



## Pangenome-based genome inference using integer programming

Ghanshyam Chandra, Md Helal Hossen, Stephan Scholz, et al.

*Genome Res.* published online August 21, 2025

Access the most recent version at doi:[10.1101/gr.280567.125](https://doi.org/10.1101/gr.280567.125)

---

**P<P** Published online August 21, 2025 in advance of the print journal.

**Accepted Manuscript** Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

**Creative Commons License** This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

**Email Alerting Service** Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

# Pangenome-based genome inference using integer programming

Ghanshyam Chandra<sup>1</sup>, Md Helal Hossen<sup>2</sup>, Stephan Scholz<sup>3,4</sup>, Alexander T Dilthey<sup>3,4</sup>, Daniel Gibney<sup>2</sup>,  
Chirag Jain<sup>1\*</sup>

<sup>1</sup>Department of Computational and Data Sciences, Indian Institute of Science, Bangalore KA 560012, India

<sup>2</sup>Department of Computer Science, The University of Texas at Dallas, TX 75080, USA

<sup>3</sup>Institute of Medical Microbiology and Hospital Hygiene, Heinrich Heine University Düsseldorf, Düsseldorf,  
Germany

<sup>4</sup>Center for Digital Medicine, Heinrich Heine University Düsseldorf, Düsseldorf, Germany

\*corresponding author: [chirag@iisc.ac.in](mailto:chirag@iisc.ac.in)

**Abstract.** Affordable genotyping methods are essential in genomics. Commonly used genotyping methods primarily support single nucleotide variants and short indels but neglect structural variants. Additionally, accuracy of read alignments to a reference genome is unreliable in highly polymorphic and repetitive regions, further impacting genotyping performance. Recent works highlight the advantage of haplotype-resolved pangenome graphs in addressing these challenges. Building on these developments, we propose a rigorous alignment-free genotyping method. Our optimization framework identifies a path through the pangenome graph that maximizes the matches between the path and substrings of sequencing reads (e.g.,  $k$ -mers) while minimizing recombination events (haplotype switches) along the path. We prove that this problem is NP-Hard and develop efficient integer-programming solutions. We benchmarked the algorithm using downsampled short-read datasets from homozygous human cell lines with coverage ranging from  $0.1\times$  to  $10\times$ . Our algorithm accurately estimates complete major histocompatibility complex (MHC) haplotype sequences with small edit distances from the ground-truth sequences, providing a significant advantage over existing methods on low-coverage inputs. While this algorithm is designed for haploid genomes, we discuss directions for extending it to diploid genotyping.

## Introduction

Many initiatives are in progress for building haplotype-resolved pangenome references of human and non-human species (Gao et al., 2023; Liao et al., 2023; Smith et al., 2023). Among many applications, pangenome graphs can enable cost-effective genotyping and imputation of a wide spectrum of variant classes beyond single nucleotide polymorphisms (SNPs) and short indels (Harris et al., 2024). Pangenome graphs represent sequence alignment of high-quality fully-phased genome assemblies of individuals from diverse populations (Baaijens et al., 2022). A pangenome graph can be represented as either cyclic or acyclic directed graph where the vertices are labeled with sequences. Paths in this graph spell the reference haplotype sequences and their recombinations. The graph-based representation is flexible enough to incorporate single-nucleotide polymorphisms (SNPs), indels (short insertions and deletions), large structural variants (SVs), nested variants, gene absence/presence, etc. (Computational Pan-Genomics Consortium, 2018).

Recent works propose the use of pangenome references to improve genotyping accuracy from short-read sequencing data (Bradbury et al., 2022; Ebler et al., 2022; Eggertsson et al., 2017; Grytten et al., 2022; Hickey et al., 2020; Letcher et al., 2021; Mun et al., 2023; Sibbesen et al., 2018). Especially for SVs, these methods are an effective alternative to the conventional genotyping methods that are based on aligning reads to a single reference because short-read alignments can be inaccurate for the reads originating from SVs (Ebert et al., 2021; Mahmoud et al., 2019). Methods such as PRG (Dilthey et al., 2015), Pangenie (Ebler et al., 2022) and KAGE (Grytten et al., 2022), utilize  $k$ -mer statistics to infer paths in the graph that correspond to the target genome. These methods compare the  $k$ -mers surrounding a variant site in the graph with the  $k$ -mer counts in the sequencing data to calculate likelihoods of reference and alternative alleles. Pangenie and KAGE also use the long-range haplotype information available in the haplotype-resolved pangenome references. The other approach used in methods such as Giraffe (Sirén et al., 2021) and Graphtyper (Eggertsson et al., 2017) involves aligning reads to a pangenome graph.

There have been efforts on improving the accuracy of read alignments to pangenome graphs as well. A large combinatorial search space in terms of the number of candidate paths in a pangenome graph increases ambiguity during read alignment. This issue has motivated methods that either impute a personalized reference genome (Vaddadi et al., 2024), sample variants (Jain et al., 2021; Pritt et al., 2018; Tavakoli et al., 2022) to obtain a smaller graph, or prioritize the use of reference haplotypes in the graph during alignment (Chandra et al., 2024; Mustafa et al., 2024; Sirén et al., 2021). Our previous work proposed haplotype-aware sequence alignment to graphs by introducing penalties for haplotype switches in an alignment (Chandra et al., 2024). A recent feature added to VG allows sampling of reference haplotypes and their recombinations

from the graph that are most relevant to the target genome using a  $k$ -mer-based greedy heuristic (Sirén et al., 2024).

Low-coverage sequencing, combined with genotyping and phasing, is a cost-effective approach to conduct large-scale genetic studies (Davies et al., 2021; Li et al., 2021; Martin et al., 2021; Rubinacci et al., 2021). In this paper, we develop a rigorous formulation and algorithms for genotyping using pangenome references. Our framework is also applicable to low-coverage short-read sequencing data (coverage  $0.1 - 1\times$ ). Similar to the standard Li and Stephens model (Li and Stephens, 2003), we view the target genome as an imperfect mosaic of the reference haplotypes. Our contributions are as following.

- We introduce a novel optimization framework to estimate the complete haplotype sequence of a haploid genome by determining an appropriate path in the pangenome graph. The objective is to maximize the number of shared substrings (e.g.,  $k$ -mers or minimizers) between the sequencing data and the sequence spelled by the path. We permit recombinations in the path, subject to a fixed penalty per recombination. We refer to this problem as *Path Inference Problem* (formally defined in Methods).
- We prove that the Path Inference Problem is NP-hard, even when restricted to binary alphabets.
- To solve this problem, we develop two integer-programming solutions which involve linear and quadratic constraints, respectively. The two solutions involve a tradeoff between runtime and memory usage.
- We demonstrate the utility of this framework by testing it on downsampled short-read datasets from five human haploid cell lines (coverage  $0.1 - 10\times$ ). For these five samples, complete major histocompatibility complex (MHC) haplotype sequences have been previously determined using long-read assembly (Houwaart et al., 2023). As our pangenome reference, we used a haplotype-resolved pangenome directed acyclic graph (DAG) of 49 MHC haplotype sequences (Li, 2022). We chose MHC region for evaluation because this is the most polymorphic and gene-rich region of the human genome (Dilthey, 2021). The length of this region is about 5 Mbp.
- Using datasets with  $0.1\times$  coverage, our algorithm outputs MHC sequences that are up to 99.96% identical to the ground-truth sequences. It compares favorably to the existing methods.

## Methods

### Overview

Our method, PHI (**P**angenome-based **H**aplotype **I**nference), takes as input a pangenome graph reference and short-read sequencing data from a target haploid genome. We assume that the given pangenome graph

has been constructed using high-quality, haplotype-resolved genome assemblies. This graph is modeled as a directed acyclic graph, where each node is labeled with a DNA sequence, and edges represent the adjacency of these sequences along known haplotypes. Any path through this graph can represent an existing haplotype or a recombinant mosaic of segments from multiple haplotypes.

In the Li-Stephens haplotype model (Li and Stephens, 2003), a new haplotype is formed by copying segments from reference haplotypes, occasionally switching between them. Inspired by this model, we aim to identify a path through the graph that best explains the observed sequencing reads by maximizing the number of shared  $k$ -mers between the reads and the sequence encoded by the chosen path. To account for recombination, our method allows switching between haplotypes in the graph, but imposes a fixed penalty for each recombination event. We refer to this computational task as the Path Inference Problem (illustrated in Figure 1).

