# Population Genomics Analyses on pangenome graphs

## Flavia Villani
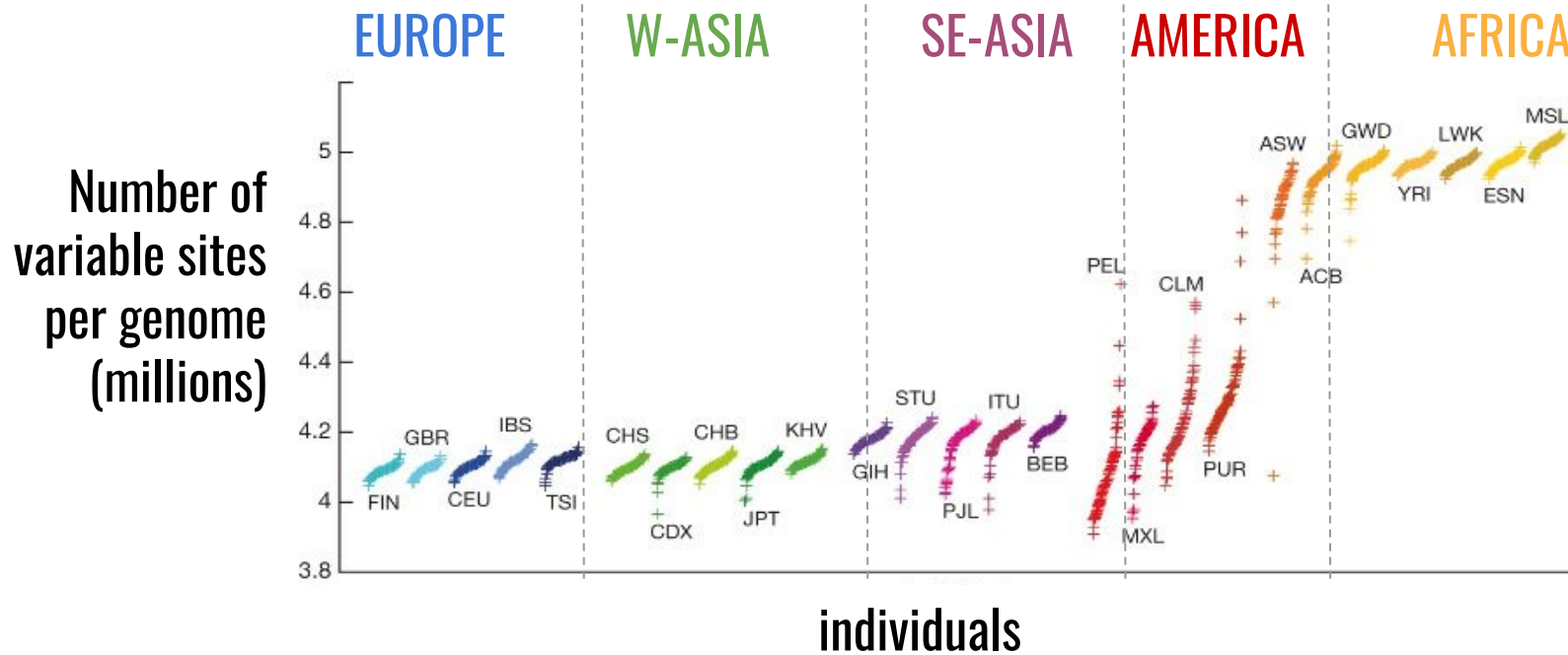
*Consiglio Nazionale delle Ricerche | Istituto di Genetica e Biofisica "Adriano Buzzati-Traverso" | Napoli*

# Population Genetics



Number of variable sites per genome (millions) vs individuals, grouped by EUROPE, W-ASIA, SE-ASIA, AMERICA, AFRICA.

4.3 M differences on average between two individuals

# Pangenomics approach for identification structural variants



Rebecca L Pollex and Robert A Hegele. "Copy number variation in the human genome and its implications for cardiovascular disease". In: Circulation 115.24 (2007),pp. 3130–3138
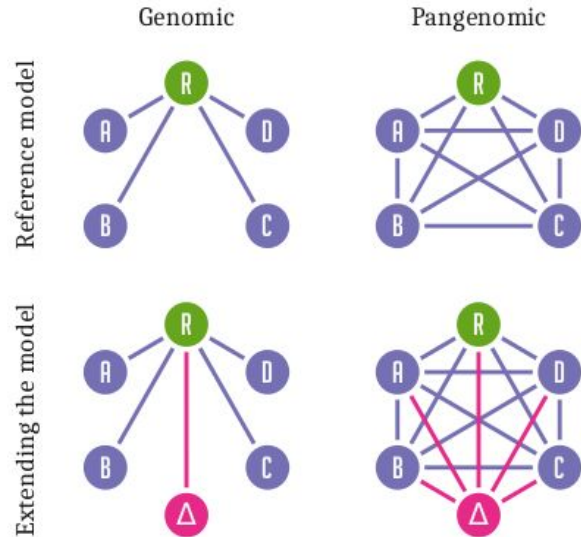
# Genomic *versus* pangenomic

