

# Aligning pangenomes with hierarchical wavefront algorithm

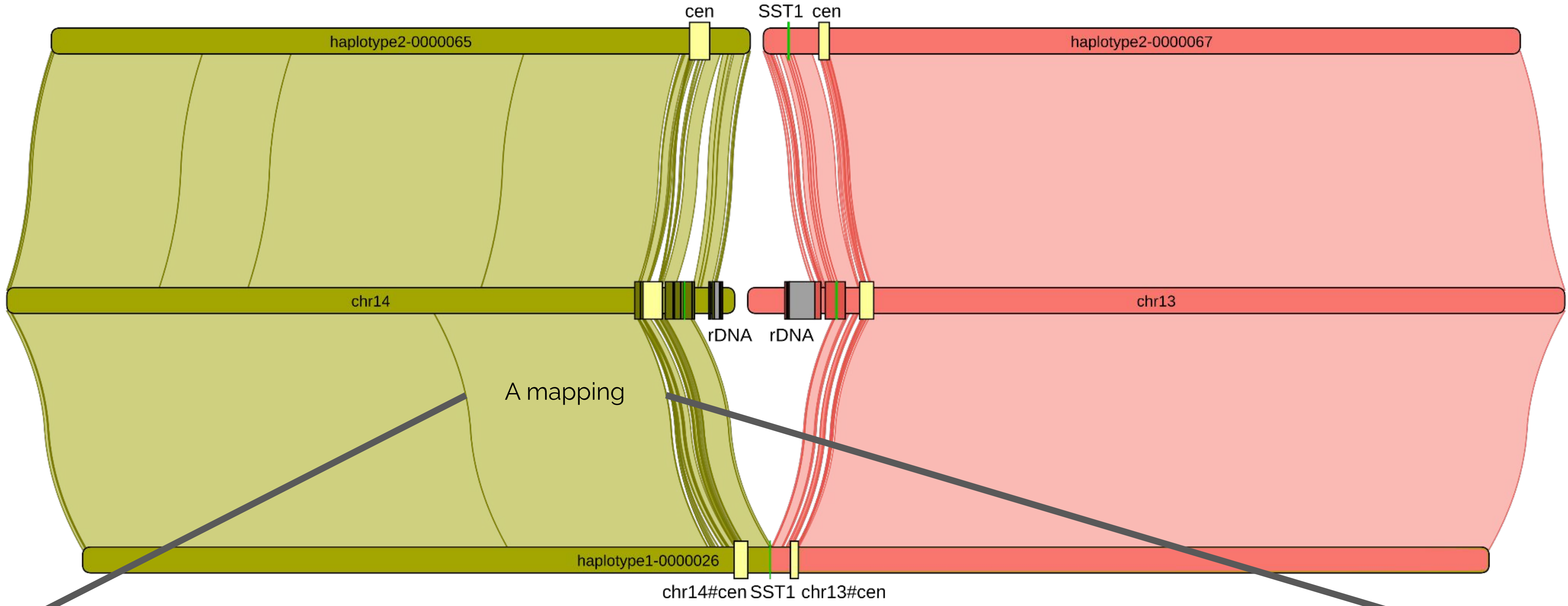
Andrea Guarracino<sup>1</sup>, Zhigui Bao<sup>2</sup>, Bryce Kille<sup>3</sup>, Njagi Mwaniki<sup>4</sup>, Santiago Marco-Sola<sup>5</sup>, and Erik Garrison<sup>1</sup>

<sup>1</sup> Department of Genetics, Genomics and Informatics, University of Tennessee Health Science Center, Memphis, TN 38163, USA; <sup>2</sup> Department of Molecular Biology, Max Planck Institute for Biology Tübingen, Max-Planck-Ring 9, 72076 Tübingen, Germany; <sup>3</sup> Department of Computer Science, Rice University, Houston, TX, United States; <sup>4</sup> Department of Computer Sciences, University of Pisa, Pisa, 56127, Italy; <sup>5</sup> Department of Computer Sciences, Barcelona Supercomputing Center, Barcelona, 08034, Spain.