We show that solving this optimization task exactly is computationally intractable (NP-hard). To address this, we develop two efficient integer programming algorithms that leverage structural properties of the problem. A key component in our method is the construction of an *expanded graph*, which is a modified version of the pangenome graph in which each reference haplotype is represented as a separate path. All paths start and end at *source* and *sink* vertices respectively. Within this expanded graph, the Path Inference Problem can be reinterpreted as a constrained network flow problem: finding a unit flow from source to sink that minimizes the combined penalties for recombinations and unmatched  $k$ -mers. In the following subsections, we describe our algorithms formally.

## Notations and Problem Formulation

Let  $G(V, E, \sigma, \mathcal{H})$  denote a directed acyclic graph (DAG) representing a haplotype-resolved pangenome reference. Function  $\sigma$  assigns a string label over alphabet  $\Sigma = \{A, C, G, T\}$  to each vertex. A path  $(u_1, u_2, \dots, u_n)$  in  $G$  spells string  $\sigma(u_1) \circ \sigma(u_2) \circ \dots \circ \sigma(u_n)$ , where  $s_1 \circ s_2$  denotes the concatenation of strings  $s_1$  and  $s_2$ .  $\mathcal{H} = \{h_1, h_2, \dots, h_{|\mathcal{H}|}\}$  denotes a set of paths in  $G$  such that each of these paths spells a reference haplotype sequence used in the pangenome reference. We refer to these paths as haplotype paths. We assume that each haplotype path is described by an array, i.e.,  $h_i[1]$  is the first vertex in  $h_i$ ,  $h_i[2]$  is the second vertex in  $h_i$ , etc. The length of a haplotype path  $h_i$ , that is, the count of vertices in  $h_i$  is denoted as  $|h_i|$ . The set of haplotype paths covering vertex  $v \in V$  is denoted as  $haps(v)$ . We assume that, for each edge  $(u, v) \in E$ , there exists a haplotype path  $h_i \in \mathcal{H}$  such that  $u$  and  $v$  are consecutive vertices in  $h_i$ . In other words, each edge is

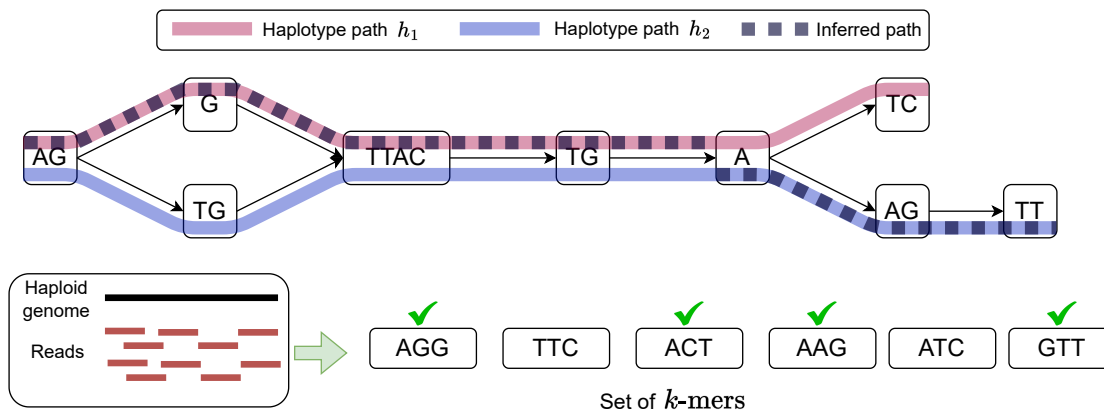


Fig. 1: **A toy example to illustrate the Path Inference Problem.** At the top, we show a haplotype-resolved pangenome graph containing two haplotype paths  $h_1$  and  $h_2$ . These paths are shown in pink and blue colors respectively. An inferred path containing a single recombination is illustrated using a thick dashed line. At the bottom, we show a set of  $k$ -mers observed in the sequencing reads assuming  $k = 3$ . A  $k$ -mer is marked with a green tick if it is a substring of the sequence spelled by the inferred path **AGGTTACTGAAGTT**. Two  $k$ -mers (**TTC** and **ATC**) are not present in this sequence. Our optimization framework identifies an optimal inferred path with minimum cost given user-defined penalties for recombinations and for missing  $k$ -mers.

supported by at least one haplotype path. The functions  $N^+(v)$  and  $N^-(v)$  denote the sets of out-neighbors and in-neighbors, respectively, of the vertex  $v$ , for all  $v \in V$ .

Genotyping a haploid genome sample can be framed as finding a path in the pangenome graph that contains the sample's variants (Paten et al., 2017). Such a path can be modeled as a mosaic path formed by recombinations between the haplotype paths of the graph. Based on this intuition, we define an *inferred path*  $\mathcal{P}$  of length  $n$  as an ordered set  $(a_1, a_2, \dots, a_n)$ , where each  $a_i$  is a two-tuple  $(u, h)$ ,  $u \in V, h \in \mathcal{H}$ . For a tuple  $a_i$ , we use  $a_i.u$  to denote the respective vertex and  $a_i.h$  to denote the respective haplotype path. By using this notation, we keep track of the haplotype path being used alongside the vertex indices. In an inferred path,  $(a_i.u, a_{i+1}.u) \in E$  for all  $i \in [1, n)$ . If  $a_i.h = a_{i+1}.h$ , then we also require that  $a_i.u$  and  $a_{i+1}.u$  must be consecutive vertices in haplotype path  $a_i.h$ . The string spelled by path  $\mathcal{P}$  is  $\sigma(a_1.u) \circ \sigma(a_2.u) \circ \dots \circ \sigma(a_n.u)$ . We denote this string as  $\sigma(\mathcal{P})$ , with a slight abuse of notation.

We say a *recombination*, or a haplotype switch, occurs between two consecutive vertices  $a_i.u$  and  $a_{i+1}.u$  in  $\mathcal{P}$  if  $a_i.h \neq a_{i+1}.h$ . We use  $\gamma(\mathcal{P})$  to denote the count of recombinations in  $\mathcal{P}$ . Since an inferred path must be a path in the pangenome graph by our definition, this setup models homologous recombination events during meiosis. Our goal is to determine the best path in the pangenome graph that is consistent with the set of strings (e.g.,  $k$ -mers) present in the sequencing reads. We define the problem as follows.

*Problem 1 (Path Inference Problem).*

- **Input:** A haplotype-resolved pangenome DAG  $G = (V, E, \sigma, \mathcal{H})$ , a set of strings  $\mathcal{S}$  from the target genome, and a non-negative integer  $c$  indicating recombination penalty.
- **Output:** An inferred path  $\mathcal{P}$  such that

$$Cost(\mathcal{P}) = c \cdot \gamma(\mathcal{P}) + \sum_{r \in \mathcal{S}} \bar{\chi}(r, \sigma(\mathcal{P}))$$

is minimized, where  $\bar{\chi}(r, \sigma(\mathcal{P})) = 0$  if string  $r$  occurs as a substring of string  $\sigma(\mathcal{P})$  and 1 otherwise.

The intuition behind our formulation is to maximize the number of string matches along the inferred path while minimizing the number of recombinations. This approach yields an inferred path that incorporates the majority of strings from  $\mathcal{S}$  as a substring with a finite number of recombinations, constrained by a recombination penalty  $c$ . Set  $\mathcal{S}$  in the above formulation can be set of either  $k$ -mers or minimizers observed in the sequencing reads. The optimization is robust to sequencing errors. Errors in reads typically introduce novel  $k$ -mers that are not present anywhere in the graph. Since such erroneous  $k$ -mers cannot be matched to any path in the graph, the associated cost is unavoidable. As a result, they do not affect the optimization outcome.

We assume a constant cost per recombination. Ideally, the model should use variable, locus-specific recombination costs based on known recombination rates (Petes, 2001). However, recombination maps analogous to those used on linear reference genomes are not yet available for pangenome graphs. Next, we consider the complexity of Path Inference Problem. We show that the problem is NP-hard.

**Theorem 1.** *Problem 1 is NP-hard. This holds for any value of  $c = |V|^{\Theta(1)}$  and even when  $\Sigma = \{0, 1\}$ .*

We refer the reader to Supplemental Note S1 for a proof. This result suffices to establish the theoretical hardness of finding an optimal inferred path under the given cost function. The hardness result holds specifically for our formulation (Problem 1) and does not apply to other possible formulations for genome inference on pangenome graphs. Note that there exist many versions of Problem 1 that could be practically useful and whose computational complexity remains open. For example, when a small (e.g., polylogarithmic) bound is given on the number of allowed recombinations in the inferred path or when the set of haplotype paths in the pangenome graph is small. We leave this investigation for future research.