Eizenga, Jordan & Novak, Adam & Sibbesen, Jonas & Heumos, Simon & Ghaffaari, Ali & Hickey, Glenn & Chang, Xian & Seaman, Josiah & Rounthwaite, Robin & Ebler, Jana & Rautiainen, Mikko & Garg, Shilpa & Paten, Benedict & Marschall, Tobias & Sirén, Jouni & Garrison, Erik. (2020). Pangenome Graphs. Annual Review of Genomics and Human Genetics.
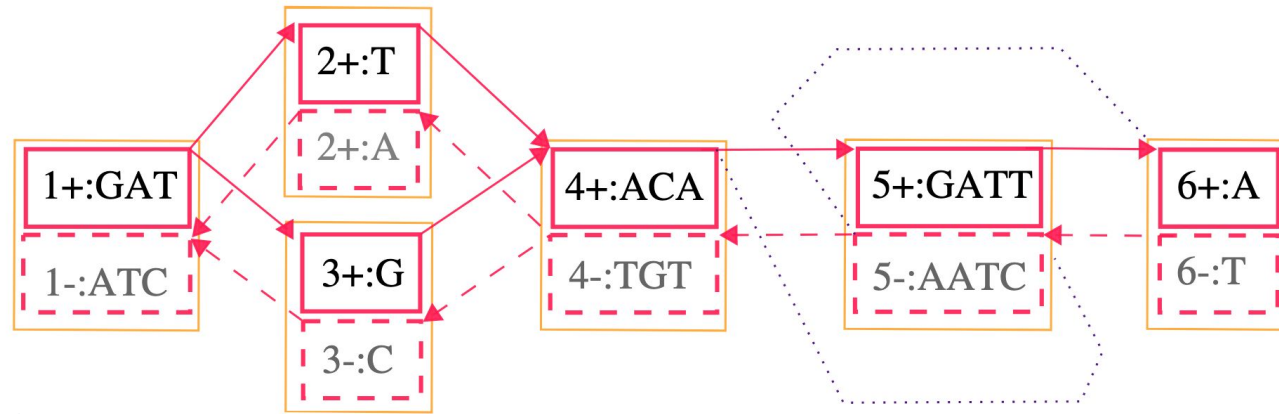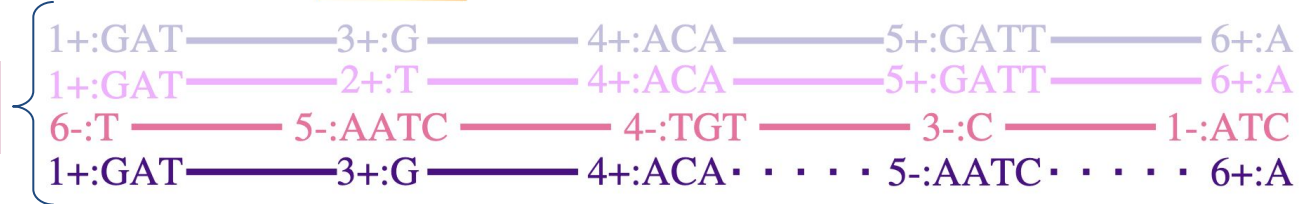
# Genomic *versus* pangenomic

Eizenga, Jordan & Novak, Adam & Sibbesen, Jonas & Heumos,
Simon & Ghaffaari, Ali & Hickey, Glenn & Chang, Xian & Seaman,
Josiah & Rounthwaite, Robin & Ebler, Jana & Rautiainen, Mikko
& Garg, Shilpa & Paten, Benedict & Marschall, Tobias & Sirén,
Jouni & Garrison, Erik. (2020). Pangenome Graphs. Annual
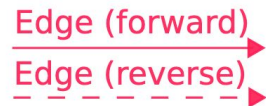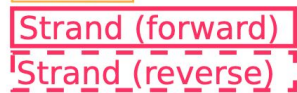Review of Genomics and Human Genetics.

