC1118:VI (FN393068.1) 38,747–49,793 and was visualized in Figure 1a using Mauve Viewer with MAFFT [7] to produce gapped alignments. A graph visualization of the graph component that contains this SV can be seen in Figure 2. Both progressiveMauve and Mugsy capture the large 38kb insertion but both miss the 12kb translocation (progressiveMauve shown in Figure 1b).

The 23kb Ref:VI deletion with 5kb translocation into EC1118:X was not found by the Polygraph, progressiveMauve or Mugsy. The PG did find a 5kb translocation from Ref:XIV in EC1118:X at the location the 5kb Ref:VI translocation should be. The 5kb translocation came from Ref:XIV 9,739–14,941 and was inserted into EC1118:X (FN393076.1) at 18,6768–19,1969. Neither Mugsy nor progressiveMauve identified this translocation.

We investigated the translocation further by mapping all gene sequences from the reference using to EC1118 with BWA [8]. In the 5kb region where the translocation occurred, we found that there were six genes that

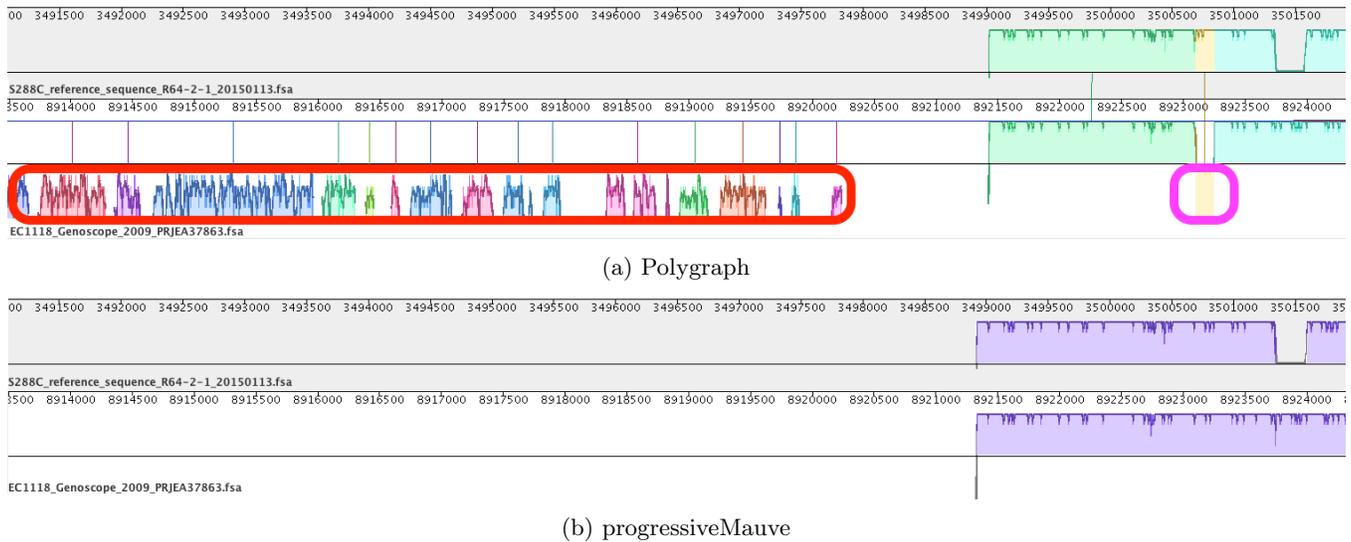


Figure 1: (a) The Polygraph identifies an inverted translocation from chromosome VIII highlighted in red and an inversion in magenta that (b) progressiveMauve does not identify. Visualized using Muave Viewer.

mapped: three from Ref:VI and three from Ref:XIV forming three putative homologous gene pairings that map to the same position in EC1118:X (Figure 3). All six mapped genes had only a handful of polymorphisms compared to the EC1118 sequence.