## Construction of an Expanded Graph

Before developing our integer programming solutions to Problem 1, it is first helpful to define an additional graph representation, which we call as *expanded graph*. In pangenome graphs, multiple haplotype paths

share vertices if the sequences are conserved, whereas in the expanded graph, we will split all haplotypes into separate paths (Figures 2A, 2B). The expanded graph enables us to model Problem 1 as a sort of network flow problem. In particular, the inferred path will be reconstructed from a flow of value one in the expanded graph. We will assign weights to edges to account for recombination penalty. Additional constraints will be used to capture how many strings in  $\mathcal{S}$  occur in the resulting inferred path.

In an inferred path  $\mathcal{P} = (a_1, \dots, a_n)$ , we call a recombination  $a_i.h \neq a_{i+1}.h$  *useful* except when  $a_i.u$  and  $a_{i+1}.u$  are consecutive vertices in haplotype path  $a_i.h$ . Lemma 1 allows us to only consider useful recombinations when finding an optimal solution to Problem 1.

**Lemma 1.** *There exists an optimal inferred path  $\mathcal{P} = (a_1, \dots, a_n)$  for Problem 1 such that all recombinations in  $\mathcal{P}$  are useful.*

*Proof.* Suppose there is an optimal inferred path  $\mathcal{P} = (a_1, \dots, a_n)$  for Problem 1 where for some  $a_i, a_i.h \neq a_{i+1}.h$  such that  $a_i.u$  and  $a_{i+1}.u$  are consecutive vertices in haplotype path  $a_i.h$ . Furthermore, suppose we start with the smallest  $i$  where this holds. We then change the haplotype path for  $a_{i+1}$  to equal  $a_i.h$ . This does not increase the overall cost, since the number of string  $\mathcal{S}$  occurring in  $\sigma(\mathcal{P})$  has not changed, and the number of recombinations either decreases or stays the same. Continuing this process from the next  $j > i$ , such that  $a_j.h \neq a_{j+1}.h$  and  $a_j.u$  and  $a_{j+1}.u$  are consecutive vertices in  $a_j.h$ , we achieve an inferred path satisfying the conditions stated in the lemma after at most  $n$  iterations.  $\square$

Next, we present a definition of the expanded graph where we will consider only the useful recombinations. For technical reasons, we preprocess each edge in  $E$ , splitting it and adding a new vertex labeled with the empty string  $\varepsilon$ . Each added vertex inherits the haplotype paths which supported the edge it was formed from. This added step is to prevent recombinations from a haplotype to itself when we build our expanded graph. Now, let  $V = \{u_1, \dots, u_n\}$ . For haplotype path  $h_j \in \mathcal{H}$ , let  $u_{h_j[i]}$  denote the  $i^{\text{th}}$  vertex in haplotype path  $h_j$ . We use  $G_E = (V_E, E_E, \sigma_E)$  to denote the expanded graph. In  $G_E$ , vertices are string-labeled and edges are weighted. Vertex set  $V_E$  is defined as:

$$V_E = \{s\} \cup \{t\} \cup \{u_{h_j[i]}^j \mid 1 \leq j \leq |\mathcal{H}|, 1 \leq i \leq |h_j|\} \quad (1)$$

The vertex set contains a source and sink vertex,  $s$  and  $t$ , respectively. The vertex set also contains a set of disjoint vertices for each haplotype path in  $\mathcal{H}$  (Figures 2A, 2B). A superscript is used to indicate which haplotype path the vertex is designated to. We refer to the ordered vertex set  $u_{h_j[1]}^j \cdots u_{h_j[|h_j|]}^j$  as a haplotype path in  $G_E$ .

We denote weighted edges in  $E_E$  as tuples of the form  $(start, end, weight)$ . The weighted edge set is

$$E_E = \{(s, u_{h_j[1]}^j, 0) \mid 1 \leq j \leq |\mathcal{H}|\} \quad (2)$$

$$\cup \{(u_{h_j[|h_j|]}^j, t, 0) \mid 1 \leq j \leq |\mathcal{H}|\} \quad (3)$$

$$\cup \{(u_{h_j[i]}^j, u_{h_j[i+1]}^j, 0) \mid 1 \leq j \leq |\mathcal{H}|, 1 \leq i < |h_j|\} \quad (4)$$

$$\cup \{(u_{h_j[i]}^j, u_k^{j'}, c) \mid 1 \leq j, j' \leq |\mathcal{H}|, \exists (u_{h_j[i]}, u_k) \in E, \text{ s.t. } i = |h_j| \text{ or } h_j[i+1] \neq u_k\} \quad (5)$$

Next, we give some intuition for each line (2)-(5) in the above construction of  $E_E$ .

- (2) Weight 0 edges are created from  $s$  to the start of each haplotype path in  $G_E$ .
- (3) Weight 0 edges are created from the end of each haplotype path in  $G_E$  to  $t$ .
- (4) Weight 0 edges are created between adjacent vertices in each haplotype path. That is, in the path for  $h_j$ , an edge is created from  $u_{h_j[i]}$  to  $u_{h_j[i+1]}$ .
- (5) Weight  $c$  edges are used to represent the useful recombinations described in Lemma 1. We call these *recombination edges*.

We use  $\epsilon$  to denote the empty string. The vertex labels are defined as follows:

$$\sigma_E(u_{h_j[i]}^j) = \sigma(u_{h_j[i]}) \text{ for } 1 \leq j \leq |\mathcal{H}|, 1 \leq i \leq |h_j| \quad (6)$$

$$\sigma_E(s) = \sigma_E(t) = \epsilon \quad (7)$$

- (6) The vertices in a haplotype path are labeled according to the corresponding vertex label in  $G$ . These labels will be used to identify matches.
- (7) The source, sink, do not require vertex labels and are hence labeled with the empty string  $\epsilon$ .

*Optimizing the Expanded Graph.* One issue with the above construction is that the number of recombination edges for a given potential recombination can be  $O(|\mathcal{H}|^2)$  in the worst case. This occurs because we maintain  $|haps(v)|$  copies of each vertex  $v \in V$ . For every edge  $(u, v) \in E$  allowing a recombination, we add  $O(|haps(u)| \cdot |haps(v)|)$  edges to the edge set  $E_E$ . Since both  $|haps(u)|$  and  $|haps(v)|$  can be at most  $|\mathcal{H}|$ , any potential recombination can result in  $O(|\mathcal{H}|^2)$  recombination edges in the worst case. We observe this issue in practice as well. An improvement is to represent a recombination by having an intermediate vertex  $w_e$  that represents the edge  $e \in E$  allowing for the recombination. We then create an edge to  $w_e$  from every vertex in a haplotype path which the recombination would start from, and edges from  $w_e$  to every vertex in a haplotype path to