# Graphic representation of a pangenome



Paths

1+:GAT — 3+:G — 4+:ACA — 5+:GATT — 6+:A
1+:GAT — 2+:T — 4+:ACA — 5+:GATT — 6+:A
6-:T — 5-:AATC — 4-:TGT — 3-:C — 1-:ATC
1+:GAT — 3+:G — 4+:ACA · · · · 5-:AATC · · · · 6+:A

Node

Strand (forward)    Edge (forward)    Edge (reversing)
Strand (reverse)    Edge (reverse)

Eizenga J. M*., Novak  A. M.*, Kobayashi E., Villani F., Cisar C., Heumos S., Hickey G., Colonna V., Paten B.  & Garrison  E. (2020).  Efficient dynamic variation graphs. *Bioinformatics Oxford.*

# Genetic variants in the linear and graphical model



**Pangenome (GFA)**

| 2:A |
| 4:ACA |
| 1:GAT |
| 3:T |

Path x1
Path y1

**Genomics (VCF)**

| #CHROM | POS | REF | ALT |
|--------|-----|-----|-----|
| x | 4 | A | T |

**Genomics standard analyses are based on linear representation of genomes**

# Goal

To develop a library of functions (vgpop) for population genetic analysis on pangenomic models

# Library vgpop

**Parsing pangenome**

bubblepop

**Population genetics**

num_sequences

num_segregatingsites

allele_frequencies

fst

**Format conversion**

gfa2vcf

seqgen2gfa+vcf

**Application**

Simulated data

Real data:
- HLA
- Sars-Cov2

https://github.com/Flavia95/VGpop

# Library vgpop

**Parsing pangenome**

bubblepop

**Population genetics**

num_sequences

num_segregatingsites

allele_frequencies

fst

**Format conversion**

gfa2vcf

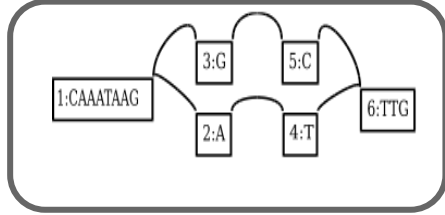seqgen2gfa+vcf
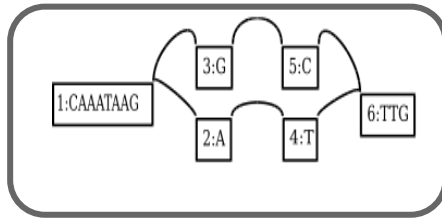
**Application**

Simulated data

Real data:
- HLA
- Sars-Cov2

https://github.com/Flavia95/VGpop
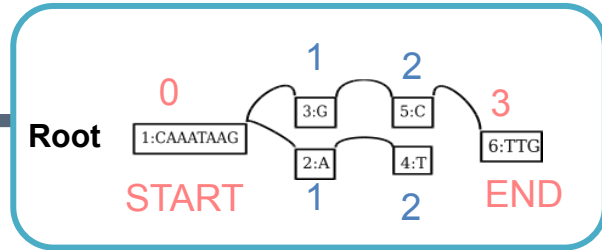
# bubblepop

## A. Graph

# bubblepop

A. Graph



B. Tree

# bubblepop



A. Graph

B. Tree

C. Bubble Calling

| | | Insertion | | Deletion | | SNV | |
|---|---|---|---|---|---|---|---|
| REF | AAAT | --- | TTTCT | GG | AGTTCTAT | T | ATAT |
| ALT | AAAT | AA | TTTCT | --- | AGTTCTAT | A | ATAT |

| POSITION | pos1 | pos2 | pos3 | pos4 |
|----------|------|------|------|------|
| **PATH1** | T | T | T | T |
| **PATH2** | A | G | A | T |
| **PATH3** | T | A | T | A |
| **PATH4** | A | T | A | A |

