



# Collaborative Cross Graphical Genome

Hang Su<sup>1</sup>, Ziwei Chen<sup>2</sup>, Jay Rao<sup>2</sup>, Fernando Pardo Manuel de Villena<sup>3</sup>, Leonard McMillan<sup>1,2</sup>

<sup>1</sup> Curriculum in Bioinformatics and Computational Biology, University of North Carolina at Chapel Hill, Chapel Hill, NC, 27599

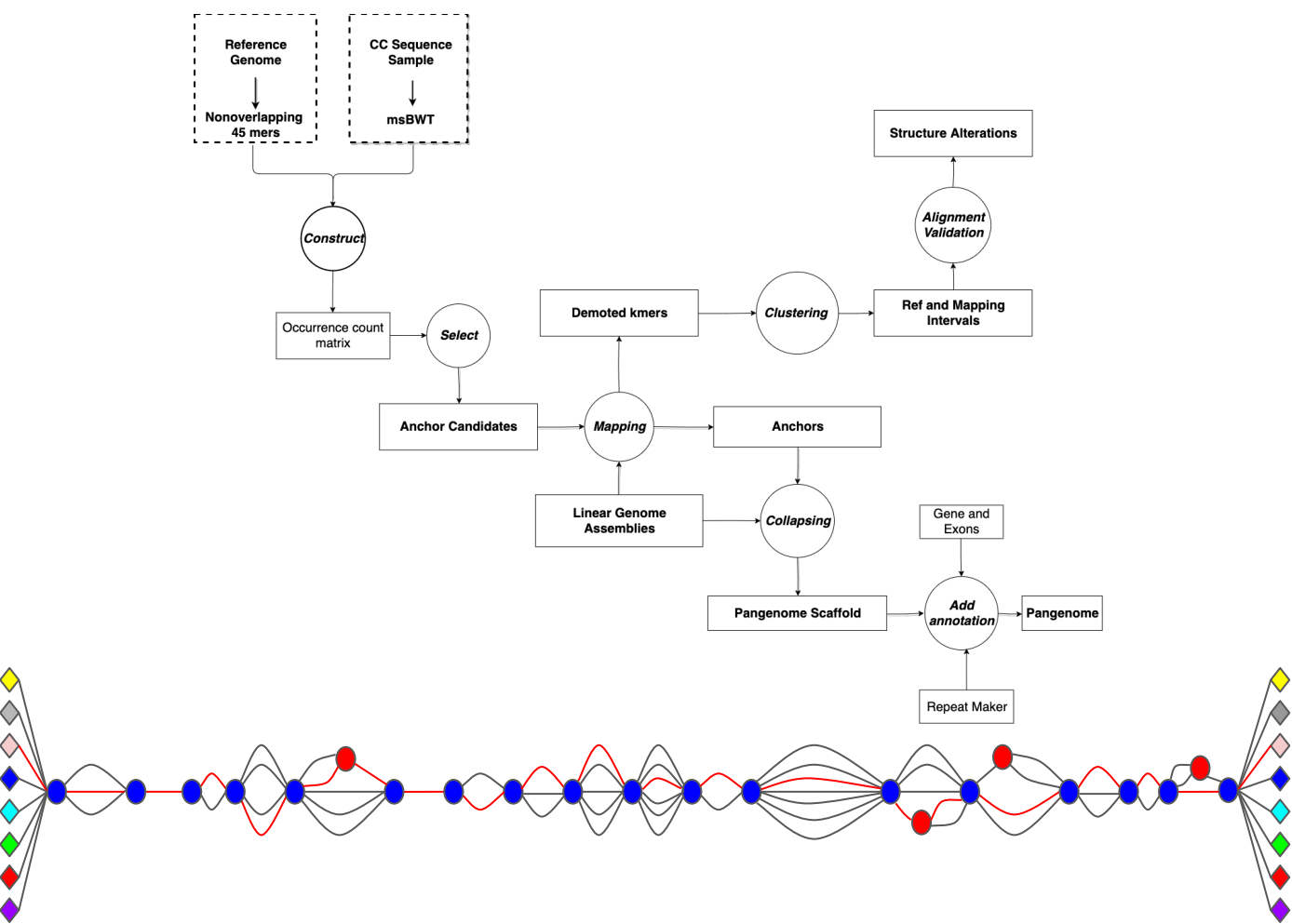
<sup>2</sup> Department of Computer Science, University of North Carolina at Chapel Hill, Chapel Hill, NC, 27599.

<sup>3</sup> Department of Genetics, University of North Carolina, Chapel Hill, NC 27599.