We then compared the 5kb regions from the genomes through multiple sequence alignment (MSA). They were extracted from:

- EC1118:X 186,768–191,969
- Ref:VI 7,829–13,038
- Ref:XIV 9,739–14,941

MSA was computed using MAFFT [7]. The most notable difference revealed through the MSA was a three base homopolymer thymine deletion in EC1118:X and Ref:XIV (Figure 4). In total, there were five base positions indicating that the Ref:XIV region is more similar to EC1118:X than Ref:VI is.

While this finding contradicts Novo et al.’s statement that the translocation originates from chromosome VI we find sufficient evidence that further investigation on the origins of the translocation is warranted. Additionally, this analysis would not be possible using Mugsy or progressiveMauve as they did not identify it.

3.1.2 Chromosome XIV

This SV is a 17kb novel insertion into Ref:XIV. We found an 18.6kb insertion from EC1118:XIV (FN393084.1) 0–18,654 at the expected location

Ref:XIV 558,235. Both progressiveMauve (18,656bps and Mugsy (18,133bps) identify this insertion as well.

3.1.3 Chromosome XV

This SV is a 65kb replacement of the last 9.7kb in the right arm of Ref:XV. We identified this insertion from EC1118:XV (FN394216.1) 1,045,161–1,110,477 replacing Ref:XV (NC_001147) 1,081,537–1,091,291.

progressiveMauve misidentifies the 9.7kb deletion as a 18.5kb deletion and finds a 6.7kb translocation from Ref:XVI 14,105–18,180 into EC1118:XV 1,036,531–1,040,665 (FN394216.1).

Mugsy identifies the 65kb insertion but misidentifies the 9.7kb deletion. Where the 9.7kb deletion should be, it finds 8 translocations from chromosomes V (FN393065.1), VI (FN393069.1), XII (FN393079.1), XIII (FN393081.1).

3.2 Multiple Genome Alignment

The Polygraph can also be used to align, compare and visualize multiple genomes. We aligned five *Escherichia coli* genomes and have visualized the alignment in Figure 5. Visualization is a convenient feature because conserved regions can be easily identified as well as heterozygous regions which is useful for identifying potential sites for phylogenetic analysis.