# Library vgpop

**Parsing pangenome**

bubblepop

**Population genetics**

num_sequences

num_segregatingsites

allele_frequencies

fst

**Format conversion**
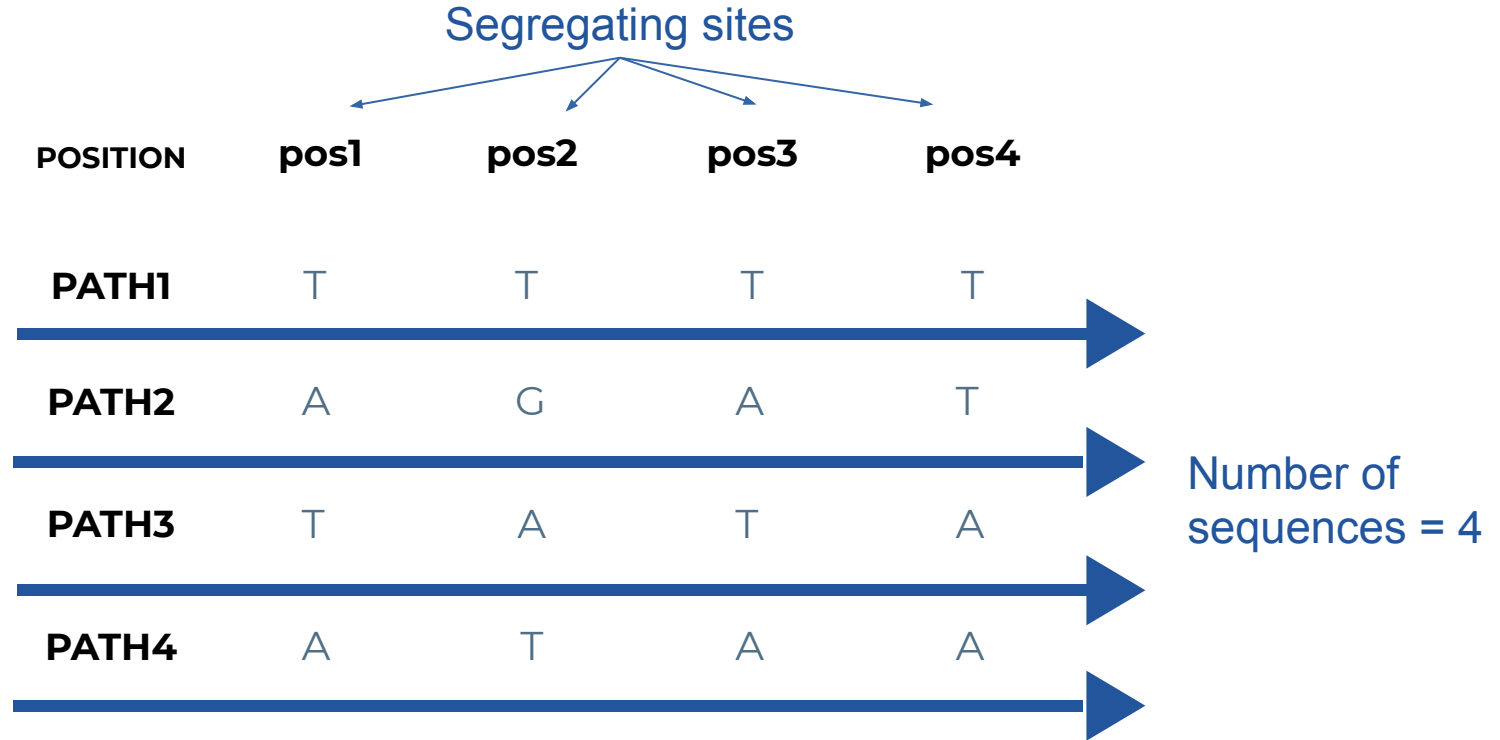
gfa2vcf

seqgen2gfa+vcf

**Application**

Simulated data

Real data:
- HLA
- Sars-Cov2

https://github.com/Flavia95/VGpop

# Segregation sites and sequences

Segregating sites

| POSITION | pos1 | pos2 | pos3 | pos4 |
|----------|------|------|------|------|
| PATH1 | T | T | T | T |
| PATH2 | A | G | A | T |
| PATH3 | T | A | T | A |
| PATH4 | A | T | A | A |

Number of sequences = 4

# Allele frequencies

2N = 20 chromosomes (APLOID)



AA AA Aa Aa Aa
aa aa aa aa aa

| ALLELE | A | a |
|---|---|---|
| ALLELE COUNTS | $n_A = 7$ | $n_a = 13$ |
| ALLELE FREQUENCIES | $f_A = \dfrac{n_A}{2N} = 0.35$ | $f_a = \dfrac{n_a}{2N} = 0.65$ |

# Wright's fixation index ($F_{st}$)



$$F_{st} \equiv \frac{\sigma^2}{\pi(1-\pi)}$$

POP1 0.45
POP2 0.40
POP3 0.50

POP1 0.45
POP2 0.45
POP3 0.45

$\sigma^2 = 0.005$
$\pi = 0.46$
$F_{st} = 0.02$

$\sigma^2 = 0.06$
$\pi = 0.46$
$F_{st} = 0.25$

Barbujani, G., & Colonna, V. (2010). Human genome diversity: frequently asked questions. *Trends in Genetics*, *26*(7), 285-295.

