

Constructing founder sets under allelic and non-allelic homologous recombination

Konstantinn Bonnet¹, Tobias Marschall¹, and Daniel Doerr¹

¹*Faculty of Medicine, Heinrich Heine University, Düsseldorf, Germany*

Abstract

Homologous recombination between the maternal and paternal copies of a chromosome is a key mechanism for human inheritance and shapes population genetic properties of our species. However, a similar mechanism can also act between different copies of the same sequence, then called *non-allelic homologous recombination* (NAHR). This process can result in genomic rearrangements—including deletion, duplication, and inversion—and is underlying many genomic disorders. Despite its importance for genome evolution and disease, there is a lack of computational models to study genomic loci prone to NAHR.

In this work, we propose such a computational model, providing a unified framework for both (allelic) homologous recombination and NAHR. Our model represents a set of genomes as a graph, where human haplotypes correspond to walks through this graph. We formulate two founder set problems under our recombination model, provide flow-based algorithms for their solution, and demonstrate scalability to problem instances arising in practice.

1 Introduction

Twenty years ago, at this conference, Esko Ukkonen introduced the problem of inferring founder sets from haplotyped SNP sequences under allelic recombination [30]. Ukkonen’s work has since inspired a wealth of research addressing various aspects and applications of founder set reconstruction ranging from the reconstruction of ancestral recombinations and pangenomics to applications in phage evolution [15, 18, 29]. In its original setting, the problem sets out from a given set of m sequences of equal length n , where characters across sequences residing at the same index position correspond to a SNP. It then asks for a smallest set of sequences, called *founder set*, such that each given sequence can be constructed through a series of crossovers between sequences of the founder set, where each segment between two successive recombinations must meet a minimum length threshold. The *Founder Set Reconstruction* problem is NP-hard in general [22], but is solvable in linear time for the special case of founder sets of size two [30, 35]. Since its introduction, various heuristics and approximations have been proposed [35, 24, 25]. A popular subject of study is a variant that restricts crossovers to coincide across the given set of sequences. The resulting problem, known as *Minimum Segmentation Problem* is polynomial [26]. In his

seminal paper, Ukkonen devised a $O(n^2m)$ algorithm for its solution which has been improved by Norri *et al.* [16] to linear time, i.e. $O(nm)$, capitalizing on recent breakthroughs in data structures [9].

Just like the founder set problem, the vast majority of population genetic analyses and genome-wide association studies have been focused on SNPs in the past decades, neglecting the more complex forms of variation—mostly for technical difficulties in detecting them. In particular, structural variants (SVs), commonly defined as variants of at least 50bp, have posed substantial challenges and studies based on short sequencing reads typically detect less than half of all SVs present in a genome [37]. Recent technological and algorithmic advances help to overcome these limitations [27]. Long read technologies now enable haplotype-resolved *de novo* assembly of human genomes [19], which in turn enables a much more complete ascertainment of SVs [10]. Earlier this year, the first complete telomere-to-telomere assembly of a human genome was announced [17], heralding a new era of genomics where high-quality, haplotype-resolved assemblies of complex repetitive genomic structures become broadly available. Presently, the Human Pangenome Reference Consortium (HPRC), is applying these techniques to generate a large panel of haplotype-resolved genome assemblies from samples of diverse ancestries [33]. These emerging data sets enable studying genetic loci involving duplicated sequence, called *segmental duplications* (SDs), which are amenable to NAHR and are therefore highly mutable and show complicated evolutionary trajectories [13, 31].

Interestingly, at loci with highly similar segments arranged in opposite orientations, such as Segment 3 in Figure 1, NAHR can lead to *inversion*, i.e. the reversal of the interior sequence (Segment 4 in Figure 1). Because of being flanked by the same sequence on both sides, such events have been largely undetectable by sequencing technologies with read lengths below the length of the duplicated sequence; in particular by conventional short read sequencing. Recent studies applying multiple technologies reveal that inversions affect tens of megabases of sequence in a typical human genome [7]. Unlike most other classes of genetic variation, inversions are often *recurrent* with high mutation rates, that is, the same events have happened multiple times in human history [21]. Depending on the structures of duplicated sequence at a particular locus, individual human haplotypes can differ in their potential for NAHR. This can have important implications for the risk for a range of genetic disorders caused by NAHR-mediated mutations [21].

For over two decades, a branch of research dedicated to the study of gene orders has produced various models for the study of genome rearrangements. These range from the classic reversal [3, 2] and transposition [4] model to composed models for two or more balanced rearrangements [32, 8], to generalized models such as the popular *double cut and join* (DCJ) model [36, 5]. As the research in this field continues, advanced models can additionally accommodate one or more types of unbalanced rearrangements, i.e., deletion, insertion, and duplication [28, 6]. Yet, none of these models adequately considers sequence similarity as a prerequisite for NAHR, which is a key molecular mechanism shaping complex loci in the human genome. In summary, there are now technological opportunities to study the population history of recalcitrant SD loci that are prone to genome rearrangements and relevant to disease, but computational models to facilitate this have so far been lacking.

In this work, we study homologous recombination in a genome model that

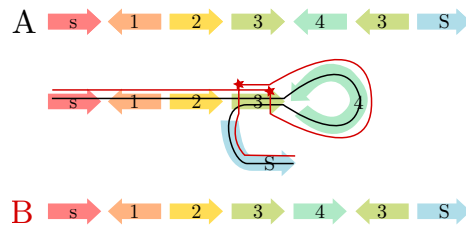


Figure 1: *Illustration of an NAHR-mediated inversion.* Haplotype A (black line) represents the original configuration, while haplotype B (red line) can be derived from A by two recombination events between inverted repeats of genomic marker 3 as indicated by the red stars.

represents DNA sequences at a level of abstraction where they are already decomposed into genomic markers with assigned homologies. Here, our notion of homology is a synonym for *high DNA sequence similarity*, as we adopt the terminology underlying the concept of homologous recombination. It should not be confused with the evolutionary meaning of this term. Our model permits recombination events to occur between homologous markers independent of their position within or between haplotypes, as long as the markers' orientations are respected. In other words, a marker can only recombine with a homologous marker along the same direction, as illustrated by Figure 1. By virtue of recapitulating the underlying molecular mechanism (NAHR), it implicitly allows for all the rearrangements it can give rise to, including deletion, duplication, and inversion.

Marker decomposition and homology assignment can be done in practice with genome graph builders such as MBG [23], minigraph [12], or pgg¹. In fact, our algorithms are based on *variation graph* or *pangenome graph*, where nodes correspond to homologous DNA segments and edges between segments correspond to observed adjacencies in a given set of haplotypes.

2 Methods

2.1 Preliminaries

A (*genomic*) *marker* m is an element of the finite universe of markers denoted by \mathcal{M} , and is associated with a fragment of a double-stranded DNA molecule. Each marker can be traversed in *forward* and *reverse* direction. A marker in forward orientation (which is the default orientation) is traversed from left to right. Overline notation \overline{m} indicates the reversal of a marker m , which is carried out relative to its orientation, i.e., $\overline{\overline{m}} = m$. Similarly, $\overline{\mathcal{M}}$ represents the set of all reverse oriented markers. We designate two forward oriented markers $\{s, S\} \subseteq \mathcal{M}$ as *terminal markers*. In what follows, we study *terminal sequences*, that is, sequences drawn from the alphabet of oriented markers $\mathcal{M} \cup \overline{\mathcal{M}}$ that start with s or \overline{S} , end in S or \overline{s} and do not contain any further terminal markers in between. A terminal sequence can be traversed in forward and reverse direction. A *haplotype* is a terminal sequence that starts with s (*source*) and ends with S (*sink*).

¹<https://github.com/pangenome/pggb>

Example 1. Consider in the following two sequences of genomic markers A and X drawn from the universe of markers $\mathcal{M} = \{\mathbf{s}, 1, 2, 3, 4, \mathbf{S}\}$, where $A = \mathbf{s}\bar{1}234\bar{3}\mathbf{S}$ and $X = \mathbf{s}\bar{1}234\bar{3}\bar{2}\bar{1}\bar{\mathbf{s}}$. Sequence A starts and ends with terminal markers \mathbf{s} and \mathbf{S} , respectively, thus constituting a haplotype drawn from \mathcal{M} . Conversely, X starts with \mathbf{s} and ends in $\bar{\mathbf{s}}$ and therefore is a terminal sequence, but not a haplotype.

Given a sequence A , $|A|$ indicates the length of A which corresponds to the number of A 's constituting elements. \bar{A} defines the *reverse complementation* of sequence A , i.e., the simultaneous reversal of the sequence and its constituting elements. The element at the i th position in sequence A is denoted by $A[i]$. A *segment* of sequence A starting at position i and ending at and including position j , is denoted by $A[i..j]$. In particular, $A[..i] := A[1..i]$ and $A[i..|A|]$ denote the *prefix* and *suffix* of A , respectively. The operator “+” indicates the concatenation of two sequences.