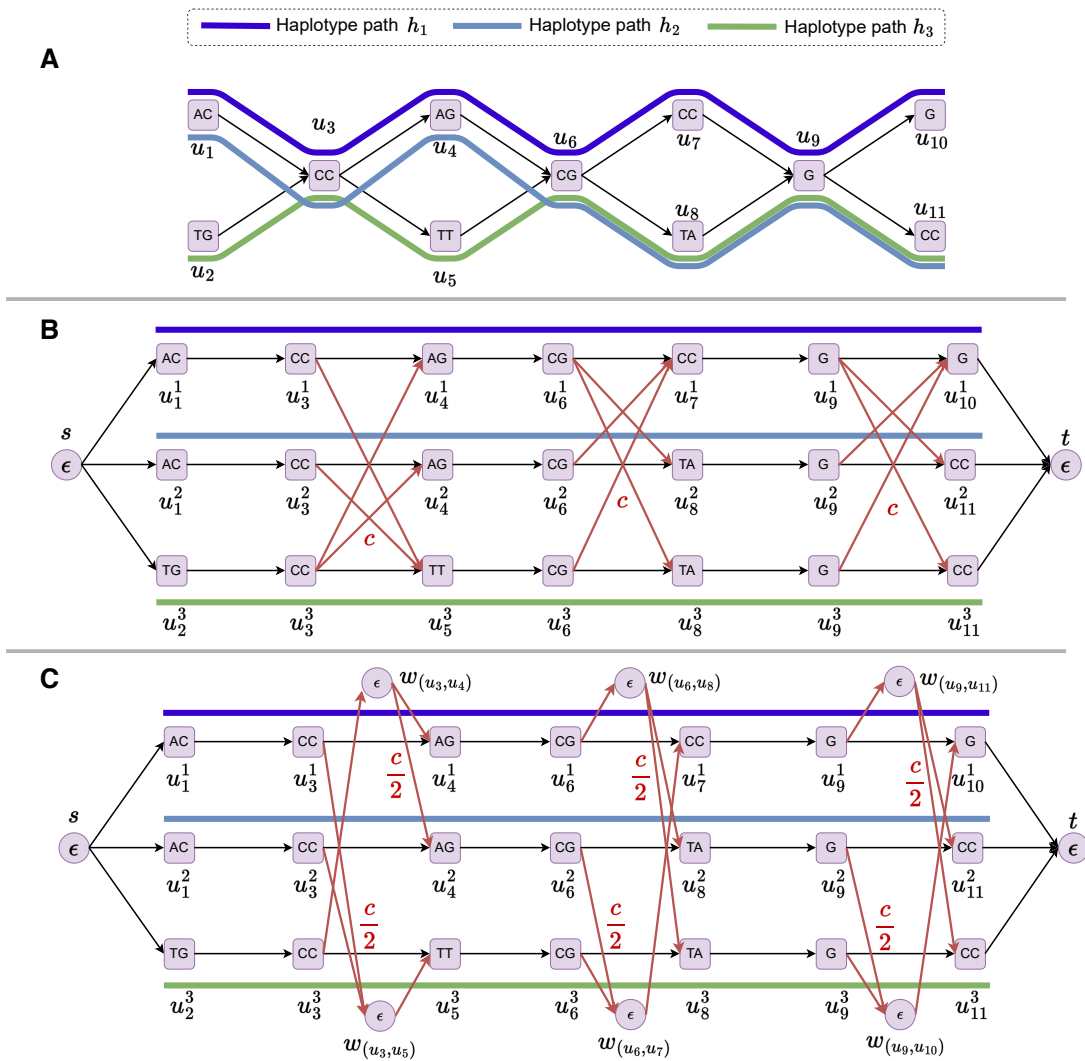


Fig. 2: **Illustration of an expanded graph.** (A) A pangenome graph with three haplotype paths  $h_1$ ,  $h_2$  and  $h_3$ . (B) The corresponding expanded graph which includes three disjoint paths, one for each haplotype path. The recombination edges are shown in purple, these edges have a weight of  $c$ . We consider only the useful recombinations (Lemma 1). The edges which are not recombination edges in the expanded graph have a weight of 0. (C) The corresponding optimized expanded graph.

which the recombination would lead to (Figure 2C). More formally, the modified vertex set becomes

$$V_E = \{s\} \cup \{t\} \cup \{u_{h_j[i]}^j \mid 1 \leq j \leq |\mathcal{H}|, 1 \leq i \leq |h_j|\} \cup \{w_e \mid e \in E\} \quad (8)$$

We also replace Line (5) in the construction of  $E_E$  with the Lines (9) and (10) as follows:  $E_E = \dots$

$$\cup \{(u_{h_j[i]}^j, w_e, c/2) \mid 1 \leq j \leq |\mathcal{H}|, \exists e = (u_{h_j[i]}, u_k) \in E, \text{ s.t. } i = |h_j| \text{ or } h_j[i+1] \neq u_k\} \quad (9)$$

$$\cup \{(w_e, u_k^j, c/2) \mid 1 \leq j \leq |\mathcal{H}|, \exists e = (u_k, u_{h_j[i]}) \in E, \text{ s.t. } i = 1 \text{ or } h_j[i-1] \neq u_k\} \quad (10)$$

We now call these edges created in Lines (9) and (10) the recombination edges. After creating the edges in  $E_E$ , we delete any  $w_e$  vertex that is isolated in  $G_E$ . Finally, for any remaining  $w_e$  vertices, we define  $\sigma_E(w_e) = \epsilon$ . Observe, that the above modification allows for the same set of useful recombinations as our initial expanded graph construction. However, per potential useful recombination, the number of edges remains  $O(|\mathcal{H}|)$  rather than  $O(|\mathcal{H}|^2)$ . Before giving the integer programming solutions, we require one additional definition.

**Definition 1 (Hits).** For a string  $r \in \mathcal{S}$ , assuming  $\max_{u \in V_E} |\sigma_E(u)| < |r|$ , a path in  $G_E$ , denoted as an ordered edge set  $((u, v), (v, w), (w, x) \dots, (y, z))$ , matches  $r$  if  $r = \sigma_E(u)' \circ \sigma_E(v) \circ \sigma_E(w) \circ \sigma_E(x) \circ \sigma_E(y) \circ \sigma_E(z)'$ , where  $\sigma_E(u)'$  is a suffix of  $\sigma_E(u)$  and  $\sigma_E(z)'$  a prefix of  $\sigma_E(z)$ . We use  $hits(r)$  to represent the set of paths matching string  $r$  in  $G_E$ .

## Integer Linear Programming (ILP) Solution

We assume that the maximum length of any vertex label is upper bounded by the length of any string in  $\mathcal{S}$ , i.e.,  $\max_{u \in V_E} |\sigma_E(u)| < \min_{r \in \mathcal{S}} |r|$ . This condition can be easily enforced in the input graph by adjusting the lengths of vertex labels, e.g., by splitting a vertex with a long label into two, while ensuring that the graph's topology is preserved. We assume  $\min_{r \in \mathcal{S}} |r| > 1$ .

The basis for our solution is to find an  $st$ -flow with a flow of 1 through the expanded graph  $G_E$ . Our integer programs will utilize binary decision variable  $x_{uv}$  for each edge. The variable  $x_{uv}$  will take the value 1 if edge  $(u, v) \in E_E$  is part of the solution flow and 0 otherwise. Because these are binary variables, the flow will always be a path. From the solution path in  $G_E$ , it is straight forward to recover the corresponding inferred path  $\mathcal{P}$ . We use binary decision variable  $z_r$  for each string  $r \in \mathcal{S}$  such that  $z_r$  will take the value 1 if the solution flow includes a subpath from  $hits(r)$ . We also use variable  $z_{r\omega}$  for each  $\omega \in hits(r), r \in \mathcal{S}$ .

Letting  $weight(u, v)$  denote the weight of an edge  $(u, v) \in E_E$ , our ILP formulation is as follows:

$$\min \sum_{(u,v) \in E_E} weight(u, v) \cdot x_{uv} + \sum_{r \in \mathcal{S}} (1 - z_r), \quad (11)$$

subject to

$$\sum_{v \in \mathcal{N}^+(u)} x_{uv} - \sum_{v \in \mathcal{N}^-(u)} x_{vu} = \begin{cases} 1 & \text{if } u = s, \\ -1 & \text{if } u = t, \\ 0 & \text{otherwise,} \end{cases} \quad \forall u \in V_E, \quad (12)$$

$$\sum_{(u,v) \in \omega} x_{uv} \geq |\omega| \cdot z_{r\omega}, \quad z_{r\omega} \in \{0, 1\}, \quad \forall \omega \in hits(r), \forall r \in \mathcal{S}, \quad (13)$$

$$\sum_{\omega \in hits(r)} z_{r\omega} = z_r, \quad z_r \in \{0, 1\}, \quad \forall r \in \mathcal{S}, \quad (14)$$

$$x_{uv} \in \{0, 1\}, \quad \forall (u, v) \in E_E \quad (15)$$

In the ILP formulation, the Objective (11) models  $Cost(\mathcal{P})$ . The summation over  $weight(u, v) \cdot x_{uv}$  imposes penalty  $c$  for each recombination. This is due to the two  $c/2$  weighted recombination edges that must traversed when the path switches between haplotype paths in  $G_E$  (Figure 2C). In the second summation, the term  $(1 - z_r)$  adds a penalty of 1 to the objective for every  $r \in \mathcal{S}$  where  $\bar{\chi}(r, \sigma(\mathcal{P})) = 1$ . Constraint (12) enforces flow conservation, allowing a unit flow from the source vertex  $s$  to the sink vertex  $t$ , ensuring that the ILP formulation selects a single path in the expanded graph.

To explain the function of Constraint (13), termed as linear string-hit constraint and (14), observe that in an optimal solution, whenever possible the variable  $z_r$  is set to 1. This is because the term  $(1 - z_r)$  in the objective function adds a penalty of 0 whenever  $z_r = 1$ . However, this is only possible when  $z_{r\omega}$  is equal to 1 for some  $\omega \in hits(r)$ . This, in turn, is only possible if  $\sum_{(u,v) \in \omega} x_{uv} = |\omega|$ , meaning  $r$  occurs as a substring in the inferred path. Also note that at most one  $z_{r\omega}$  variable can equal 1 in Constraint (14). Other  $z_{r\omega'}$  variables, where  $\omega, \omega' \in hits(r)$  and  $\omega \neq \omega'$ , can have a value of 0, even if  $\sum_{(u,v) \in \omega'} x_{uv} = |\omega'|$ , justifying the use of equality in Constraint (14).