# Library vgpop

**Parsing pangenome**

bubblepop

**Population genetics**

num_sequences

num_segregatingsites

allele_frequencies

fst

**Format conversion**

gfa2vcf

seqgen2gfa+vcf

**Application**
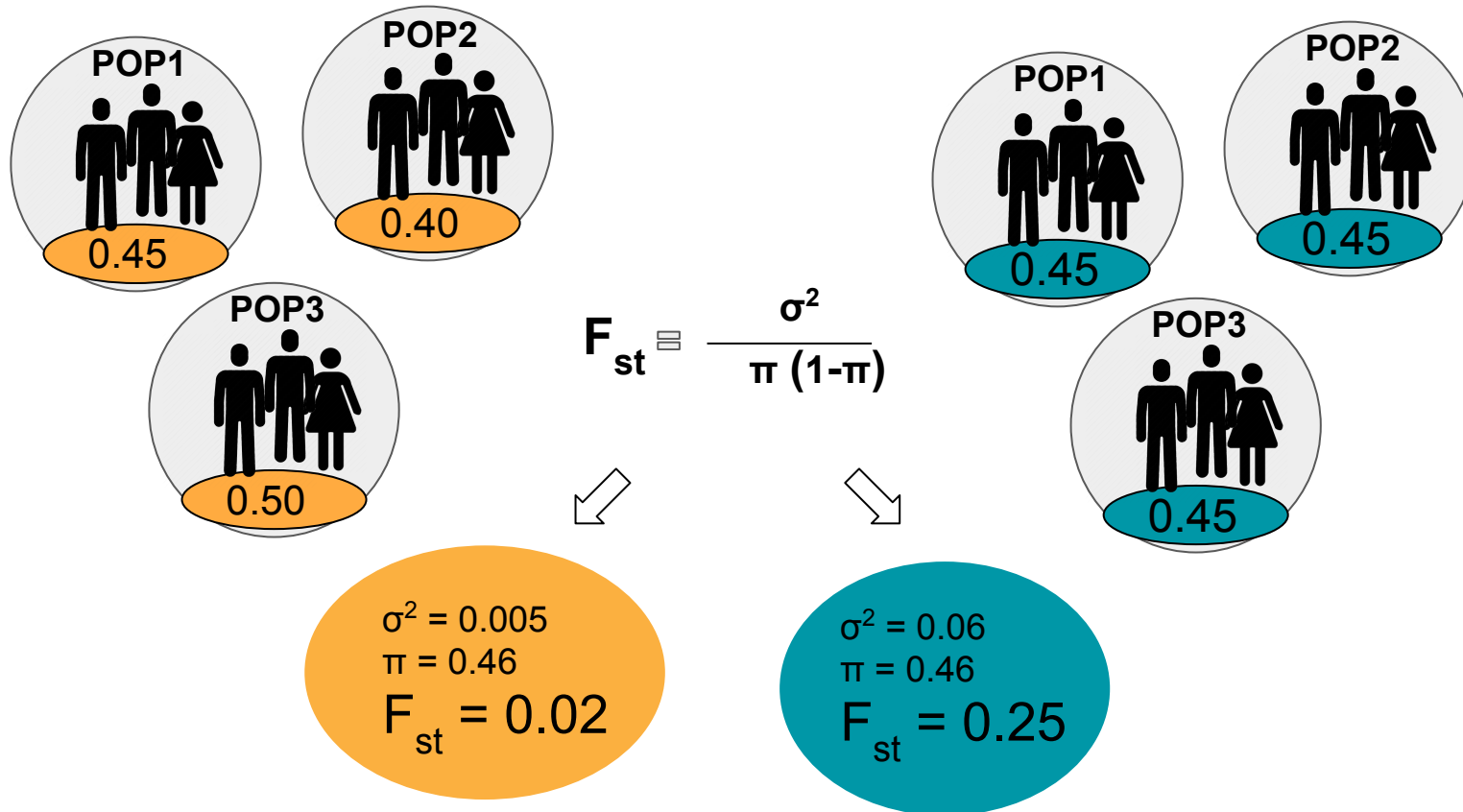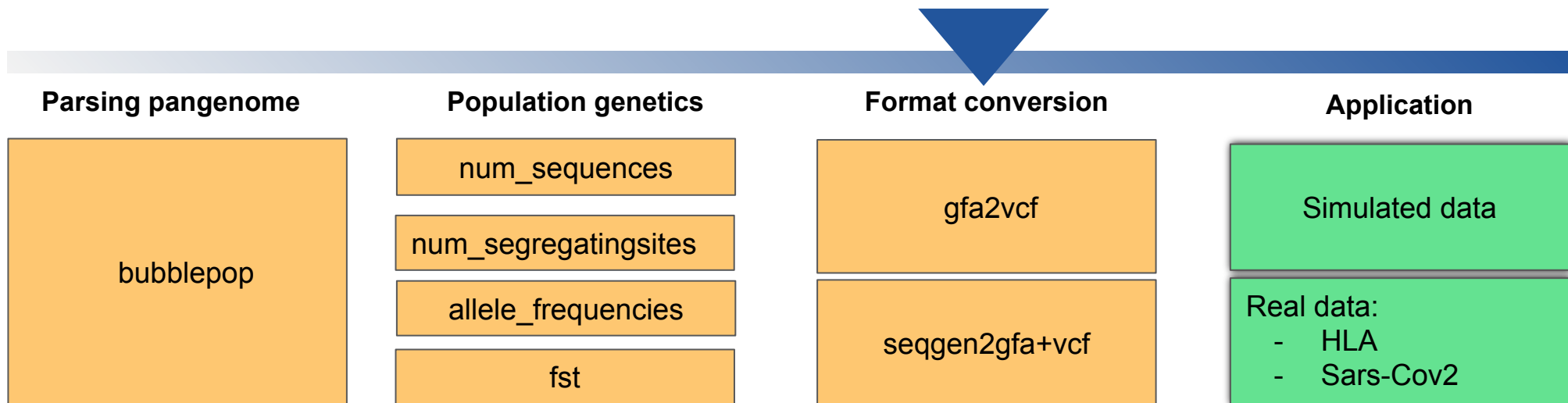
Simulated data

Real data:
- HLA
- Sars-Cov2

https://github.com/Flavia95/VGpop

# Format conversion

## Pangenomic model (GFA)



**gfa2vcf** →

## Linear model (VCF)

| #CHROM | POS | ID | REF | ALT | INFO |
|--------|-----|-----|-----|-----|----------|
| x | 9 | . | G | A | TYPE=snv |
| x | 10 | . | C | T | TYPE=snv |

## Simulation sequences (Seq-Gen)

```
2 10
Taxon1  ATCTTTGTAG
Taxon2  ATCCTAGTAG
```

**seqgen2gfa+vcf** →

## Pangenomic model (GFA)

```
H       VN:Z:1
S       1       CACTA
S       2       ATTA
L       1       +       2       + 0M
P       x       1+,2+   0M
```

## Linear model (VCF)

| #CHROM | POS | ID | REF | ALT | INFO |
|--------|-----|-----|-----|-----|----------|
| x | 2 | . | G | A | TYPE=snv |
| x | 3 | . | C | T | TYPE=snv |

# Implementation of vgpop in Rust

Rust is a programming language focused on performance and safety.