## 1) Introduction

- **Sequence alignment** is a core task in bioinformatics.
- Due to the computational complexity of exact alignment, **heuristic methods** are closely tied to the sequence types being considered.
- **Seeding and extension** strategies are driven by seeds that are chained and filtered to find candidate regions for base-level alignment.
- Increasing read lengths and availability of big whole-genome assemblies rendered such approaches **overly sensitive**.
- The dramatic **increase in scale** poses significant challenges for current methods.

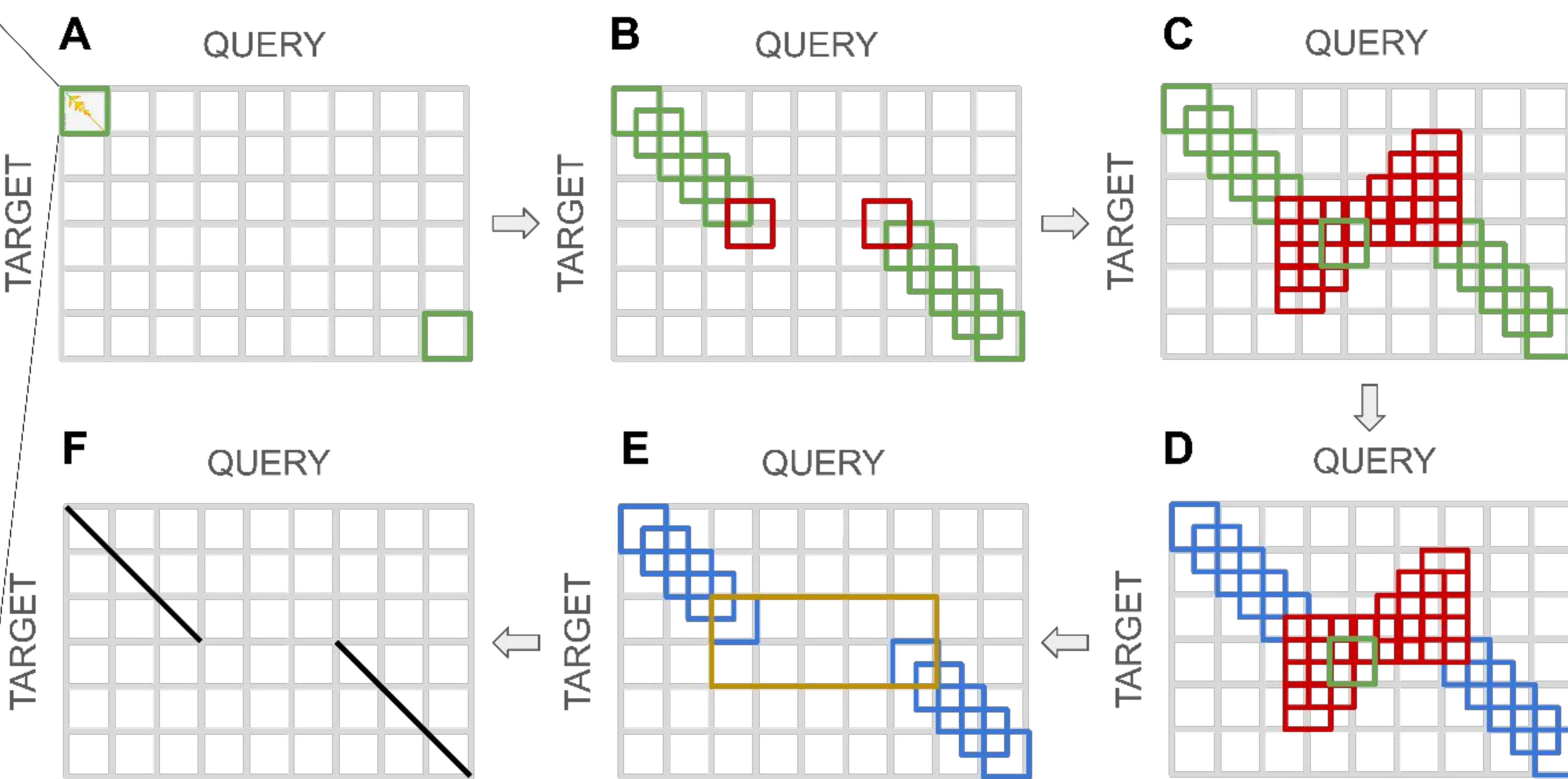
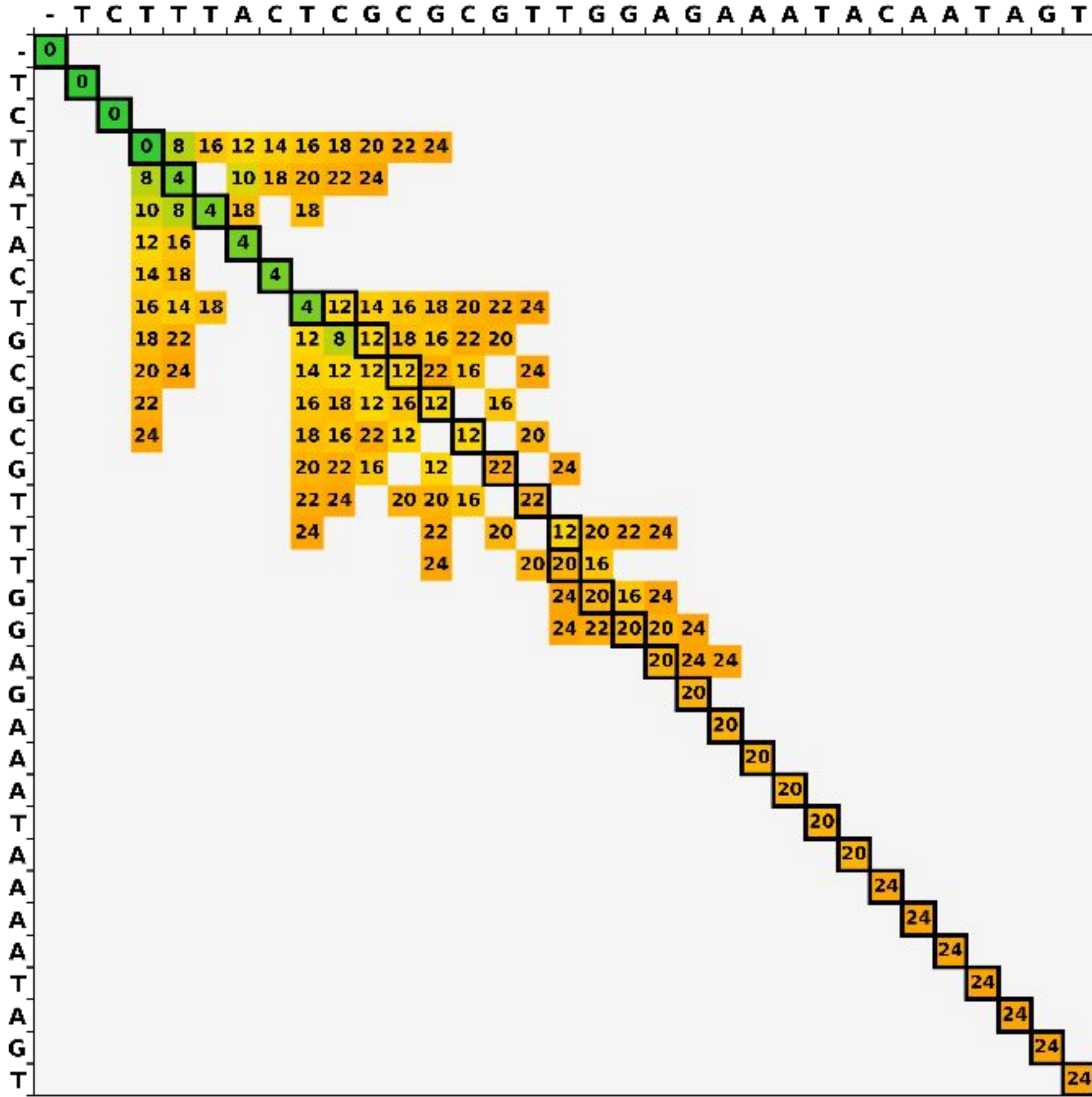
## Sequence mappings between human **acrocentric chromosomes**