Example 1 (con't). The length of A is $|A| = 7$; its reverse complement is $\bar{A} = \bar{\mathbf{S}}34\bar{3}\bar{2}\bar{1}\bar{\mathbf{s}}$; $A[4..6]$ is a segment of A and corresponds to sequence $3\bar{4}\bar{3}$; The segments $X[..6] = \mathbf{s}\bar{1}234\bar{3}$ and $A[7..] = \mathbf{S}$ are a prefix and a suffix of X and A , respectively; The concatenation of prefix $X[..6]$ and suffix $A[7..]$ results in haplotype $X[..6] + A[7..] = \mathbf{s}\bar{1}234\bar{3}\mathbf{S}$.

A *recombination* is an operation that acts on a shared oriented marker m of any two terminal sequences A and B : let $A[i] = B[j] = m$; recombination $\chi(A, B, i, j)$ produces terminal sequence $C = A[..i] + B[j + 1..]$. For a given set of haplotypes \mathcal{H} , $\text{span}(\mathcal{H})$ denotes the *span*, i.e., the set of all *haplotypes* generated by applying χ on haplotypes \mathcal{H} and their resulting terminal sequences. Accordingly, we also say that “ \mathcal{H} is a *generating set* of $\text{span}(\mathcal{H})$ ”. Conversely, given a (possibly infinite) set of haplotypes \mathcal{S} such that $\mathcal{H} \subseteq \mathcal{S}$, then \mathcal{H} is a generating set of \mathcal{S} iff $\text{span}(\mathcal{H}) = \mathcal{S}$.

Example 1 (con't). Recombination $\chi(A, \bar{A}, 4, 2)$ produces terminal sequence $X = \mathbf{s}\bar{1}234\bar{3}\bar{2}\bar{1}\bar{\mathbf{s}}$. Subsequent recombination $\chi(X, A, 6, 6)$, produces haplotype $B = \mathbf{s}\bar{1}234\bar{3}\mathbf{S}$. If $\{A\}$ is a given set of haplotypes, then $\text{span}(\{A\}) = \{A, B\}$.

In this paper, we study the following two problems:

Problem 1 (Founder Set). Given a set of haplotypes \mathcal{H} , find a generating set $\mathcal{F} \subseteq \text{span}(\mathcal{H})$ such that $\sum_{A \in \mathcal{F}} |A|$ is minimized.

We call a solution to Problem 1 a *founder set* and its members *founder sequences*.

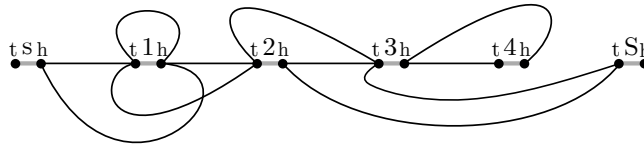
Problem 2. Given a set of haplotypes \mathcal{H} , find a founder set \mathcal{F} that minimizes the number of recombinations applied to haplotypes \mathcal{H} and their intermediate terminal sequences in constructing \mathcal{F} .

2.2 Constructing founder sets

Variation graph construction. We solve Problem 1 by studying the *variation graph* $G_{\mathcal{H}} = (V, E \cup \vec{E})$ of the given set of haplotypes \mathcal{H} . Graph $G_{\mathcal{H}}$ is an undirected edge-colored multigraph where each edge can have one of two colors corresponding to their membership in edge sets E and \vec{E} . In constructing $G_{\mathcal{H}}$,

each marker m of the universe of forward-oriented markers \mathcal{M} is represented by a tuple of its *extremities* (m^t, m^h) also called “tail” and “head” of m , respectively, and its reverse-oriented counterpart \bar{m} is represented as (m^h, m^t) . Node set V of graph $G_{\mathcal{H}}$ corresponds to the set of all marker extremities, and each marker $m \in \mathcal{M}$ gives rise to one *marker edge* $\{m^t, m^h\} \in \vec{E}$. Further, any two (not necessarily distinct) nodes $m_1^b, m_2^c \in V$ are connected by one *adjacency edge* $\{m_1^b, m_2^c\} \in E$ iff there exists a sequence $A \in \mathcal{H} \cup \bar{\mathcal{H}}$ with $A = ..m_1 m_2..$ such that $m_1 = (m_1^a, m_1^b)$, $m_2 = (m_2^c, m_2^d)$ and $\{a, b\} = \{c, d\} = \{t, h\}$.

Example 2. Let $H_1 = s\bar{1}234\bar{3}S$, $H_2 = s111234\bar{3}S$, $H_3 = s\bar{1}234\bar{3}2\bar{3}4\bar{3}S$, and $H_4 = s\bar{1}2S$, then the variation graph $G_{\mathcal{H}}$ of $\mathcal{H} = \{H_1, H_2, H_3, H_4\}$ is as follows, with marker edges drawn in gray and adjacency edges in black:



Proposition 1. Let $G_{\mathcal{H}}$ be the variation graph of haplotypes \mathcal{H} and \mathcal{X} be the set of all alternating walks between terminal markers s^t and S^h in $G_{\mathcal{H}}$, then $\text{span}(\mathcal{H}) = \mathcal{X}$.

Proof.

\Rightarrow Observe that no recombination can create a new pair of consecutive markers $m_1 m_2$ that is not contained in any sequence $A \in \mathcal{H} \cup \bar{\mathcal{H}}$. Therefore, each haplotype $B \in \text{span}(\mathcal{H})$ is a succession of consecutive markers drawn from sequences in $\mathcal{H} \cup \bar{\mathcal{H}}$, i.e., B can be delineated in $G_{\mathcal{H}}$ by following adjacency edges corresponding to its succession of consecutive markers.

\Leftarrow If each alternating walk $X = (s^t, s^h, \dots, S^t, S^h) \in \mathcal{X}$ in variation graph $G_{\mathcal{H}}$ corresponds to a haplotype $B \in \text{span}(\mathcal{H})$, then X must be producible through a series of recombinations of haplotypes \mathcal{H} and their recombinants. We show this by construction:

- (a) Pick some haplotype $A \in \mathcal{H}$ and initialize $i \leftarrow 1$;
- (b) Let $B \in \mathcal{H} \cup \bar{\mathcal{H}}$ be a sequence such that for some position j , $B[j..j+1] = m_1 m_2$ with $m_1 = X[i..i+1]$ and $m_2 = X[i+2..i+3]$. Then $A \leftarrow \chi(A, B, i/2, j)$.
- (c) Increase i by 2 and repeat step **b** unless $i = |X| - 3$.

Observe that by construction of the variation graph $G_{\mathcal{H}}$, a suitable sequence $B \in \mathcal{H} \cup \bar{\mathcal{H}}$ must exist in each iteration of step **b**.

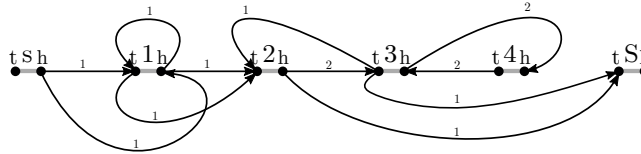
□

Defining flows on variation graphs. We determine a minimum set of founder sequences by solving a network flow problem in variation graph $G_{\mathcal{H}}$ where flow is allowed to travel along adjacency edges in either direction. In doing so, we find a non-negative flow $\phi : V \times V \rightarrow \mathbb{N}$ such that the total flow $\sum_{u,v \in V} \phi(u, v)$ of graph $G_{\mathcal{H}}$ is minimized and satisfies the following constraints:

$$\begin{array}{ll}
 \forall u, v \in V & \phi(u, v) \in \mathbb{N} \quad (\text{constrain flow to integer}) \\
 \forall (u, v) \in \{(u', v') \mid u', v' \in V : \{u', v'\} \notin E\} & \phi(u, v) = 0 \quad (\text{constrain travel of flow}) \\
 \forall v \in V & i(v) := \sum_{u \in V} \phi(u, v) \quad (\text{incoming flow}) \\
 & o(v) := \sum_{u \in V} \phi(v, u) \quad (\text{outgoing flow}) \\
 \forall \{u, v\} \in E & \phi(u, v) + \phi(v, u) \geq 1 \quad (\text{flow coverage}) \\
 & o(s^t) = i(S^h) = 0 \quad (\text{flow direction } s^t \rightarrow S^h) \\
 \forall m \in \mathcal{M} \setminus \{s, S\} & i(m^t) = o(m^h) \quad (\text{flow conservation}) \\
 & i(m^h) = o(m^t)
 \end{array}$$

Note that the flow can travel in both directions of an edge $\{u, v\} \in E$ and that $\phi(u, v) = \phi(v, u)$ does not hold true in general. The only node pairs of the graph that are *unbalanced*, i.e., do not satisfy flow conservation, are (s^t, S^h) and (S^t, S^h) .