A weakness of the proposed ILP formulation is that the number of string-hit constraints equals the total number of string matches, that is,  $\sum_{r \in \mathcal{S}} hits(r)$ . We design another formulation with quadratic constraints in which fewer constraints are needed.

## Integer Quadratic Programming (IQP) Solution

In our IQP formulation, Objective (11), and Constraints (12), and (14) and (15) remain unchanged from the ILP formulation. Constraints in (13) are replaced by quadratic constraints defined as

$$\sum_{\omega \in hits(r)} \left( 1 - |\omega| + \sum_{(u,v) \in \omega} x_{uv} \right) \cdot z_{r\omega} = z_r, \quad \forall r \in \mathcal{S}, \quad (16)$$

We call Constraint (16) the quadratic string-hit constraint. Again, due to Constraint (14) at most one  $z_{r\omega}$  variable can be 1. The expression  $1 - |\omega| + \sum_{(u,v) \in \omega} x_{uv}$  sums to 1 when the subpath  $\omega$  is contained in the flow. In this case  $z_r$  will take the value 1 and no penalty is paid in the objective. Conversely, if some of the edges for  $\omega$  are not in the flow, the expression will sum to  $\leq 0$ . If this is the case for each  $\omega \in hits(r)$ , then Constraint (16) can only be satisfied by setting  $z_r = 0$  and  $z_{r\omega} = 0$  for each  $\omega \in hits(r)$ . Since  $z_r = 0$ , a penalty is paid in the objective. The total number of quadratic string-hit constraints is  $|\mathcal{S}|$ . In our experiments, we observe that IQP formulation solves the problem faster, albeit while requiring more memory.

As a further improvement, we relax the variables  $x_{uv}$  for all  $(u, v) \in E_E$  to continuous values  $x_{uv} \in [0, 1]$  in Constraint (15), following Lemma 2.

**Lemma 2.** *An optimal solution  $\phi_{cont}$  to the IQP (or ILP) with relaxed Constraint (15) where variables  $x_{uv}$  lie within the continuous interval  $[0, 1]$  can be transformed in polynomial time to an optimal solution  $\phi$  satisfying  $x_{uv} \in \{0, 1\}$  for all  $(u, v) \in E_E$ .*

*Proof.* First, observe that  $z_r = 1$  if and only if all edges in some  $\omega \in hits(r)$  have their corresponding variables set to 1. This follows from Constraints (13) and (16), and the fact that at most one  $z_{r\omega}$  can be 1 for a given  $r$ , by Constraint (14).

If  $z_r = 0$  for all  $r \in \mathcal{S}$  in  $\phi_{cont}$ , then  $\phi$  can be trivially obtained as a single haplotype path in  $G_E$  without recombination penalties. In such a case, all edge variables are assigned either 0 or 1.

For the remaining cases, we introduce the following terms:

- $\omega \in hits(r)$  is a *used hit-subpath* if  $z_{r\omega} = 1$ .
- A flow between vertices  $u$  and  $v$  can be decomposed into  $uv$ -paths each assigned some positive flow and called *flow subpaths*.
- $\omega$  is the *first used hit-subpath* if there is a flow subpath from vertex  $s$  to the first vertex of  $\omega$  without passing through another used hit-subpath.
- $\omega$  is the *last used hit-subpath* if there is a flow subpath from the last vertex of  $\omega$  to vertex  $t$  without passing through another used hit-subpath.

- $\omega$  and  $\omega'$  are *consecutive used hit-subpaths* if there is a flow subpath between them without passing through a third used hit-subpath, where  $\omega' \neq \omega$  and  $\omega' \in \text{hits}(r)$ .

Now, if  $z_r = 1$  in  $\phi_{cont}$  for some  $r \in \mathcal{S}$ , there exists a used hit-subpath. We obtain  $\phi$  as following. The flow used to reach the first hit-subpath avoids recombination penalties by following a single haplotype path. Similarly, the flow from the end vertex on the last used hit-subpath to  $t$  avoids recombinations penalties by staying on a single haplotype path. Next, consider two consecutive used hit-subpaths  $\omega$  and  $\omega'$ , with  $u$  and  $v$  as their respective end and start vertices. If  $u$  and  $v$  are on different haplotype paths, any flow subpaths between  $u$  and  $v$  must minimize the recombination penalty. The same minimum recombination cost can be achieved by replacing the potentially multiple fractional flow subpaths with a single path that incurs the same recombination penalty. We can select any flow subpath from  $u$  to  $v$  and assign its edge variables to 1. Edge variables on edges used on the flow from  $u$  to  $v$  and not on this selected path are set to 0.  $\square$

## Results

**Implementation Details.** We implemented our ILP and IQP solutions in C++ using Gurobi (v11.0.2) solver. We call our software as PHI. The user can provide a pangenome reference as either a graph (GFA format) or as a list of phased variants (VCF format). Given short-read or long-read sequencing data of either a haploid or a homozygous genome, PHI outputs the haplotype sequence associated with the optimal inferred path from the graph in FASTA format.

Given a set of reads, we compute  $(w, k)$  window minimizers (Roberts et al., 2004) for identifying our *hits* (Definition 1). By default, we use  $w = 25$  and  $k = 31$ . This choice of  $k$  strikes a balance, as longer  $k$ -mers have a higher likelihood of including sequencing errors, while shorter  $k$ -mers can increase the number of false-positive matches. The set of window minimizers corresponds to the input set of strings  $\mathcal{S}$  in Problem 1. In practice, most minimizers will have a match with two or more haplotypes. We expect that a small fraction of the minimizers will match to a single haplotype due to haplotype-specific variation. Such minimizers will help steer the optimization toward identifying the correct haplotype paths.

Computing minimizer matches between two strings is faster than computing minimizer matches on a pangenome graph. For this reason, we find minimizer matches between reads and the sequences spelled by all the haplotype paths in the graph. This means  $\text{hits}(r)$  includes only those subpaths that are completely contained in some haplotype path in  $G_E$  (Definition 1). This restriction to  $\text{hits}(r)$  also prevents us from needing to perform the additional edge splitting step described in Methods. In PHI, we set the default value

of recombination penalty parameter  $c$  as 100. We discuss the effect of this parameter later. We ran all our experiments on AMD EPYC 7763 processors with 512 GB RAM. We used 32 threads in all experiments.

**Datasets.** We evaluated our algorithm by estimating MHC sequences of five haplotypes (APD, DBB, MANN, QBL, SSTO) from homozygous human cell lines. Recently, Houwaart *et al.* (Houwaart *et al.*, 2023; Scholz, 2024) published complete assemblies of these MHC sequences using long and short-read sequencing. The average length of these assemblies is 4.99 Mbp. We downloaded the five short-read sequencing datasets available from this study. To evaluate our algorithm using varying sequencing coverage, we down-sampled each short-read dataset to obtain coverage of  $0.1\times$ ,  $0.5\times$ ,  $1\times$ ,  $2\times$ ,  $5\times$ , and  $10\times$ . We also used the full datasets for evaluation (coverage  $12.9 - 18.2\times$ ). We used the complete assemblies of five MHC haplotypes as ground-truth to evaluate the accuracy of our estimated sequences. To quantify the accuracy, we measured edit distance between each estimated sequence and the corresponding ground-truth sequence.

We built a haplotype-resolved pangenome graph of 49 complete MHC sequences (Li, 2022) using Minigraph-Cactus (Hickey *et al.*, 2023). These sequences were extracted from phased assemblies of 24 diploid human samples (Liao *et al.*, 2023) and the CHM13 reference (Nurk *et al.*, 2022). Using Minigraph-Cactus, we obtained the pangenome reference in a VCF format file. We subjected this file to further simplification steps (<https://github.com/eblerjana/genotyping-pipelines/tree/main/prepare-vcf-MC>) to ensure compatibility with various tools. We show sequence similarity statistics between the complete MHC assemblies of five haplotypes (APD, DBB, MANN, QBL, SSTO) and the 49 pangenome reference haplotypes in Supplemental Table S1.