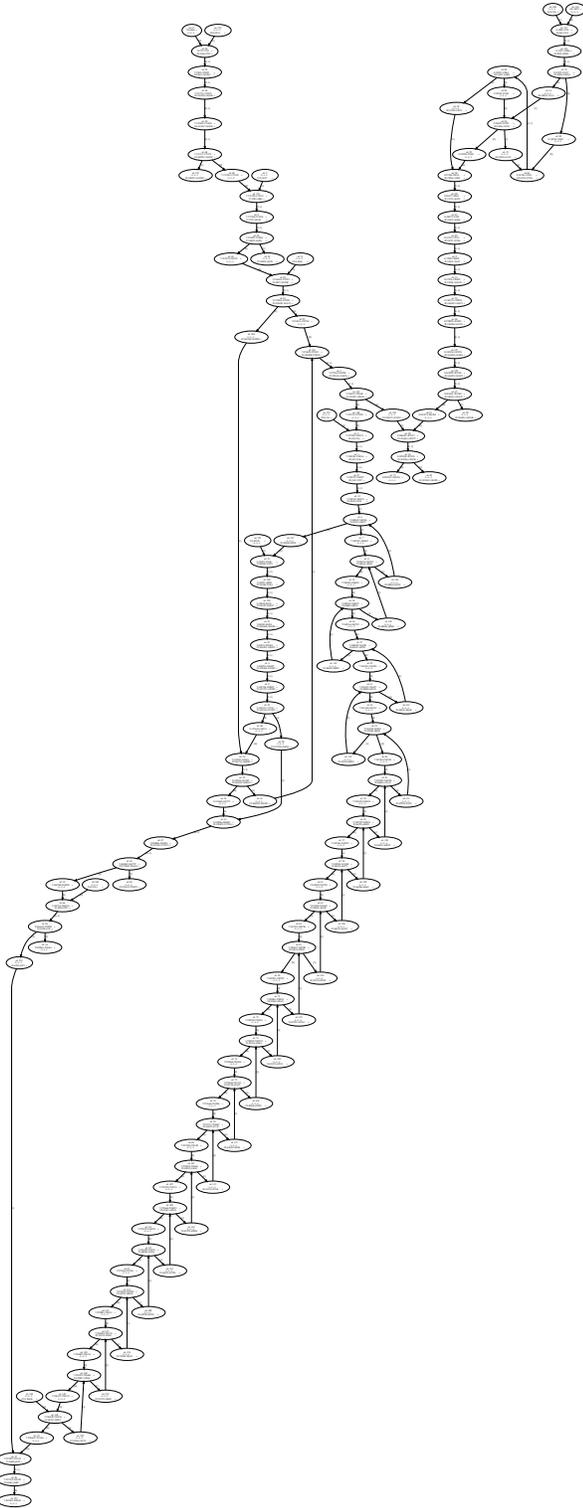


Figure 2: A graph component from the Polygraph for yeast chromosome VI structural variants.



Figure 3: Mapping of the three genes from Ref:VI (YFL059W, YFL060C and YFL061W) on the top row and Ref:XIV (YNL333W, YNL334C and YNL335W) on the bottom in IGV [12].

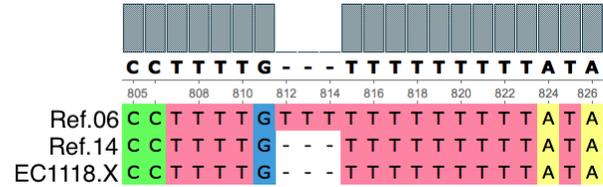


Figure 4: Multiple sequence alignment of EC1118:X, Ref:VI and Ref:XIV in UGENE [11].

4 Discussion

The Polygraph is able to identify numerous structural variants between the two yeast genomes beyond what Novo et al. as well as progressiveMauve and Mugsy were able to identify. Additionally, the resulting graph is small and traversal algorithms can easily be applied. Visual inspection of the PG is simple with yeast-sized genomes and is also human-decipherable. Deeper analysis is easily accomplished as precise genomic coordinates are displayed for each vertex in the graph indicating putative homologous regions.

Of the three structural variants that were indicated by Novo et al., we were able to identify two without caveat with better results compared to Mugsy and progressiveMauve. The main drawback to this fine-grained analysis is runtime. The Polygraph takes significantly more time to run compared to the other software packages. Because the PG is a new algorithm that employs some parallelism there are still many areas where our code efficiency could be increased.

5 Conclusion

In this work we have demonstrated the utility of the Polygraph, a new data structure designed for whole genome comparison and analysis. We have demonstrated the construction and refinement algorithms that can simplify a graph representing two genomes enough to be human-understandable when visualized. We also demonstrated the utility of the Polygraph by applying it to the yeast genome for identifying SVs. We also demonstrated results of the PG when applied to more than two genomes. While superior results are observed, runtime is much longer than similar packages and