## 2) Method

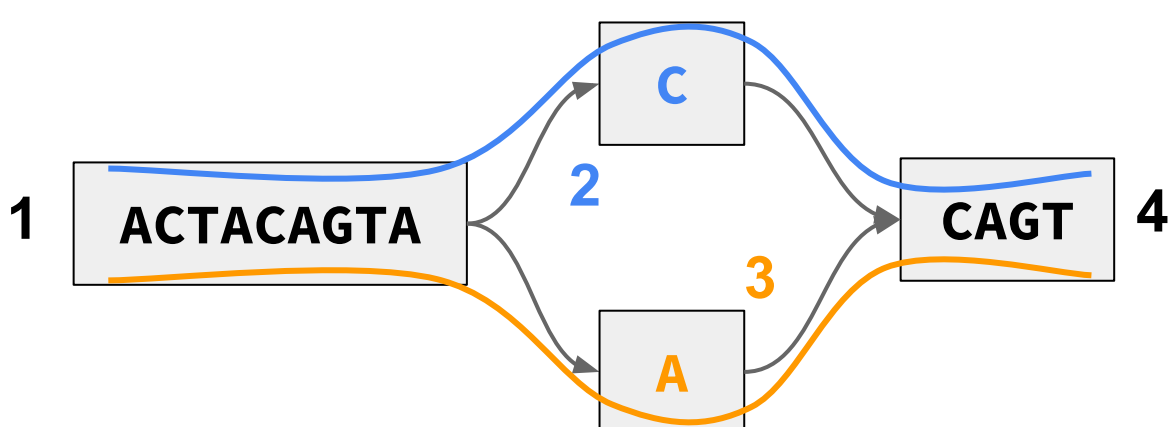
- We improve alignment performance for this new scale by **sparsifying** key information that drives mapping and alignment.
- We present **WFMASH**<sup>1</sup>, a **sequence aligner** for large pangenomes.
- WFMASH leverages a 2-step strategy:
  - a) It first applies a locality-sensitive hashing algorithm, adapted from MashMap<sup>2,3</sup>, to determine **syntenic region boundaries** (aka mappings) between long DNA sequences;
  - b) Then, a **hierarchical implementation of the wavefront alignment (WFA) algorithm**<sup>4,5</sup> allows computing the base-level alignment of the identified regions.

## 3) Hierarchical wavefront alignment algorithm scheme



## 4) Results

- With a **plant pangenome** made with 10 **tomato** haploid assemblies, WFMASH computes the all-vs-all alignment ~2.7X faster than MINIMAP2<sup>6</sup>, on average, while requiring ~5.6X less memory and aligning longer sequences.
- With an **inter-species pangenome** made with 16 **primates** assemblies (from *Homo sapiens*, *Gorilla gorilla*, *Pan paniscus*, *Pan troglodytes*, *Pongo abelii*, *Pongo pygmaeus*, *Symphalangus syndactylus*), WFMASH computes the all-vs-all alignment in 70 hours using 128 threads on an AMD EPYC 7742 64-Core processor, with a memory peak of 202 GB.
- We use WFMASH in the Pangenome Graph Building (PGGB) pipeline<sup>7,8</sup> to build unbiased **pangenome variation graphs**.



## Alignment between two 70 kbp long syntenic regions of two human chromosomes 13