- ❖ Great **ecosystem** (Cargo, crates.io, docs.rs).
- ❖ Much **safer** than C++ while having a similar **speed.**
- ❖ Friendly and helpful **community.**
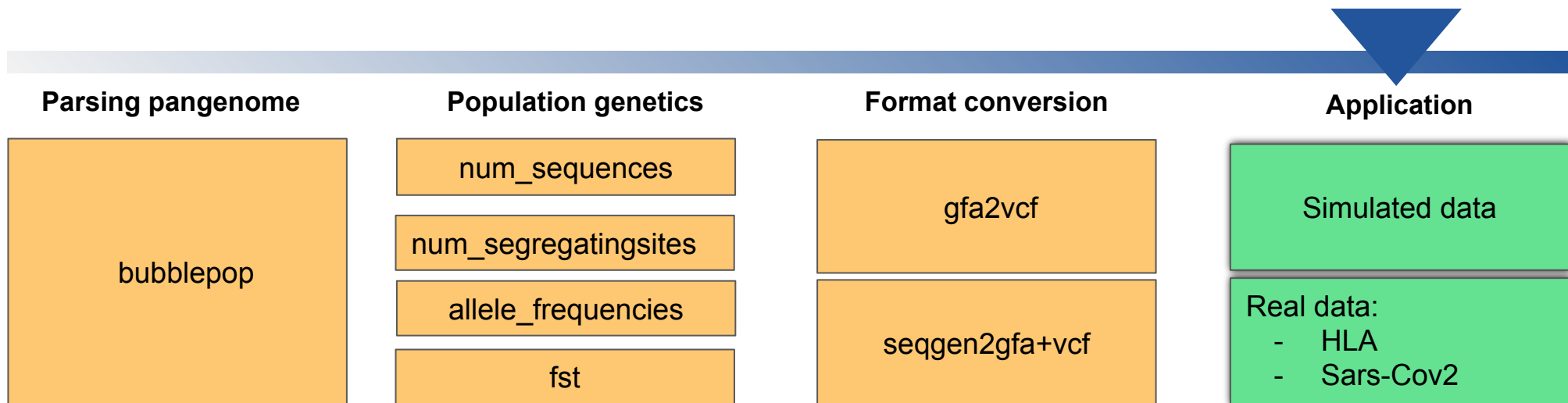- ❖ Used in many open source projects, such as **Firefox.**

https://www.rust-lang.org/

Francesco Porto
Gianluca Della Vedova

https://github.com/HopedWall/rs-gfatovcf

# Library vgpop

**Parsing pangenome**

bubblepop

**Population genetics**

num_sequences

num_segregatingsites

allele_frequencies

fst

**Format conversion**

gfa2vcf

seqgen2gfa+vcf

**Application**

Simulated data

Real data:
- HLA
- Sars-Cov2

https://github.com/Flavia95/VGpop

# F$_{st}$ on simulated data

# Workflow

```
┌─────────────────┐
│       ms        │
└─────────────────┘
         │
         ▼
┌─────────────────┐        ┌─────────────────┐
│   Simulation    │──────▶ │    Seq-Gen      │
│  variable sites │        └─────────────────┘
└─────────────────┘                 │
                                    ▼
                           ┌─────────────────┐
                           │   Simulation    │
                           │   sequences     │
                           └─────────────────┘
```

vgpop  code

Existent code

# Workflow

# $F_{ST}$ on 100 replicate use *vgpop e vcftools*



*allele_freq*
*fst*

# Library vgpop

**Parsing pangenome**

| |
|---|
| bubblepop |

| |
|---|
| bubblecall |

**Population genetics**

| |
|---|
| num_sequences |

| |
|---|
| num_segregatingsites |

| |
|---|
| allele_frequencies |

| |
|---|
| fst |

**Format conversion**

| |
|---|
| gfa2vcf |

| |
|---|
| seqgen2gfa+vcf |

**Application**

| |
|---|
| Simulated data |

Real data:
- HLA
- Sars-Cov2

https://github.com/Flavia95/VGpop

# Allele frequencies on HLA

# Gene HLA-E

Pangenome
9 sequences

# Gene HLA-E



Pangenome
9 sequences

REF

$$\text{Freq} = \frac{6}{9}$$

| GENE | PANGENOME | POSITION | REF | ALT | FREQ |
|------|-----------|----------|-----|-----|------|
| HLA-E | HLAE-3133 | 551 | T | C | 0.67 |

# Variant discovery in HLA with rust implementation

❖ From 12 sequences
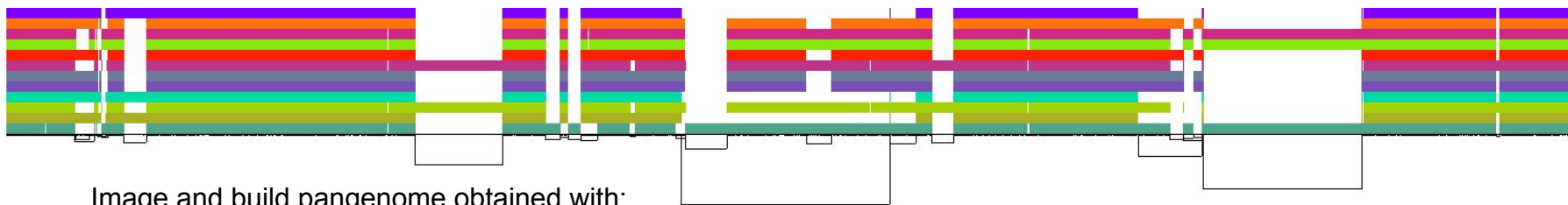❖ Size: 163416 nucleotides
❖ Run time: ~0.1s
❖ Variants found: 7505