Example 1 (cont'd). *The drawing below illustrates a flow solution on variation graph $G_{\mathcal{H}}$, with the direction and amount of flow along adjacency edges indicated by labeled arrowed arcs.*

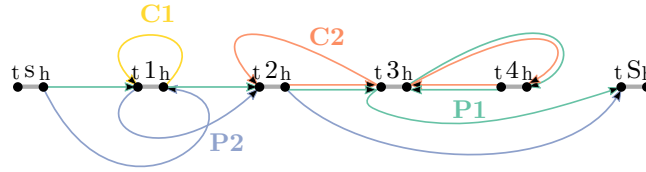


Deriving haplotypes from flows. By applying the Flow Decomposition Theorem [1, p. 80f], any *flow*, i.e., solution to the above-specified constraints, is decomposable into a set of alternating paths going from source s^t to sink S^h and a set of alternating cycles. Ahuja *et al.* [1] give a simple and efficient algorithm that does so in polynomial time and which we describe below, adapted to our circumstances. The idea is to perform a random walk in the graph from source to sink or within a cycle, thereby consuming flow along adjacency edges until all flow is depleted. The proof of the algorithm remains unchanged to that given by Ahuja *et al.*, thus is not repeated here.

1. Set $u \leftarrow s^t$.
2. Setting out from current node u , traverse the incident marker edge to some node v , choose any neighbor w of v for which $\phi(v, w) > 1$. Follow the adjacency edge to w and decrease the flow $\phi(v, w)$ by 1. Set $u \leftarrow w$.
3. As long as $u \neq S^t$ do as follows: if u has been visited in the traversal before, then extract the corresponding alternating cycle from the recorded sequence and report it. Proceed with the traversal by repeating step 2.
4. However, if $u = S^t$, follow the marker edge to S^h and report the recorded sequence as a path.

5. If s^h is incident to edges with positive flow, proceed with step 1. Otherwise, there still might be strictly positive flow remaining in the graph corresponding to unreported cycles. In that case, pick any node $u \leftarrow m^a$ such that for some node w , $\phi(m^b, w) > 0$, $\{a, b\} = \{t, h\}$ and $m \in \mathcal{M}$, and proceed with step 2.

Example 1 (cont'd). *The components of the flow solution on variation graph $G_{\mathcal{H}}$ comprise two cycles $C1$ and $C2$, and two (s^t, S^h) -paths $P1$ and $P2$, as illustrated below.*



What remains is the integration of cycles into walks that then correspond to the haplotypes of the founder set. The integration is facilitated by a graph structure, the *component graph*. The component graph $G' = (V', E', l)$ is an edge-labeled, directed multigraph, where, in its initial construction, each alternating (s^t, S^h) -path and each cycle reported during flow decomposition is represented by a distinct node of V' . In the component graph G' , each cycle c of the flow decomposition sharing one or more markers with another component c' is connected by one or more directed edges (c, c') to that component, with each edge's label $l(c, c')$ corresponding to one distinct shared marker, oriented according to their succession in c (which may not be the same as in c'). The component graph is then successively deconstructed until empty as follows:

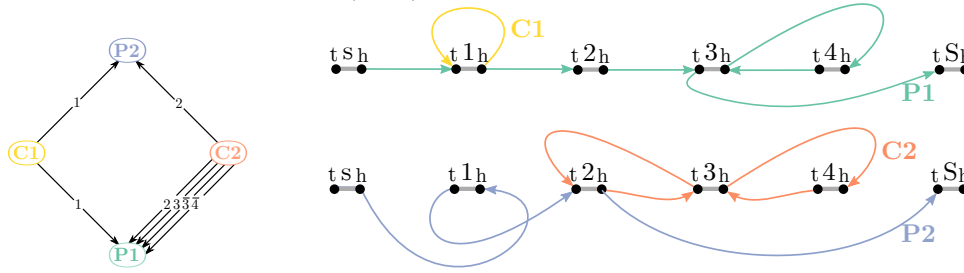
1. Remove and report all (s^t, S^h) -walks with in-degree 0 from node set V'^2 .
2. Pick a cycle $c \in V'$ with in-degree 0, or, if none such exists, any arbitrary cycle $c \in V'$.
3. Pick an outgoing edge $(c, c') \in E'$ such that c' is a (s^t, S^h) -walk. If no such c' exists, c is only adjacent to cycles, out of which one c' is picked at will. Let $(m^a, m^b) \leftarrow l(c, c')$, $\{a, b\} = \{t, h\}$. If marker m is embedded in c' in same orientation, i.e. $c' = ..m^a m^b ..$, then linearize c in m , i.e., $c = m^b c_1 .. c_{k-1} m^a$, and integrate it into c' such that $c' \leftarrow ..m^a m^b c_1 .. c_{k-1} m^a m^b ..$. Otherwise, integrate the reversed linearization of c , i.e. $c' \leftarrow ..m^b m^a c_{k-1} .. c_1 m^b m^a ..$. Remove cycle c and its outgoing edges from component graph G' .
4. Proceed with step 1 until no more components remain and all (s^t, S^h) -walks are reported.

The search for components with in-degree 0 can be efficiently implemented through preorder traversal of G' . Note that each cycle must have at least one outgoing edge and that ultimately all cycles *must be* integrable into a (s^t, S^h) -walk, otherwise this would imply that $G_{\mathcal{H}}$ contains a disconnected, circular component that is not reachable by an alternating path from source s^t to sink S^h ,

²By construction, (s^t, S^h) -walks have out-degree 0, i.e., those with in-degree 0 are singleton in G' .

thus contradicting the correctness of $G_{\mathcal{H}}$'s construction. The reported (s^t, S^h) -walks represent the wanted haplotypes of the founder set.

Example 1 (cont'd). *The plot below depicts the component graph of components $C1$, $C2$, $P1$, and $P2$ (left) and the final two (s^t, S^h) -walks that collectively represent a founder set of \mathcal{H} (right).*



Theorem 1. *Any flow that minimizes the total flow $\sum_{u,v \in V} \phi(u,v)$ of variation graph $G_{\mathcal{H}} = (V, E \cup \vec{E})$ of a given set of haplotypes \mathcal{H} gives rise to a solution to Problem 1.*

Proof. It is sufficient to show that every flow is decomposable into a set of haplotypes (\Rightarrow) and every founder set represents a valid flow (\Leftarrow).

- \Rightarrow Any flow of variation graph $G_{\mathcal{H}}$ is decomposable into a set of haplotypes \mathcal{H}' , as demonstrated above. Observe that the above-listed flow constraints enforce the derived haplotypes \mathcal{H}' to cover the entire graph $G_{\mathcal{H}}$ and consequently $G_{\mathcal{H}'} = G_{\mathcal{H}}$. This implies that $\text{span}(\mathcal{H}') = \text{span}(\mathcal{H})$, i.e., \mathcal{H}' is a generating set of $\text{span}(\mathcal{H})$. Therefore, the sum of lengths of haplotypes derived from a flow solution is an upper bound of Problem 1.
- \Leftarrow Any set of haplotypes $\mathcal{H}' \subseteq \text{span}(\mathcal{H})$ that covers each consecutive pair of markers $m_1 m_2$ in haplotypes \mathcal{H} at least once (either in forward orientation $m_1 m_2$ or in reverse orientation $\overline{m_2 m_1}$) represents a valid flow of $G_{\mathcal{H}}$. To construct a flow from \mathcal{H}' , set $\phi(m_1^a, m_2^b)$ to the number of occurrences of consecutive markers $m_1 m_2$ in haplotypes of \mathcal{H}' with $m_1 = (m_1^a, m_1^b)$ and $m_2 = (m_2^c, m_2^d)$, $\{a, b\} = \{c, d\} = \{t, h\}$. Observe that by construction, flow is integer, travels from source s^t to sink S^h and satisfies coverage and conservation constraints.

□

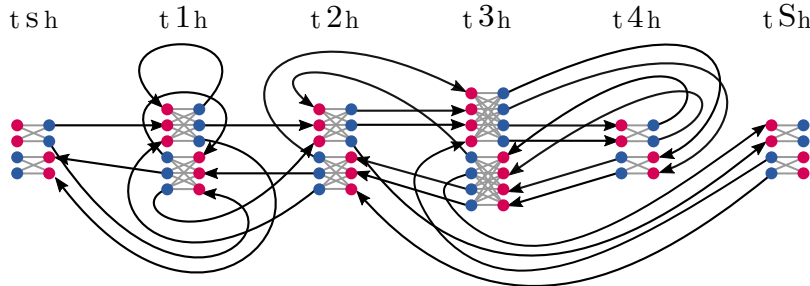
2.3 Minimizing recombinations in founder sequences

We now present an algorithm towards solving Problem 2, i.e., the problem of finding a founder set that minimizes the number of recombinations needed for its construction from a given set of haplotypes \mathcal{H} . Our approach is exact under the assumption that the overall multiplicity of each pair of consecutive markers in the founder set of a solution to Problem 2 is known, yet the pair's particular orientation in a founder sequence may be unresolved. To this end, we presume a given function $\hat{\phi}(m_1, m_2)$ acting as oracle for the overall multiplicity of any given pair of consecutive oriented markers $m_1, m_2 \in \mathcal{M} \cup \overline{\mathcal{M}}$. More

specifically, $\hat{\phi}(m_1, m_2)$ reports the total number of occurrences of $m_1 m_2$ and $\overline{m_2 m_1}$ in a solution to Problem 2. In addition, we make use of function $\mu(m) := \sum_{m' \in \mathcal{M} \cup \overline{\mathcal{M}}} \hat{\phi}(m, m')$ to retrieve the multiplicity of any marker $m \in \mathcal{M} \cup \overline{\mathcal{M}}$. Our solution makes use of the *flow graph* that is defined in the subsequent paragraph. We calculate a matching in the flow graph that describes a set of founder sequences, each corresponding to a succession of segments of haplotypes \mathcal{H} . The objective of the matching is to minimize the total number of these segments across all founder sequences which is equivalent to minimizing the number of recombinations for their construction from haplotype set \mathcal{H} .

