# A pangenome for the expanded BXD family of mice

Flavia Villani<sup>1</sup>, David G. Ashbrook<sup>1</sup>, Andrea Guarracino<sup>2</sup>, Simon Heumos<sup>3</sup>, Christian Fischer<sup>1</sup>, Hao Chen<sup>1</sup>, ,Pjotr Prins<sup>1</sup>, Erik Garrison<sup>1</sup>, Robert W. Williams<sup>1</sup> and Vincenza Colonna<sup>1,4</sup>

<sup>1</sup>UTHSC, GGI, Memphis, TN, USA; <sup>2</sup>University of Tor Vergata, Biology, Rome, Italy; <sup>3</sup>Quantitative Biology Center (QBiC), University of Tübingen, Tübingen, Germany; <sup>4</sup>CNR, IGB, Naples, Italy



#### seqT seq10 seq9 Sample2 seq8 seq7 seq6 seq5 seq4 seq3 seq2

seq2 seq3 seq4 seq5 seq6 seq6 seq9 seq9 seq10 seq11

#### (A) ALL-VS-ALL ALIGNMENTS WITH WFMASH

A hierarchical implementation of the

#### **(B) GRAPH INDUCTION** WITH <u>SEQWISH</u>

2. Build a pangenome with the PanGenome Graph Builder (PGGB)

D

D

To build a variation graph, we collapse the components of the graph connected by alignments and retrace the paths through the original graph.

В

В



smoothed graph

#### **(C)** GRAPH NORMALIZATION WITH <u>SMOOTHXG</u>

A partial order alignment is run for each part of the graph, resulting in a normalized graph with base-level resolution of all

all genomes are related to a reference genome.

In a **pangenomic** setting<sup>1</sup>, we model relationships between all the direct genomes in our analysis.

https://pangenome.github.io/

wavefront algorithm<sup>2</sup> allows us to obtain base-level global alignments for all mappings.



Q

variant classes.

**PGGB** has three distinct phases (A) all-to-all alignment with <u>wfmash</u>, (**B**) graph induction with seqwish, and (C) normalization with smoothxg, which produces the resulting graph shown.

### 3. A pangenome for the BXD family

- 152 BXD strains sequenced (10x technology).
- 148 BXD strains assembled (30-60 thousand contigs).
- Contig length peak at 31kb.



## A graphical representation of the pangenome of 148 family members (all isogenic inbred)



**PHASE 2: VARIANT CALLING** 



PHASE 3: VG vs TRUTH SETS



In DBA/2J we called from the pangenome with  $\underline{vq}$  218,139 simple microvariants (SNPs, MNPs, and INDELs), that is 45,208 and 39,937 more variants compared to GATK applied to 10x and PACBIO sequence data, respectively.

- Identification of complex variants.
- Use nanopore technology. \*