**Other Methods.** We compared PHI with two existing pangenome-based genotyping tools - (i) VG (v1.60) (Sirén *et al.*, 2024) and (ii) PanGenie (v3.1) (Ebler *et al.*, 2022). VG supports sampling of relevant haplotypes from a pangenome graph by comparing  $k$ -mer counts in the reads and  $k$ -mers of a reference haplotype. The selection of haplotypes is done locally in fixed-length non-overlapping blocks. Recombinations may be introduced to create contiguous haplotypes across the blocks. The number of samples can be specified by the user. Accordingly, VG's haplotype sampling feature can be adapted for haplotype sequence estimation by simply setting the number of desired samples to one. Next, PanGenie supports short-read genotyping using a haplotype-resolved pangenome graph. PanGenie uses a hidden Markov model, which is similar to the standard Li and Stephens model (Li and Stephens, 2003). PanGenie compares  $k$ -mer counts in the reads with the  $k$ -mers present in the graph to compute genotype likelihoods. PanGenie exhibited better genotyping accuracy and speed than other genotyping tools (Ebler *et al.*, 2022). Our sequencing datasets are

derived from homozygous cell lines, therefore we ignored the heterozygous genotype calls made by PanGenie (Supplemental Table S3). We incorporated PanGenie’s predicted genotypes in the reference sequence to obtain the haplotype sequence. We list our commands to run PHI, VG and PanGenie in Supplemental Table S2.

**Genotyping performance.** We evaluated PHI, VG and PanGenie methods in their ability to infer the MHC sequences from short read datasets of varying coverage (see Figure 3A-E). Using low coverage datasets ( $0.1 - 2\times$ ), PHI exhibits significantly higher accuracy. VG and PanGenie methods may not be suitable for low-coverage sequencing. For example, the distribution of  $k$ -mer counts at low coverage can be unreliable. Distinguishing  $k$ -mers originating from unique versus repetitive regions, as required by PanGenie and VG, is also challenging at low-coverage. Using coverage of  $5\times$  or more, the results of VG and PHI are comparable. PanGenie also produces comparable results using full datasets. We note that the integer programming (IQP) approach used in PHI requires more time and memory compared to the methods used in VG and PanGenie. PHI used up to 1.5 hours and 137 GB RAM in a single experiment. In contrast, VG and PanGenie required  $< 5$  minutes and  $< 50$  GB memory. It may be possible to optimize PHI by incorporating efficient heuristics. We show detailed performance statistics for PHI, including its runtime and memory usage in Supplemental Table S4.

**Effect of parameter  $c$ .** We further explored the effect of parameter  $c$  in PHI. The user-specified parameter  $c$  determines the cost per recombination event. Having a good balance between allowing too many recombination events versus not allowing recombinations is important. By default, PHI uses  $c = 100$ . We arrived at this default value by testing three different settings:  $c = 10$ ,  $c = 100$ , and  $c = 1000$ . For this experiment, we used the five full short-read datasets as input to PHI. We present the edit distances, counts of recombination events, and the lengths of the estimated sequences in Table 1. We observe that using  $c = 100$  results in the lowest edit distances consistently across all five haplotypes. In the following, we argue why having  $c = 10$  or  $c = 1000$  results in less accurate outputs.

A large value of  $c$  such as 1000 can be interpreted as significantly limiting the flexibility to recombine, thus forcing the optimal output to use fewer recombinations than required. Accordingly, we find that the count of recombinations using  $c = 1000$  reduces using all five datasets (Table 1). A large value of  $c$  is appropriate in those scenarios where the target haplotype is known to closely resemble a haplotype included in the pangenome graph. Otherwise, it will result in low output accuracy as we see in this experiment. On the other hand, we find that a very low value of  $c$  such as 10 allows an excessive number of recombination

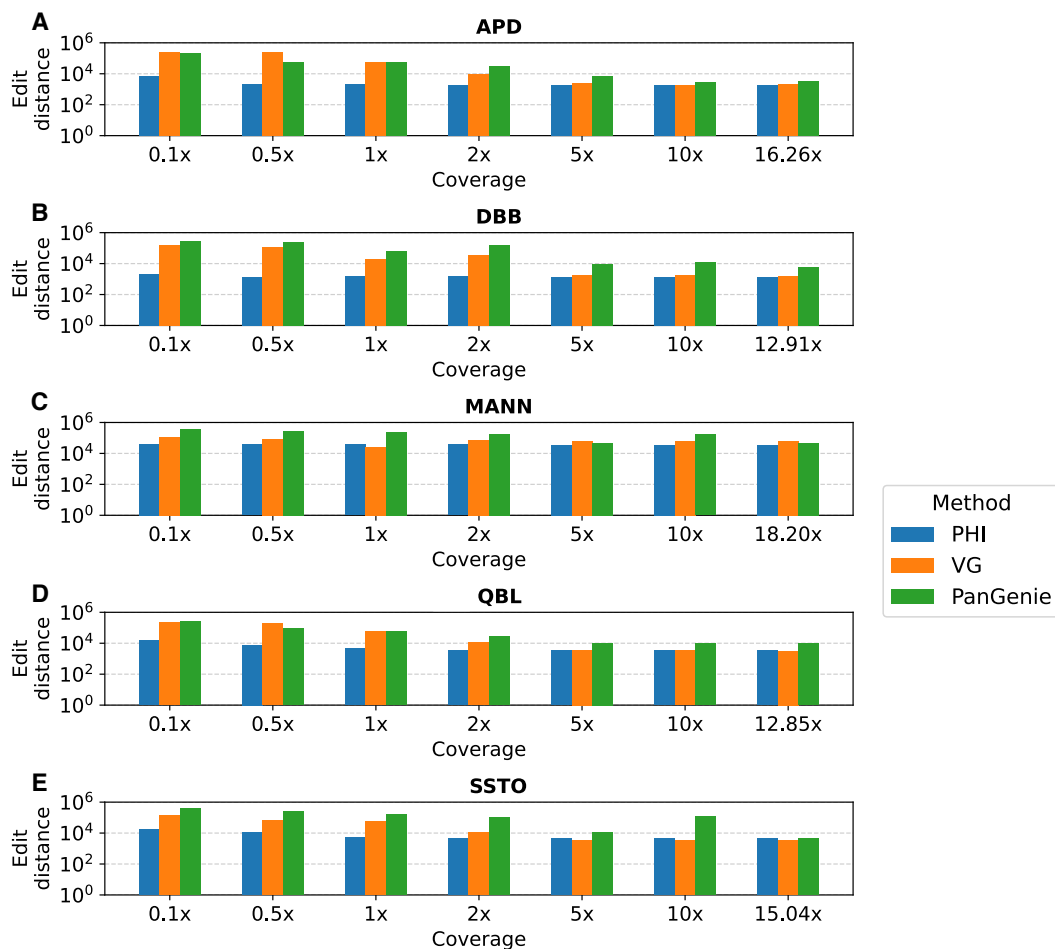


Fig. 3: **Accuracy of the haplotype sequences estimated by PHI, VG, and PanGenie.** We used short-read datasets from MHC sequences of five haplotypes (APD, DBB, MANN, QBL, SSTO). The x-axes indicate the coverage of short-read data. The y-axes indicate the edit distance between the estimate haplotype sequence and the ground-truth sequence on a logarithmic scale.

events. In such situations, the algorithm greedily chooses paths in the graph to maximize minimizer matches. This greedy choice affects the accuracy when there are false positive minimizer matches due to repeats and sequencing errors. Due to this reason, we find that having  $c = 10$  results in longer output sequences than expected. Overall, this experiment also justifies the design of our problem formulation and the importance of maintaining a balance between the count of minimizer matches and the count of recombination events.

Table 1: **Effect of modifying parameter  $c$  on PHI’s output.** The table presents output statistics obtained by running PHI using three different choices of the recombination penalty parameter  $c$ : 10, 100, and 1000. We calculated (i) accuracy, that is, edit distance values between the estimated haplotype sequences and the ground truth assemblies, (ii) count of recombinations in the estimated haplotype sequences, and (iii) lengths of the estimated sequences. The last column shows the length of the ground truth sequences.