Flow graph construction. The flow graph $G_{\mathcal{H}, \hat{\phi}} = (V_{\hat{\phi}}, E_{\hat{\phi}} \cup \overrightarrow{E_{\hat{\phi}}})$ is a directed edge-colored multigraph graph with adjacency edges $E_{\hat{\phi}}$ and marker edges $\overrightarrow{E_{\hat{\phi}}}$, where each marker extremity m^a , $m \in \mathcal{M}$ and $a \in \{t, h\}$, gives rise to $2 \cdot \mu(m)$ elements in node set $V_{\hat{\phi}}$, representing $\mu(m)$ many “in” (i) and $\mu(m)$ many “out” (o) nodes. That is, $V_{\hat{\phi}} = \{i_{m^a}^x \mid m \in \mathcal{M}, a \in \{t, h\}, x \in 1..\mu(m)\} \cup \{o_{m^a}^x \mid m \in \mathcal{M}, a \in \{t, h\}, x \in 1..\mu(m)\}$. Each out node $u \in V_{\hat{\phi}} \setminus (\{o_{m^h}^x \mid 1..\mu(m)\} \cup \{o_{m^t}^x \mid 1..\mu(m)\})$ is incident to *one and only one* directed adjacency edge (u, v) connecting u to some in node v thereby realizing one occurrence of its representing pair of consecutive oriented markers in a founder sequence. Conversely, each forward-oriented marker $m \in \mathcal{M}$ contributes $\mu(m)^2$ many directed marker edges that connect in/tail nodes with out/head nodes, i.e., $\{(i_{m^t}^x, o_{m^h}^y) \mid x, y \in 1..\mu(m)\}$. Analogously, each reverse-oriented marker $\overline{m} \in \overline{\mathcal{M}}$ contributes $\mu(m)^2$ many in/head-to-out/tail-directed marker edges $\{(i_{m^h}^x, o_{m^t}^y) \mid x, y \in 1..\mu(m)\}$.

Example 0 (cont'd). The flow graph $G_{\mathcal{H}, \hat{\phi}}$ for the given set of haplotypes $\mathcal{H} = \{H_1, H_2, H_3, H_4\}$ and a given $\hat{\phi}$ is as follows:



In nodes and out nodes are highlighted in red and blue, respectively. For clarity, the direction of marker edges (gray edges; directed from in to out node) is omitted in the illustration.

Graph decomposition. A perfect matching of marker edges in flow graph $G_{\mathcal{H}, \hat{\phi}}$ produces a set of alternating walks and alternating cycles through $G_{\mathcal{H}, \hat{\phi}}$, yet only half of the graph is eligible to form a solution to Problem 2. More precisely, for each marker $m \in \mathcal{M}$, exactly half of the number of its associated nodes in $V_{\hat{\phi}}$ must be saturated in the matching that we seek, the other half as well as their incident edges must remain unsaturated. Further, we aim to admit only matchings that consist entirely of alternating $(i_{s^t}^x, o_{s^h}^y)$ -walks, because only those correspond to valid haplotypes of $\text{span}(\mathcal{H})$.

³ $\hat{\phi}$ and μ are symmetric w.r.t. the relative orientation of markers, $\hat{\phi}(m_1, m_2) = \hat{\phi}(\overline{m_2}, \overline{m_1})$ and $\mu(m) = \mu(\overline{m})$

At last, we aim to assign to each saturated node $v \in V_{\hat{\phi}}$ a position in some haplotype A of given haplotype set \mathcal{H} . That way, we are able to determine whether the incident adjacency edge serves as continuation of the associated segment in A , or whether the incident saturated marker edge implies a recombination between two distinct segments.

The *integer linear program* (ILP) shown in Algorithm 1 implements the above-stated constraints.

Matching constraints. Each edge and node of flow graph $G_{\mathcal{H},\hat{\phi}}$ is associated with binary variables of \mathbf{x} and \mathbf{y} , respectively, that determine their saturation in a solution (cf. domains D.1 and D.2). Constraint C.01 ensures that each saturated marker edge is incident to saturated nodes. Perfect matching constraints, i.e., constraints that impose each saturated node being incident to exactly one marker edge, are implemented by constraint C.02. Similarly, constraint C.03 ensures that an adjacency edge is saturated iff its incident nodes are saturated. In other words, constraints C.01-C.03 together ensure that each component of the saturated graph corresponds to an alternating path or cycle component (the latter being prohibited by further constraints). The following two constraints C.04 and C.05 control the overall size of the saturated graph. In doing so, they ensure that, in a solution to Problem 2, the number of saturated nodes and adjacency edges matches the postulated multiplicity of markers $\mu(m)$, $m \in \mathcal{M} \cup \overline{\mathcal{M}}$, and pairs of consecutive markers $\hat{\phi}(m_1, m_2)$, $m_1, m_2 \in \mathcal{M} \cup \overline{\mathcal{M}}$, respectively.

Path constraints. Constraints C.05-C.08 force each component of the saturated graph to start and end in nodes associated with source s^t and sink S^h , respectively, thereby ruling out any cycles. To this end, they make use of a set of integer variables \mathbf{f} (cf. Domain D.03) that define an increasing flow within each saturated component that is bounded by constant T corresponding to the total flow of the graph, i.e., $T := \sum_{m \in \mathcal{M}} \mu(m)$. In each saturated marker edge, the flow is increased by 1 while along each adjacency edge, flow is kept constant. This prevents the formation of saturated cycles, because their flow would be infinite. Lastly, constraint C.08 preclude paths from starting in S^h or ending in s^t , leaving only one option for any saturated component open, that is, the formation of a (s^t, S^h) -path.

Haplotype assignment. Each node in a solution to the ILP is associated with exactly one position in a haplotype in \mathcal{H} , recorded by binary variables \mathbf{c} . Moreover, any marker edge whose incident pair of nodes is associated with the same position of the same haplotype corresponds to a conserved segment, i.e, no recombination within this marker has taken place. Each marker edge corresponding to a conserved segment contributes a score unit to the objective function. These score units are encoded by binary variables \mathbf{t} (cf. domain D.05). Constraint C.09 ensures that each marker is associated with exactly one position j in a haplotype A of set $\mathcal{H} \cup \overline{\mathcal{H}}$, while C.10 confines incident nodes of adjacency edges to represent a consecutive marker pair $A[j..j+1]$. At last, constraint C.11 allows \mathbf{t} variables of marker edges to take on value 1 only if that marker edge is saturated and its incident nodes are associated with the same haplotype position.

By maximizing the sum over \mathbf{t} variables, the objective minimizes the total number of segments needed to decompose the calculated founder sequences into segments from haplotypes $\mathcal{H} \cup \overline{\mathcal{H}}$ that are delimited by recombination events.

Algorithm 1 An ILP solution to Problem 2.