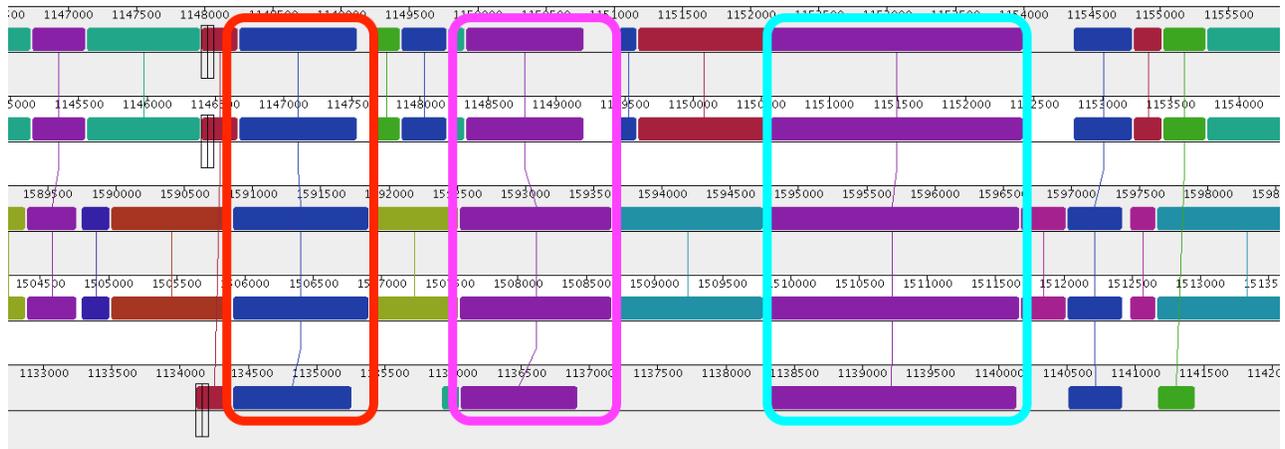


Figure 5: Multiple genome alignment of five *E. coli* genomes with three different homologous regions highlighted in red, magenta and cyan in the Mauve Viewer.

requires additional work.

Our results show that the Polygraph is a viable data structure for comparing genomes. New methods for leveraging new data are necessary, especially as sequencing technology improves and genome assemblies for individuals become prevalent. Using the Polygraph, structural variants can be found, visualized and analyzed easily. As the Polygraph is extended to handle more genomes it can be used for whole genome phylogenetic tree reconstruction as well as identify complex genomic variations for disease association studies.

References

- [1] Samuel V Angiuoli and Steven L Salzberg. Mugsy: fast multiple alignment of closely related whole genomes. *Bioinformatics*, 27(3):334–342, 2010.
- [2] Aaron E Darling, Bob Mau, and Nicole T Perna. progressivemaue: multiple genome alignment with gene gain, loss and rearrangement. *PloS one*, 5(6):e11147, 2010.
- [3] Arthur L Delcher, Simon Kasif, Robert D Fleischmann, Jeremy Peterson, Owen White, and Steven L Salzberg. Alignment of whole genomes. *Nucleic acids research*, 27(11):2369–2376, 1999.
- [4] Arthur L Delcher, Adam Phillippy, Jane Carlton, and Steven L Salzberg. Fast algorithms for large-scale genome alignment and comparison. *Nucleic acids research*, 30(11):2478–2483, 2002.
- [5] M Stanley Fujimoto, Cole Lyman, Anton Suvorov, Paul Bodily, Quinn Snell, Keith Crandall, Seth Bybee, and Mark Clement. Genome polymorphism detection through relaxed de bruijn graph construction. In *Bioinformatics and Bioengineering (BIBE), 2017 IEEE 17th International Conference on*, pages 212–216. IEEE, 2017.
- [6] Guillaume Holley, Roland Wittler, and Jens Stoye. Bloom filter trie: an alignment-free and reference-free data structure for pan-genome storage. *Algorithms for Molecular Biology*, 11(1):3, 2016.
- [7] Kazutaka Katoh and Daron M Standley. Mafft multiple sequence alignment software version 7: improvements in performance and usability. *Molecular biology and evolution*, 30(4):772–780, 2013.
- [8] Heng Li and Richard Durbin. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*, 25(14):1754–1760, 2009.
- [9] Guillaume Marçais and Carl Kingsford. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*, 27(6):764–770, 2011.
- [10] Maite Novo, Frédéric Bigey, Emmanuelle Beyne, Virginie Galeote, Frédéric Gavory, Sandrine Mallet, Brigitte Cambon, Jean-Luc Legras, Patrick Wincker, Serge Casaregola, et al. Eukaryote-to-eukaryote gene transfer events revealed by the genome sequence of the wine yeast *saccharomyces cerevisiae* ec1118. *Proceedings of the National Academy of Sciences*, pages pnas-0904673106, 2009.
- [11] Konstantin Okonechnikov, Olga Golosova, Mikhail Fursov, and Ugene Team. Unipro ugene: a

unified bioinformatics toolkit. *Bioinformatics*, 28(8):1166–1167, 2012.

- [12] James T Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S Lander, Gad Getz, and Jill P Mesirov. Integrative genomics viewer. *Nature biotechnology*, 29(1):24, 2011.
- [13] Todd J Treangen, Brian D Ondov, Sergey Koren, and Adam M Phillippy. The harvest suite for rapid core-genome alignment and visualization of thousands of intraspecific microbial genomes. *Genome biology*, 15(11):524, 2014.