Haplotype	Edit distance			No. of recombinations			Haplotype length (Mbp)			
	$c = 10$	$c = 10^2$	$c = 10^3$	$c = 10$	$c = 10^2$	$c = 10^3$	Estimated sequence by PHI		Ground-truth	
							$c = 10$	$c = 10^2$	$c = 10^3$	
APD	144835	1810	9994	36	10	4	5.07	4.93	4.92	4.93
DBB	67290	1377	2111	13	4	2	5.11	5.05	5.05	5.05
MANN	131564	35940	38872	37	14	5	5.15	5.04	5.04	5.03
QBL	125655	3343	13034	46	17	4	5.03	4.90	4.90	4.90
SSTO	125770	4637	13600	68	24	5	5.17	5.04	5.05	5.05

**Effect of our optimizations.** In PHI, we implemented both ILP-based and IQP-based solutions to solve the optimization problem. Using either solution, Gurobi solves Problem 1 to optimality. We benchmarked our ILP and IQP solutions to compare their runtime and memory-usage (see Figure 4A-D). On low-coverage datasets ( $0.1-1\times$ ), the runtimes are comparable. At higher coverage, the IQP solution runs faster, which is likely due to fewer string-hit constraints used (Methods). Although, it requires approximately 1.5 times more memory. This may be because Gurobi requires additional storage to handle quadratic constraints. Accordingly, while using PHI, the user can choose between ILP and IQP using a command line argument based on the available memory. If no choice is provided, the IQP solution is used by default. We also evaluated the advantage of relaxing edge variables to continuous values (Lemma 2) by comparing it to another version of our code where we set the edge variables to be discrete. Relaxation of variables decreases runtime of the IQP solution by a factor of 1.6 on average (Supplemental Figure S2). Not much effect on the runtime is observed in the ILP solution (Supplemental Figure S3).

**Impact of graph expansion with the addition of more genomes.** We evaluated the impact of pangenome graph expansion on PHI’s genotyping accuracy as well as runtime. To do this, we created five

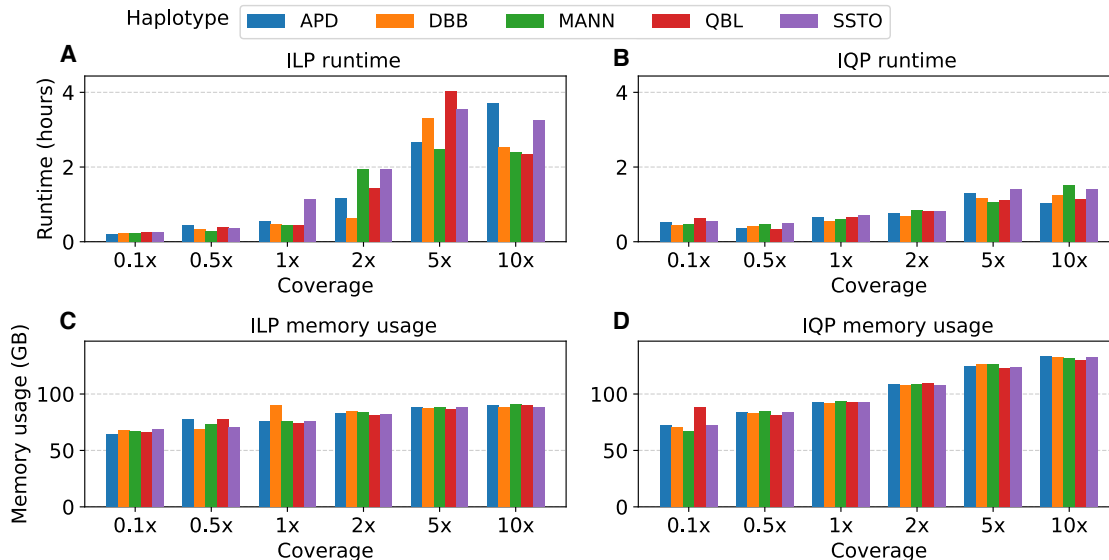


Fig. 4: **Performance comparison between the ILP and IQP solutions implemented in PHI.** We compared their runtime and memory-usage using short-read sequencing datasets sampled from five haplotypes.

versions of our pangenome graph, each containing an increasing number of reference haplotypes, added progressively. The first graph comprises a single diploid sample (chosen randomly from 24 diploid samples) plus CHM13 reference, therefore, it has three reference haplotypes in total. The second graph includes two more diploid samples (chosen randomly from the remaining 23), therefore, it has seven reference haplotypes in total. Similarly, third, fourth and fifth graphs contain 13, 25 and 49 reference haplotypes, respectively. The fifth graph is equivalent to the graph used in previous experiments as well. This results in five different graphs that have 3, 7, 13, 25, and 49 reference haplotypes respectively.

We repeated our experiments with full short-read datasets using these five graphs and present results in Figure 5. We observe that edit distances between the estimated sequences and the ground truth sequences decrease with the increasing number of reference haplotypes. This is expected because more haplotypes are available to choose from when we compute our inferred path in the graph. We also observe an increase in runtime and memory usage. Runtime appears to increase superlinearly and memory appears to increase linearly with the number of reference haplotypes. This is because the size of expanded graph and the number of minimizer matches increase leading to more variables and constraints in our integer program.

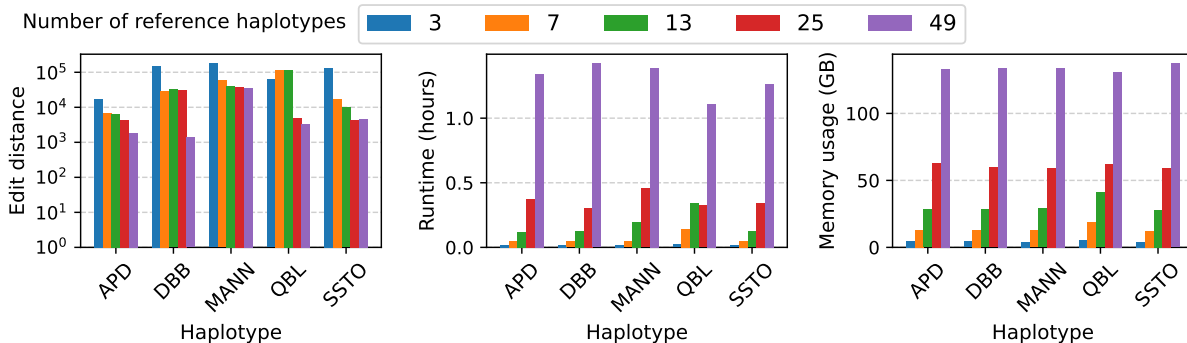


Fig. 5: **Assessment of PHI's performance with the increasing number of genomes in pangenome graph.** The left figure shows the accuracy in terms of edit distance between the output sequences and ground-truth sequences. The middle and right figure show the runtime and memory-usage respectively.

## Discussion

Genotyping using pangenome graphs is equivalent to finding a walk in the graph that contains the sample's variants (Paten et al., 2017). If the sample is diploid, this becomes equivalent to finding a pair of paths. Drawing inspiration from this idea, we proposed a rigorous framework to infer a path through the graph, such that the sequence spelled by the path is consistent with the sequencing data in terms of the shared  $k$ -mers between them, while permitting a limited number of recombinations in the path, each incurring a fixed penalty. This optimization problem requires considering all possible paths in the graph. We proved that this problem is NP-Hard and subsequently gave efficient integer programming solutions. As part of our methodology, we introduced the expanded graph data structure on which we could compute an appropriate  $st$ -flow of 1. Experimental results demonstrate the advantage of the proposed ILP/IQP approaches for accurate genome inference, especially with low-coverage data (coverage  $0.1 - 1\times$ ). Thus, our algorithm can facilitate affordable genotyping and association studies of complex and repeat-rich regions of the genome.

Although our approach is currently tailored to haploid samples, it could generalize to diploid samples. This may be accomplished by finding an  $st$ -flow of 2 through the expanded graph and modifying some constraints. How well this approach genotypes and phases the genome would be interesting to explore. Another limitation of this work is that we do not capture uncertainty. For example, there may be multiple inferred paths with minimum cost. Other future directions to explore include obtaining a better runtime bound for the bounded recombination case. This could be achieved by modifying Problem 1 to have an upper limit on the count of recombinations, effectively limiting the number of recombinations in the inferred haplotype. This problem version is yet to be explored.

Our problem formulation employs a simple cost model for penalizing recombination events by applying a fixed penalty for each event. In PHI, we set this parameter to 100 by default based on our empirical observations on MHC pangenome graph. But on other pangenome graphs, users may need to tune the parameter using the procedure described in Results (Subsection: Effect of parameter  $c$ ).

Further refinement of the proposed optimization framework is necessary to align it with established statistical models in population genetics and evolutionary biology. [Rosen et al. \(2017\)](#) developed an efficient algorithm for estimating haplotype likelihoods on graphs, extending the Li-Stephens model to graph-based representations. Lastly, pangenome graphs are expected to grow in the number of genomes, therefore, scaling the current approach to a large number of haplotype paths may be important. We leave these extensions to future work.