Objective:

Maximize

$$\sum_{\substack{(i_{m^a}^x, o_{m^b}^{x'}) \in \overrightarrow{E_{\hat{\phi}}}, \\ A[j] = (m^a, m^b)}} t_{i_{m^a}^x o_{m^b}^{x'}}^{A[j]}$$

Constraints:

$$(C.01) \quad y_u + y_v \geq 2 x_{uv} \quad \forall (u, v) \in \overrightarrow{E_{\hat{\phi}}}$$

$$(C.02) \quad \sum_{(u,v) \in \overrightarrow{E_{\hat{\phi}}}} x_{uv} = y_u \quad \forall \text{ in nodes } u \in V_{\hat{\phi}}$$

$$\sum_{(u,v) \in \overrightarrow{E_{\hat{\phi}}}} x_{uv} = y_v \quad \forall \text{ out nodes } v \in V_{\hat{\phi}}$$

$$(C.03) \quad x_{uv} = y_u \quad \forall (u, v) \in E_{\hat{\phi}}$$

$$x_{uv} = y_v$$

$$(C.04) \quad \sum_{x=1}^{\mu(m)} y_{i_{m^a}^x} + y_{o_{m^a}^x} = \mu(m) \quad \forall m \in \mathcal{M}, a \in \{t, h\}$$

$$(C.05) \quad \sum_{\substack{x, x' \text{ s.t.} \\ (o_{m_1^b}^x, i_{m_2^c}^{x'}) \in E_{\hat{\phi}}}} x_{o_{m_1^b}^x i_{m_2^c}^{x'}} = \hat{\phi}(m_1, m_2) \quad \forall (m_1^b, m_2^c) \text{ s.t. } \hat{\phi}(m_1, m_2) > 0, \\ m_1 = (m_1^a, m_1^b), m_2 = (m_2^c, m_2^d), \\ \{a, b\} = \{c, d\} = \{t, h\}$$

$$(C.06) \quad f_u = f_v \quad \forall (u, v) \in E_{\hat{\phi}}$$

$$(C.07) \quad f_v - f_u + T x_{uv} \leq T + 1 \quad \forall (u, v) \in \overrightarrow{E_{\hat{\phi}}}$$

$$(C.08) \quad f_v = 0 \quad \forall v \in \{o_{s^t}^x \mid x \in 1..\mu(s)\} \cup \{i_{s^h}^x \mid x \in 1..\mu(S)\}$$

$$(C.09) \quad \sum_{\substack{A \in \mathcal{H} \\ A[j]=m}} c_v^{A[j]} = 1 \quad \forall v \in V_{\hat{\phi}}, v \text{ associated with extremities of marker } m$$

$$(C.10) \quad c_{o_{m_1^b}^x}^{A[j]} = c_{i_{m_2^c}^{x'}}^{A[j+1]} \quad \forall (o_{m^b}^x, i_{m^a}^{x'}) \in E_{\hat{\phi}}, A \in \mathcal{H} \cup \overline{\mathcal{H}}, \\ i \in 1..|A| - 1, \text{ s.t. } A[j..j+1] = (m_1^a, m_1^b)(m_2^c, m_2^d)$$

$$(C.11) \quad x_{i_{m^a}^x o_{m^b}^{x'}} + c_{i_{m^a}^x}^{A[j]} + c_{o_{m^b}^{x'}}^{A[j]} \geq 3 t_{i_{m^a}^x o_{m^b}^{x'}}^{A[j]} \quad \forall (i_{m^a}^x, o_{m^b}^{x'}) \in \overrightarrow{E_{\hat{\phi}}}, A \in \mathcal{H} \cup \overline{\mathcal{H}}, \\ i \in 1..|A|, \text{ s.t. } A[j] = (m^a, m^b)$$

Domains:

$$(D.01) \quad x_{uv} \in \{0, 1\} \quad \forall (u, v) \in E_{\hat{\phi}} \cup \overrightarrow{E_{\hat{\phi}}}$$

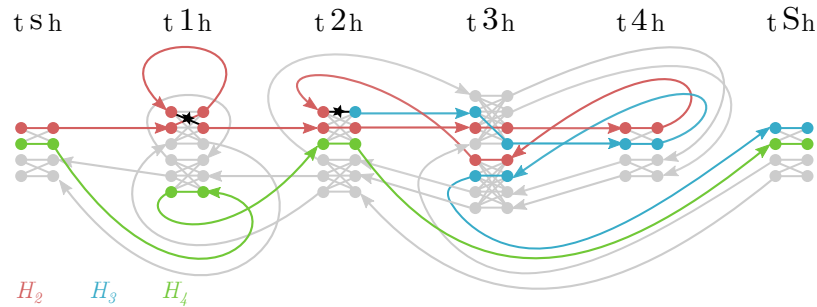
$$(D.02) \quad y_v \in \{0, 1\} \quad \forall v \in V_{\hat{\phi}}$$

$$(D.03) \quad 1 \leq f_v \leq T \quad \forall v \in V_{\hat{\phi}}$$

$$(D.04) \quad c_{i_{m^a}^x}^{A[j]}, c_{o_{m^a}^x}^{A[j]}, c_{i_{m^b}^x}^{A[j]}, c_{o_{m^b}^x}^{A[j]} \in \{0, 1\} \quad \forall A \in \mathcal{H} \cup \overline{\mathcal{H}}, j \in 1..|A|, A[j] = (m^a, m^b), \\ x \in 1..\mu(m)$$

$$(D.05) \quad t_{i_{m^a}^x o_{m^b}^{x'}}^{A[j]} \in \{0, 1\} \quad \forall A \in \mathcal{H} \cup \overline{\mathcal{H}}, j \in 1..|A|, A[j] = (m^a, m^b), \\ x \in 1..\mu(m)$$

Example 0 (cont'd). *The following plot illustrates a matching that is solution to Algorithm 1 for $G_{\mathcal{H},\hat{\phi}}$. Unsaturated nodes and edges are grayed out, haplotype assignments implied by colored paths. The solution features two recombinations, marked by “★” along their associated marker edges.*



3 Results

We implemented our methods in the programming language Rust [14] and used Gurobi [11] as the solver. Our software is *open source* and publicly available on <https://github.com/marschall-lab/hrfs>. To run Algorithm 1 on a given set of haplotypes \mathcal{H} , we estimated the overall multiplicity $\hat{\phi}(m_1, m_2)$ of pairs of consecutive markers $m_1 m_2$ from a network flow solution to Problem 1 on \mathcal{H} . Note that this dispenses Algorithm 1 from being exact in our applications.

All experiments were run on a de.NBI cloud computing machine. For benchmarking purposes, we ran Gurobi single-threaded and recorded wall clock time (in seconds) and *proportional set size* (PSS, in Mb) for memory usage. Optimization time was capped at 30 minutes, beyond which the solver must capitulate and return its best-effort solution found thus far.

3.1 Experimental data

We benchmarked the performance of our algorithms by conducting experiments on both simulated data and a real-world data set. The former presumed a simulator, capable of generating haplotypes with duplicated and inverted markers that can produce intricate homologous recombinations while providing control over the degree of complexity. To this end, we implemented our own, that constructs a single haplotype sequence sampled at random to serve as seed. This seed sequence is adjustable by the following parameters: (i) number of distinct markers, i.e., the size of its variation graph, (ii) ratio of duplications, i.e., the number of additional edges inducing duplications in a walk of the graph, (iii) ratio of inversions, i.e., the proportion of inverted orientations within the set of duplications, and lastly (iv) the number of haplotypes that are input to subsequent founder set reconstruction. The latter are generated by performing random walks in the seed sequence’s variation graph and retaining only those leading from source to sink. Our simulator enables us to explore various parameterizations that match different situations in biological data.

In addition to simulated data, we applied our methods on a biological data set from the human 1p36.13 locus described by Porubsky *et al.* [20] to demonstrate their computational capabilities in realistic instances.

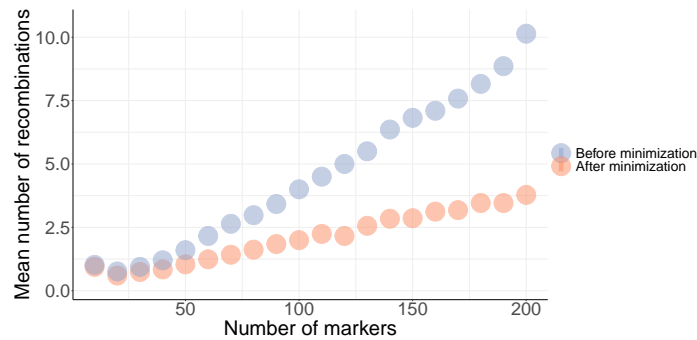


Figure 2: *Mean number of recombinations by the size of the graph.* Experiments were ran with values ranging from 10 to 200 in for the number of markers, in increments of 10. The ratio of duplications and of inversions was fixed to 10%, and number of haplotypes to 10. Each colored dot represents the mean number of recombinations over 50 replicates for one parameter set, after random assignment trials (blue) and after optimization (red).

3.2 Simulation experiments

To assess the impact of parameter configurations on the results, we ran a number of different experiments wherein all but one parameters are fixed. A reasonable choice of constants seemed to be 100 distinct markers, 10% of duplications, 10% of inversions and 10 haplotypes, motivated by our data on the 1p36.13 locus and statistics compiled by Porubsky *et al.* [20].

Reduction in number of recombinations. To evaluate the efficacy of our solution to Problem 2, we compared the number of recombinations returned by Algorithm 1 to that in a solution obtained by our network flow algorithm for Problem 1. While the former is the immediate output of Algorithm 1, additional efforts needed to be made in order to retain the latter. In doing so, we estimate the number of recombinations in the flow solution by random assignment of corresponding segments in the original haplotype set and taking the one with the lowest number in 100k trials. Figure 2 summarizes the outcome of this experiment. Over all, Algorithm 1 found a solution with fewer recombinations in all instances but a few where Gurobi returned barely best-effort solutions after reaching the time limit of 30 minutes, all of which exhibited a gap of at least 100%. The parameter settings in those cases were extremal.