Image and build pangenome obtained with:
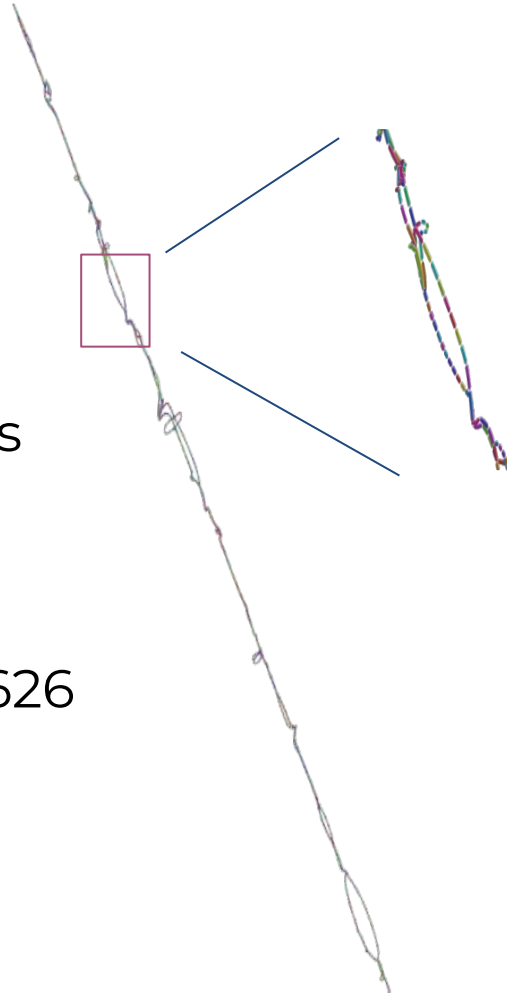https://github.com/pangenome/pggb

Code available at:
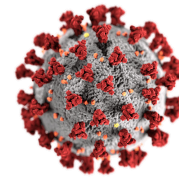https://github.com/HopedWall/rs-gfatovcf

# Variant discovery in Sars-Cov2 with rust implementation



- ❖ From 15127 genomes
- ❖ 1.2 Gbytes
- ❖ 78571 fragments
- ❖ Run time: ~16m
- ❖ Variants found: 294626

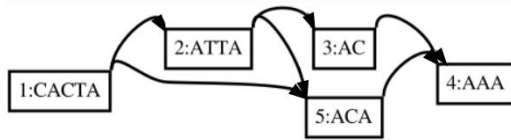**COVID-19 PubSeq**

Data available at
http://covid19.genenetwork.org/

Andrea Guarracino
Pjotr Prins

# Conclusion and next steps

## vgpop

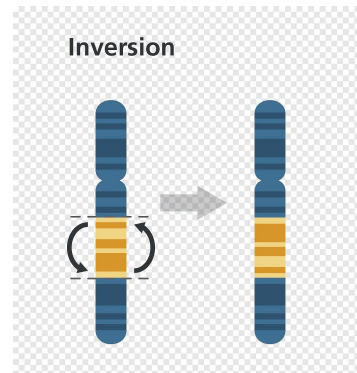Software for population genetics analyses on pangenomes

## Rust



Adding parallel computing to increase performances

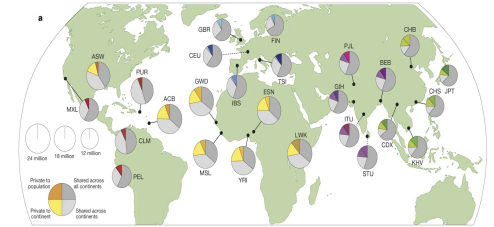https://crates.io/crates/gfautil

## Structural variation

Little considered in the standard population genetics analysis



## Population genomics analyses

Based on haplotype and on the differentiation of frequencies between populations

**IGB-CNR (US)**
Vincenza Colonna
Silvia Buonaiuto
Gianluca Damaggio
Giuliana D'Angelo

**University of Milano Bicocca (Italy)**
Francesco Porto
Gianluca Della Vedova

**University of Rome Tor Vergata (Italy)**
Andrea Guarracino

**Department of Genetics, Genomics and Informatics (UTHSC)**
Pjotr Prins
Robert W. Williams
Christian Fischer

**UCSC (US)**
Erik Garrison

*THANKS FOR YOUR ATTENTION!*