## Software availability

The software PHI is available at GitHub (<https://github.com/at-cg/PHI>) and as Supplemental Code. The scripts to reproduce the results are available at GitHub (<https://github.com/at-cg/PHI/tree/master/data>) and as Supplemental Scripts.

## Competing interest statement

The authors report no competing interests.

## Acknowledgements

This research is funded in part by the DBT/Wellcome Trust India Alliance Fellowship (grant number IA/I/23/2/506979), the Intel India Research Fellowship, the National Institutes of Health of the USA (NIH-NIAID U01 AI090905), and the Jürgen Manchot Foundation. We utilized computing resources available at the Indian Institute of Science and the U.S. National Energy Research Scientific Computing Center.

## Author contributions

All authors contributed to the development of the methods and proofs. G.C. implemented the software and conducted the experiments.

## References

- Baaijens JA, Bonizzoni P, Boucher C, Della Vedova G, Pirola Y, Rizzi R, and Sirén J. 2022. Computational graph pangenomics: a tutorial on data structures and their applications. *Natural Computing* pp. 1–28.
- Bradbury PJ, Casstevens T, Jensen SE, Johnson L, Miller Z, Monier B, Romay M, Song B, and Buckler ES. 2022. The practical haplotype graph, a platform for storing and using pangenomes for imputation. *Bioinformatics* **38**: 3698–3702.
- Chandra G, Gibney D, and Jain C. 2024. Haplotype-aware sequence alignment to pangenome graphs. *Genome Research* **34**: 1265–1275.
- Computational Pan-Genomics Consortium. 2018. Computational pan-genomics: status, promises and challenges. *Briefings in bioinformatics* **19**: 118–135.
- Davies RW, Kucka M, Su D, et al.. 2021. Rapid genotype imputation from sequence with reference panels. *Nature Genetics* **53**: 1104–1111.
- Dilthey A, Cox C, Iqbal Z, Nelson MR, and McVean G. 2015. Improved genome inference in the MHC using a population reference graph. *Nature genetics* **47**: 682–688.
- Dilthey AT. 2021. State-of-the-art genome inference in the human MHC. *The International Journal of Biochemistry & Cell Biology* **131**: 105882.
- Ebert P, Audano PA, et al.. 2021. Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**.
- Ebler J, Ebert P, Clarke WE, Rausch T, Audano PA, Houwaart T, Mao Y, Korbel JO, Eichler EE, Zody MC, et al.. 2022. Pangenome-based genome inference allows efficient and accurate genotyping across a wide spectrum of variant classes. *Nature genetics* **54**: 518–525.
- Eggertsson HP, Jonsson H, Kristmundsdottir S, et al.. 2017. Graphtyper enables population-scale genotyping using pangenome graphs. *Nature genetics* **49**: 1654–1660.
- Gao Y, Yang X, Chen H, Tan X, Yang Z, Deng L, Wang B, Kong S, Li S, Cui Y, et al.. 2023. A pangenome reference of 36 chinese populations. *Nature* **619**: 112–121.
- Grytten I, Dagestad Rand K, and Sandve GK. 2022. Kage: Fast alignment-free graph-based genotyping of SNPs and short indels. *Genome Biology* **23**: 209.
- Harris L, McDonagh EM, Zhang X, Fawcett K, Foreman A, Daneck P, Sergouniotis PI, Parkinson H, Mazzarotto F, Inouye M, et al.. 2024. Genome-wide association testing beyond SNPs. *Nature Reviews Genetics* pp. 1–15.

- Hickey G, Heller D, Monlong J, Sibbesen JA, Sirén J, Eizenga J, Dawson ET, Garrison E, Novak AM, and Paten B. 2020. Genotyping structural variants in pangenome graphs using the vg toolkit. *Genome biology* **21**: 1–17.
- Hickey G, Monlong J, Ebler J, Novak AM, Eizenga JM, Gao Y, Marschall T, Li H, and Paten B. 2023. Pangenome graph construction from genome alignments with minigraph-cactus. *Nature Biotechnology* pp. 1–11.
- Houwaart T, Scholz S, Pollock NR, et al.. 2023. Complete sequences of six major histocompatibility complex haplotypes, including all the major MHC class II structures. *HLA* **102**: 28–43.
- Jain C, Tavakoli N, and Aluru S. 2021. A variant selection framework for genome graphs. *Bioinformatics* **37**: i460–i467.
- Letcher B, Hunt M, and Iqbal Z. 2021. Gramtools enables multiscale variation analysis with genome graphs. *Genome biology* **22**: 1–27.
- Li H. 2022. Sample graphs and sequences for testing sequence-to-graph alignment.
- Li JH, Mazur CA, Berisa T, and Pickrell JK. 2021. Low-pass sequencing increases the power of GWAS and decreases measurement error of polygenic risk scores compared to genotyping arrays. *Genome research* **31**: 529–537.
- Li N and Stephens M. 2003. Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. *Genetics* **165**: 2213–2233.
- Liao WW, Asri M, Ebler J, Doerr D, Haukness M, Hickey G, Lu S, Lucas JK, Monlong J, Abel HJ, et al.. 2023. A draft human pangenome reference. *Nature* **617**: 312–324.
- Mahmoud M, Gobet N, Cruz-Dávalos DI, Mounier N, Dessimoz C, and Sedlazeck FJ. 2019. Structural variant calling: the long and the short of it. *Genome biology* **20**: 1–14.
- Martin AR, Atkinson EG, Chapman SB, Stevenson A, Stroud RE, Abebe T, Akena D, Alemayehu M, Ashaba FK, Atwoli L, et al.. 2021. Low-coverage sequencing cost-effectively detects known and novel variation in underrepresented populations. *The American Journal of Human Genetics* **108**: 656–668.
- Mun T, Vaddadi NSK, and Langmead B. 2023. Pangenomic genotyping with the marker array. *Algorithms for Molecular Biology* **18**: 2.
- Mustafa H, Karasikov M, Mansouri Ghiasi N, Rättsch G, and Kahles A. 2024. Label-guided seed-chain-extend alignment on annotated de bruijn graphs. *Bioinformatics* **40**: i337–i346.
- Nurk S, Koren S, Rhie A, Rautiainen M, et al.. 2022. The complete sequence of a human genome. *Science* **376**: 44–53.

- Paten B, Novak AM, Eizenga JM, and Garrison E. 2017. Genome graphs and the evolution of genome inference. *Genome research* **27**: 665–676.
- Petes TD. 2001. Meiotic recombination hot spots and cold spots. *Nature Reviews Genetics* **2**: 360–369.
- Pritt J, Chen NC, and Langmead B. 2018. Forge: prioritizing variants for graph genomes. *Genome biology* **19**: 1–16.
- Roberts M, Hayes W, Hunt BR, Mount SM, and Yorke JA. 2004. Reducing storage requirements for biological sequence comparison. *Bioinformatics* **20**: 3363–3369.
- Rosen Y, Eizenga J, and Paten B. 2017. Modelling haplotypes with respect to reference cohort variation graphs. *Bioinformatics* **33**: i118–i123.
- Rubinacci S, Ribeiro DM, et al.. 2021. Efficient phasing and imputation of low-coverage sequencing data using large reference panels. *Nature Genetics* **53**: 120–126.
- Scholz S. 2024. Complete sequences of six major histocompatibility complex haplotypes rev2.
- Sibbesen JA, Maretty L, Consortium DPG, and Krogh A. 2018. Accurate genotyping across variant classes and lengths using variant graphs. *Nature genetics* **50**: 1054–1059.
- Sirén J, Monlong J, Chang X, Novak AM, Eizenga JM, Markello C, Sibbesen JA, Hickey G, Chang PC, Carroll A, et al.. 2021. Pangenomics enables genotyping of known structural variants in 5202 diverse genomes. *Science* **374**: abg8871.
- Sirén J, Eskandar P, Ungaro MT, et al.. 2024. Personalized pangenome references. *Nature Methods* .
- Smith TP, Bickhart DM, Boichard D, Chamberlain AJ, Djikeng A, Jiang Y, Low WY, Pausch H, Demyda-Peyrás S, Prendergast J, et al.. 2023. The bovine pangenome consortium: democratizing production and accessibility of genome assemblies for global cattle breeds and other bovine species. *Genome biology* **24**: 139.
- Tavakoli N, Gibney D, and Aluru S. 2022. Haplotype-aware variant selection for genome graphs. In *Proceedings of the 13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*, pp. 1–9.
- Vaddadi K, Mun T, and Langmead B. 2024. Minimizing reference bias with an impute-first approach. *bioRxiv* pp. 2023–11.