Across all experiments and with a fixed ratio of duplications, inversions and number of haplotypes, the mean estimated number of recombinations both in the initial founder set and after minimization increases linearly with the number of markers, by approximately 4.2 and 2.0 per 100 markers respectively, reaching circa 10 and 3.8 for 200 markers.

Flow solution benchmark. Computing solutions with our network flow algorithm proved to be in almost all of our experiments near-instantaneous. By varying the number of distinct markers, the algorithm's performance begins to deteriorate only with very large instances beyond 100k distinct markers and becomes excruciating for instances above 1M markers. When varying other parameters, we fixed the number of distinct markers to 100k rather than 100. Un-

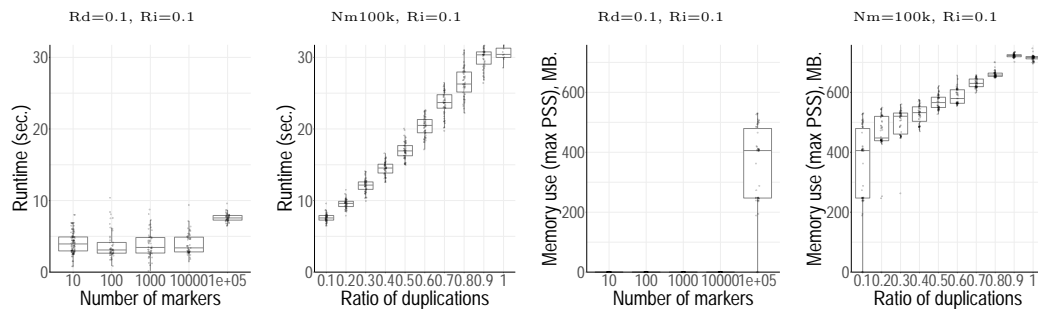


Figure 3: *Flow computation performance.* The plots show the distribution of the runtime in seconds (leftmost two plots) and of peak proportional set size (PSS) in Megabytes (rightmost two plots) of the solver. The abscissa corresponds to different values for one simulation parameter, all others being fixed to the values indicated above. The abbreviations read as follows: Nm, number of markers; Rd, ratio of duplications; and Ri, ratio of inversions.

der 100k markers, execution completes after a mean wall clock time of 3.4 ± 2.0 seconds. In 95% of all experiments, the solver’s runtime was too short to make sufficient measurements for benchmarking memory usage; the maximum PSS for the remaining ones measured at 78MB. Over the 100k mark, both the graph size and duplication ratio begin to reduce performance, with an average runtime of 19.7 ± 8.7 s. We measured peak memory consumption at 758MB across all conditions, which also occurred only at the very extremes of 100k distinct markers and a 100% ratio of duplications (Figure 3).

Recombination minimization benchmark. As shown previously, Algorithm 1 successfully reduces the number of recombinations in solutions to Problem 1. However, its runtime increases dramatically with only moderate increments of any but one parameter of our simulator, the ratio of inversions; it does not play any role in performance. For the remaining three, going beyond instances of 200 distinct markers, 20% of duplications, or 40 haplotypes typically does not allow for the optimization to finish in a reasonable amount of time (Figure 4). A similar but much less pronounced trend is seen with memory usage, which still remains relatively low. Peak memory usage was again observed at extreme parameter values with a PSS of 1072MB with 50 haplotypes.

3.3 Application: Locus 1p36.13

We obtained data from 68 human haplotypes (two per 34 individuals) at the 1p36.13 locus from Porubsky *et al.* [20] and the CHM13 T2T human reference sequence [17]. The sequences comprise only eight distinct markers, terminal markers included. The sequences are attributed to five super populations, out of which 18 are of African origin (AFR), 16 of Eastern Asian (EAS), 12 of Admixed American (AMR), 12 of European (EUR), and 10 are South Asian (SAS). Their variation graph is densely connected with 26 edges (Figure 5). The 68 haplotypes display a high degree of genetic diversity, with haplotype sequences differing in order, orientation, and copy number of the markers. Haplotype lengths in terms of the number of markers vary from 15 to 26, with a median

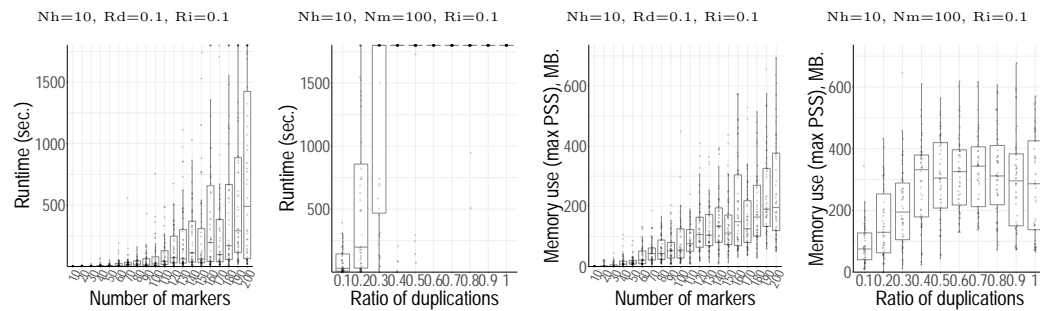


Figure 4: *Solver performance for the minimization of the number of recombinations.* The plots are analogous to the ones in Figure 3, this time measuring performance for the solution to Problem 2. The abbreviations read as follows: N_h , number of haplotypes; N_m , number of markers; R_d , ratio of duplications; and R_i , ratio of inversions.

of 19.

Our network flow algorithm determined that the data set can be generated from a single founder sequence. Our randomized algorithm for calculation of the minimum number of recombinations in a solution to Problem 1 asserted 15 recombinations after 1M trials, while Algorithm 1 obtained an optimal solution that revealed only 9 recombinations. Minimization completed in 60.3 seconds with a peak PSS of 225MB. Note that there exists multiple other co-optimal solutions.

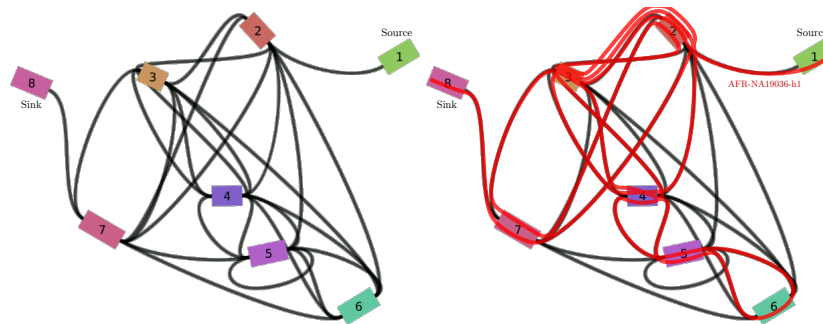


Figure 5: *Graphical representation of the variation graph for the 1p36.13 locus data.* On the left, a 2D plot rendered by Bandage [34]. Markers are represented as numbered colored rectangles, and the undirected edges connecting them as black curves. Markers 1 and 8 correspond respectively to the source and the sink of the graph. The right plot shows the walk through the graph corresponding to the sequence of haplotype NA19036.H1, a sample of African origin in the experimental data. The sample's sequence in the previously established notation is: 123456543273243278.

4 Conclusion

The advent of sequencing technology and genome assembly methodology to reconstruct full human genomes enables research into previously inaccessible segmental duplication loci. This exciting opportunity entails a demand for explanatory models that can infer evolutionary relationships and histories of complex repetitive genomic regions. In this work, we propose a model capable of explaining a broad range of balanced and unbalanced genome rearrangements. Our experiments on simulated data and on the 1p36.13 locus demonstrate that our algorithmic solutions to the founder set problem and the problem of minimizing recombinations in founder sets are capable of processing realistic instances.

Importantly, the model we are proposing is based on a molecular mechanism with a well-established role in shaping segmental duplication architecture. In our view, many past models of genome rearrangements have not sufficiently captured biological reality and there is an important need for further research aiming to incorporate knowledge of molecular mechanisms into such models. For instance, we envision future models that additionally include mechanisms like non-homologous end joining (NHEJ) and mobile element insertions. Furthermore, actual rates at which NAHR occurs depend on factors like the length of the duplicated sequence, the sequence similarity, as well as the presence of specific sequence motifs.

5 Acknowledgements

The authors kindly thank Feyza Yilmaz for providing the haplotype data of the 1p36.13 locus.

6 Competing interests

This work was supported in part by the National Institutes of Health grant 1U01HG010973 to T.M., by the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 956229, and by the BMBF-funded de.NBI Cloud within the German Network for Bioinformatics Infrastructure (de.NBI) (031A532B, 031A533A, 031A533B, 031A534A, 031A535A, 031A537A, 031A537B, 031A537C, 031A537D, 031A538A).

References

- [1] AHUJA, R. K., MAGNANTI, T. L., AND ORLIN, J. B. *Network Flows: Theory, Algorithms, and Applications*, 1 ed. Prentice Hall, Feb. 1993.
- [2] BADER, D. A., MORET, B. M., AND YAN, M. A linear-time algorithm for computing inversion distance between signed permutations with an experimental study. In *1st International Workshop on Algorithms in Bioinformatics (WABI 2001)* (Berlin, Heidelberg, 2001), Algorithms in Bioinformatics, Springer Berlin Heidelberg, pp. 365–376.
- [3] BAFNA, V., AND PEVZNER, P. A. Genome rearrangements and sorting by reversals. *SIAM Journal on Computing* 25, 2 (1996), 272–289.

- [4] BAFNA, V., AND PEVZNER, P. A. Sorting by transpositions. *SIAM Journal on Discrete Mathematics* 11, 2 (1998), 224–240.
- [5] BERGERON, A., MIXTACKI, J., AND STOYE, J. A unifying view of genome rearrangements. In *6th International Workshop on Algorithms in Bioinformatics (WABI 2006)* (Berlin, Heidelberg, 2006), P. Bucher and B. M. E. Moret, Eds., vol. 4175 of *Algorithms in Bioinformatics*, Springer Berlin Heidelberg, pp. 163–173.
- [6] BOHNENKÄMPER, L., BRAGA, M. D., DOERR, D., AND STOYE, J. Computing the rearrangement distance of natural genomes. *Journal of Computational Biology* 28, 4 (2021), 410–431.
- [7] CHAISSON, M. J. P., SANDERS, A. D., ZHAO, X., MALHOTRA, A., PORUBSKY, D., RAUSCH, T., GARDNER, E. J., RODRIGUEZ, O. L., GUO, L., COLLINS, R. L., FAN, X., WEN, J., HANDSAKER, R. E., FAIRLEY, S., KRONENBERG, Z. N., KONG, X., HORMOZDIARI, F., LEE, D., WENGER, A. M., HASTIE, A. R., ANTAKI, D., ANANTHARAMAN, T., AUDANO, P. A., BRAND, H., CANTSILIERIS, S., CAO, H., CERVEIRA, E., CHEN, C., CHEN, X., CHIN, C.-S., CHONG, Z., CHUANG, N. T., LAMBERT, C. C., CHURCH, D. M., CLARKE, L., FARRELL, A., FLORES, J., GALEEV, T., GORKIN, D. U., GUJRAL, M., GURYEV, V., HEATON, W. H., KORLACH, J., KUMAR, S., KWON, J. Y., LAM, E. T., LEE, J. E., LEE, J., LEE, W.-P., LEE, S. P., LI, S., MARKS, P., VIAUD-MARTINEZ, K., MEIERS, S., MUNSON, K. M., NAVARRO, F. C. P., NELSON, B. J., NODZAK, C., NOOR, A., KYRIAZOPOULOU-PANAGIOTOPOULOU, S., PANG, A. W. C., QIU, Y., ROSANIO, G., RYAN, M., STÜTZ, A., SPIERINGS, D. C. J., WARD, A., WELCH, A. E., XIAO, M., XU, W., ZHANG, C., ZHU, Q., ZHENG-BRADLEY, X., LOWY, E., YAKNEEN, S., MCCARROLL, S., JUN, G., DING, L., KOH, C. L., REN, B., FLICEK, P., CHEN, K., GERSTEIN, M. B., KWOK, P.-Y., LANS-DORP, P. M., MARTH, G. T., SEBAT, J., SHI, X., BASHIR, A., YE, K., DEVINE, S. E., TALKOWSKI, M. E., MILLS, R. E., MARSCHALL, T., KORBEL, J. O., EICHLER, E. E., AND LEE, C. Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat. Commun.* 10, 1 (Apr. 2019), 1784.
- [8] DIAS, Z., AND MEIDANIS, J. Genome rearrangements distance by fusion, fission, and transposition is easy. In *spire* (2001), Citeseer, pp. 250–253.
- [9] DURBIN, R. Efficient haplotype matching and storage using the positional burrows–wheeler transform (pbwt). *Bioinformatics* 30, 9 (2014), 1266–1272.
- [10] EBERT, P., AUDANO, P. A., ZHU, Q., RODRIGUEZ-MARTIN, B., PORUBSKY, D., BONDER, M. J., SULOVARI, A., EBLER, J., ZHOU, W., MARI, R. S., YILMAZ, F., ZHAO, X., HSIEH, P., LEE, J., KUMAR, S., LIN, J., RAUSCH, T., CHEN, Y., REN, J., SANTAMARINA, M., HÖPS, W., ASHRAF, H., CHUANG, N. T., YANG, X., MUNSON, K. M., LEWIS, A. P., FAIRLEY, S., TALLON, L. J., CLARKE, W. E., BASILE, A. O., BYRSKA-BISHOP, M., CORVELO, A., EVANI, U. S., LU, T.-Y., CHAISSON, M. J. P., CHEN, J., LI, C., BRAND, H., WENGER, A. M.,

- GHAREGHANI, M., HARVEY, W. T., RAEDER, B., HASENFELD, P., REGIER, A. A., ABEL, H. J., HALL, I. M., FLICEK, P., STEGLE, O., GERSTEIN, M. B., TUBIO, J. M. C., MU, Z., LI, Y. I., SHI, X., HASTIE, A. R., YE, K., CHONG, Z., SANDERS, A. D., ZODY, M. C., TALKOWSKI, M. E., MILLS, R. E., DEVINE, S. E., LEE, C., KORBEL, J. O., MARSCHALL, T., AND EICHLER, E. E. Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* (Feb. 2021).
- [11] GUROBI OPTIMIZATION, L. Gurobi optimizer reference manual, 2019. <http://www.gurobi.com>.
- [12] LI, H., FENG, X., AND CHU, C. The design and construction of reference pangenome graphs with minigraph. *Genome biology* 21, 1 (2020), 1–19.
- [13] MARQUES-BONET, T., GIRIRAJAN, S., AND EICHLER, E. E. The origins and impact of primate segmental duplications. *Trends Genet.* 25, 10 (Oct. 2009), 443–454.
- [14] MATSAKIS, N. D., AND KLOCK II, F. S. The rust language. In *ACM SIGAda Ada Letters* (2014), vol. 34(3), ACM, pp. 103–104.
- [15] NORRI, T., CAZAUX, B., DÖNGES, S., VALENZUELA, D., AND MÄKINEN, V. Founder reconstruction enables scalable and seamless pangenomic analysis. *Bioinformatics* 37, 24 (July 2021), 4611–4619.
- [16] NORRI, T., CAZAUX, B., KOSOLOBOV, D., AND MÄKINEN, V. Minimum Segmentation for Pan-genomic Founder Reconstruction in Linear Time. In *18th International Workshop on Algorithms in Bioinformatics (WABI 2018)* (Dagstuhl, Germany, 2018), vol. 113 of *Leibniz International Proceedings in Informatics (LIPIcs)*, Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, pp. 15:1–15:15.
- [17] NURK, S., KOREN, S., RHIE, A., RAUTIAINEN, M., BZIKADZE, A. V., MIKHEENKO, A., VOLLGER, M. R., ALTEMOSE, N., URALSKY, L., GERSHMAN, A., AGANEZOV, S., HOYT, S. J., DIEKHANS, M., LOGSDON, G. A., ALONGE, M., ANTONARAKIS, S. E., BORCHERS, M., BOUFFARD, G. G., BROOKS, S. Y., CALDAS, G. V., CHEN, N.-C., CHENG, H., CHIN, C.-S., CHOW, W., DE LIMA, L. G., DISHUCK, P. C., DURBIN, R., DVORKINA, T., FIDDES, I. T., FORMENTI, G., FULTON, R. S., FUNGTAMMASAN, A., GARRISON, E., GRADY, P. G. S., GRAVES-LINDSAY, T. A., HALL, I. M., HANSEN, N. F., HARTLEY, G. A., HAUKNES, M., HOWE, K., HUNKAPILLER, M. W., JAIN, C., JAIN, M., JARVIS, E. D., KERPEDJIEV, P., KIRSCH, M., KOLMOGOROV, M., KORLACH, J., KREMITZKI, M., LI, H., MADURO, V. V., MARSCHALL, T., MCCARTNEY, A. M., MCDANIEL, J., MILLER, D. E., MULLIKIN, J. C., MYERS, E. W., OLSON, N. D., PATEN, B., PELUSO, P., PEVZNER, P. A., PORUBSKY, D., POTAPOVA, T., ROGAEV, E. I., ROSENFELD, J. A., SALZBERG, S. L., SCHNEIDER, V. A., SEDLAZECK, F. J., SHAFIN, K., SHEW, C. J., SHUMATE, A., SIMS, Y., SMIT, A. F. A., SOTO, D. C., SOVIĆ, I., STORER, J. M., STREETS, A., SULLIVAN, B. A., THIBAUD-NISSEN, F., TORRANCE, J., WAGNER, J., WALENZ, B. P., WENGER, A.,

- WOOD, J. M. D., XIAO, C., YAN, S. M., YOUNG, A. C., ZARATE, S., SURTI, U., MCCOY, R. C., DENNIS, M. Y., ALEXANDROV, I. A., GERTON, J. L., O'NEILL, R. J., TIMP, W., ZOOK, J. M., SCHATZ, M. C., EICHLER, E. E., MIGA, K. H., AND PHILLIPPY, A. M. The complete sequence of a human genome. *Science* 376, 6588 (Apr. 2022), 44–53.
- [18] PARIDA, L., MELÉ, M., CALAFELL, F., BERTRANPETIT, J., AND CONSORTIUM, G. Estimating the ancestral recombinations graph (arg) as compatible networks of snp patterns. *Journal of Computational Biology* 15, 9 (2008), 1133–1153.
- [19] PORUBSKY, D., EBERT, P., AUDANO, P. A., VOLLGER, M. R., HARVEY, W. T., MARIJON, P., EBLER, J., MUNSON, K. M., SORENSEN, M., SULOVARI, A., HAUKNES, M., GHAREGHANI, M., LANSDORP, P. M., PATEN, B., DEVINE, S. E., SANDERS, A. D., LEE, C., CHAISSON, M. J. P., KORBEL, J. O., EICHLER, E. E., MARSCHALL, T., AND HUMAN GENOME STRUCTURAL VARIATION CONSORTIUM. Fully phased human genome assembly without parental data using single-cell strand sequencing and long reads. *Nat. Biotechnol.* (Dec. 2020).
- [20] PORUBSKY, D., HÖPS, W., ASHRAF, H., HSIEH, P., RODRIGUEZ-MARTIN, B., YILMAZ, F., EBLER, J., HALLAST, P., MARIA MAGGIOLINI, F. A., HARVEY, W. T., HENNING, B., AUDANO, P. A., GORDON, D. S., EBERT, P., HASENFELD, P., BENITO, E., ZHU, Q., (HGSVC), H. G. S. V. C., LEE, C., ANTONACCI, F., STEINRÜCKEN, M., BECK, C. R., SANDERS, A. D., MARSCHALL, T., EICHLER, E. E., AND KORBEL, J. O. Recurrent inversion polymorphisms in humans associate with genetic instability and genomic disorders. *Cell* (2022).
- [21] PORUBSKY, D., HÖPS, W., ASHRAF, H., HSIEH, P., RODRIGUEZ-MARTIN, B., YILMAZ, F., EBLER, J., HALLAST, P., MARIA MAGGIOLINI, F. A., HARVEY, W. T., HENNING, B., AUDANO, P. A., GORDON, D. S., EBERT, P., HASENFELD, P., BENITO, E., ZHU, Q., HUMAN GENOME STRUCTURAL VARIATION CONSORTIUM (HGSVC), LEE, C., ANTONACCI, F., STEINRÜCKEN, M., BECK, C. R., SANDERS, A. D., MARSCHALL, T., EICHLER, E. E., AND KORBEL, J. O. Recurrent inversion polymorphisms in humans associate with genetic instability and genomic disorders. *Cell* (May 2022).
- [22] RASTAS, P., AND UKKONEN, E. Haplotype inference via hierarchical genotype parsing. In *7th International Workshop on Algorithms in Bioinformatics (WABI 2007)* (Berlin, Heidelberg, 2007), R. Giancarlo and S. Hannenhalli, Eds., Algorithms in Bioinformatics, Springer Berlin Heidelberg, pp. 85–97.
- [23] RAUTIAINEN, M., AND MARSCHALL, T. MBG: Minimizer-based sparse de Bruijn Graph construction. *Bioinformatics* 37, 16 (01 2021), 2476–2478.
- [24] ROLI, A., BENEDETTINI, S., STÜTZLE, T., AND BLUM, C. Large neighbourhood search algorithms for the founder sequence reconstruction problem. *Computers & Operations Research* 39, 2 (2012), 213–224.

- [25] ROLI, A., AND BLUM, C. Tabu search for the founder sequence reconstruction problem: A preliminary study. In *Distributed Computing, Artificial Intelligence, Bioinformatics, Soft Computing, and Ambient Assisted Living* (Berlin, Heidelberg, 2009), S. Omatu, M. P. Rocha, J. Bravo, F. Fernández, E. Corchado, A. Bustillo, and J. M. Corchado, Eds., Springer Berlin Heidelberg, pp. 1035–1042.
- [26] SCHWARTZ, R., CLARK, A. G., AND ISTRAIL, S. Methods for inferring block-wise ancestral history from haploid sequences. In *2nd International Workshop on Algorithms in Bioinformatics (WABI 2002)* (Berlin, Heidelberg, 2002), Algorithms in Bioinformatics, Springer Berlin Heidelberg, pp. 44–59.
- [27] SEDLAZECK, F. J., LEE, H., DARBY, C. A., AND SCHATZ, M. C. Piercing the dark matter: bioinformatics of long-range sequencing and mapping. *Nat. Rev. Genet.* (Mar. 2018).
- [28] SHAO, M., LIN, Y., AND MORET, B. M. E. An exact algorithm to compute the double-cut-and-join distance for genomes with duplicate genes. *Journal of Computational Biology* 22, 5 (2015), 425–435.
- [29] SWENSON, K. M., GUERTIN, P., DESCHÊNES, H., AND BERGERON, A. Reconstructing the modular recombination history of staphylococcus aureus phages. *BMC bioinformatics* 14, 15 (2013), 1–9.
- [30] UKKONEN, E. Finding founder sequences from a set of recombinants. In *2nd International Workshop on Algorithms in Bioinformatics (WABI 2002)* (Berlin, Heidelberg, 2002), Algorithms in Bioinformatics, Springer Berlin Heidelberg, pp. 277–286.
- [31] VOLLGER, M. R., GUITART, X., DISHUCK, P. C., MERCURI, L., HARVEY, W. T., GERSHMAN, A., DIEKHANS, M., SULOVAR, A., MUNSON, K. M., LEWIS, A. P., HOEKZEMA, K., PORUBSKY, D., LI, R., NURK, S., KOREN, S., MIGA, K. H., PHILLIPPY, A. M., TIMP, W., VENTURA, M., AND EICHLER, E. E. Segmental duplications and their variation in a complete human genome. *Science* 376, 6588 (Apr. 2022), eabj6965.
- [32] WALTER, M. E. M., DIAS, Z., AND MEIDANIS, J. Reversal and transposition distance of linear chromosomes. In *Proceedings. String Processing and Information Retrieval: A South American Symposium (Cat. No. 98EX207)* (1998), IEEE, pp. 96–102.
- [33] WANG, T., ANTONACCI-FULTON, L., HOWE, K., LAWSON, H. A., LUCAS, J. K., PHILLIPPY, A. M., POPEJOY, A. B., ASRI, M., CARSON, C., CHAISSON, M. J. P., CHANG, X., COOK-DEEGAN, R., FELSENFELD, A. L., FULTON, R. S., GARRISON, E. P., GARRISON, N. A., GRAVES-LINDSAY, T. A., JI, H., KENNY, E. E., KOENIG, B. A., LI, D., MARSCHALL, T., MCMICHAEL, J. F., NOVAK, A. M., PURUSHOTHAM, D., SCHNEIDER, V. A., SCHULTZ, B. I., SMITH, M. W., SOFIA, H. J., WEISSMAN, T., FLICEK, P., LI, H., MIGA, K. H., PATEN, B., JARVIS, E. D., HALL, I. M., EICHLER, E. E., HAUSSLER, D., AND HUMAN PANGENOME REFERENCE CONSORTIUM. The human pangenome project:

- a global resource to map genomic diversity. *Nature* *604*, 7906 (Apr. 2022), 437–446.
- [34] WICK, R. R., SCHULTZ, M. B., ZOBEL, J., AND HOLT, K. E. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics* *31*, 20 (06 2015), 3350–3352.
- [35] WU, Y., AND GUSFIELD, D. Improved algorithms for inferring the minimum mosaic of a set of recombinants. In *Combinatorial Pattern Matching* (Berlin, Heidelberg, 2007), B. Ma and K. Zhang, Eds., Springer Berlin Heidelberg, pp. 150–161.
- [36] YANCOPOULOS, S., ATTIE, O., AND FRIEDBERG, R. Efficient sorting of genomic permutations by translocation, inversion and block interchange. *Bioinformatics* *21*, 16 (2005), 3340–3346.
- [37] ZHAO, X., COLLINS, R. L., LEE, W.-P., WEBER, A. M., JUN, Y., ZHU, Q., WEISBURD, B., HUANG, Y., AUDANO, P. A., WANG, H., WALKER, M., LOWTHER, C., FU, J., GERSTEIN, M. B., DEVINE, S. E., MARSCHALL, T., KORBEL, J. O., EICHLER, E. E., CHAISSON, M. J. P., LEE, C., MILLS, R. E., BRAND, H., AND TALKOWSKI, M. E. Expectations and blind spots for structural variation detection from long-read assemblies and short-read genome sequencing technologies. *Am. J. Hum. Genet.* (Mar. 2021).