## Abstract

Inexpensive and fast genome sequencing has enabled the assembly of multiple intra-specific genomes that capture more genetic diversity than any single reference assembly. It has been widely suggested that pangenome reference assemblies, which incorporate multiple genomes into a single representation, are the path forward, but there are few standards for, or instances of practical pangenome representations suitable for large eukaryotic genomes. We present a graph-based pangenome representation and an instance of it for a widely-used recombinant-inbred mouse genetic reference population known as the Collaborative Cross. Our pangenome representation leverages existing standards for genomic sequence representations with backward-compatible extensions to describe graph topology and genome-specific annotations along paths. It packs 82 genomes in a single graph representation which directly captures important notions relating genomes such as identity-by-descent between mouse strains, and highly variable genomic regions. The introduction of special anchor nodes with sequence data provide a valid position-reference framework that divided large eukaryotic genomes into homologous segments and addresses most of the graph-based position reference issues, such as *monotonicity*, *backward-compatibility*, *spatiality*. Paralleled edges between a pair of anchors contains all variants and provide effective way for orthogonal genome comparison and visualization. Furthermore, our graph structure allows annotations to be placed in multiple genomic contexts and simplifies their maintenance as the assembly improves. The CC reference pangenome provides an open framework for combining and incorporating new genome sequences and an effective representation for searching, annotating, comparing and visualizing biological features at different scales, which also facilitates tool chain development for downstream analysis.



**CCGG:** A single *directed, k-partite* graph representation for 82 mouse genomes, including 8 founder genomes and 74 CC genomes which are mosaic of the 8 founders

**Nodes:** Anchors: *conserved, unique 45mers* with *monotonically* increasing order  
Floating nodes, Source nodes and Sink nodes

**Edges:** Sequences between nodes, sequence diversity.

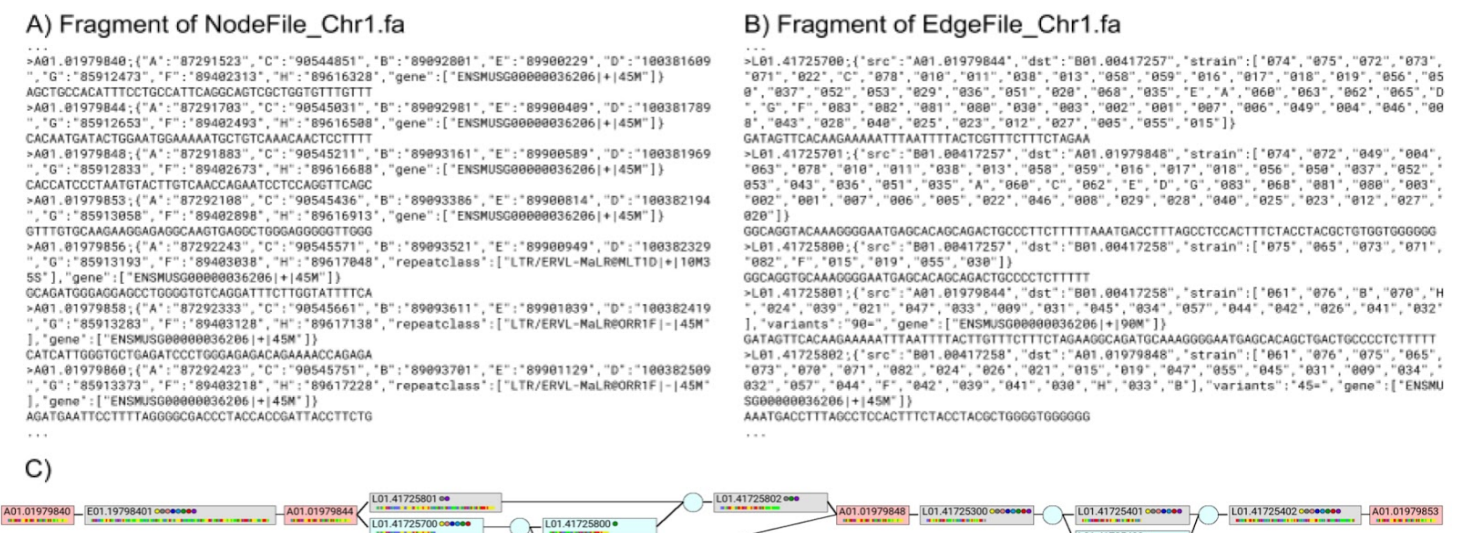
**Path:** A series of nodes and edges sharing common source and destination anchors

**Gap:** Sequences between two anchors

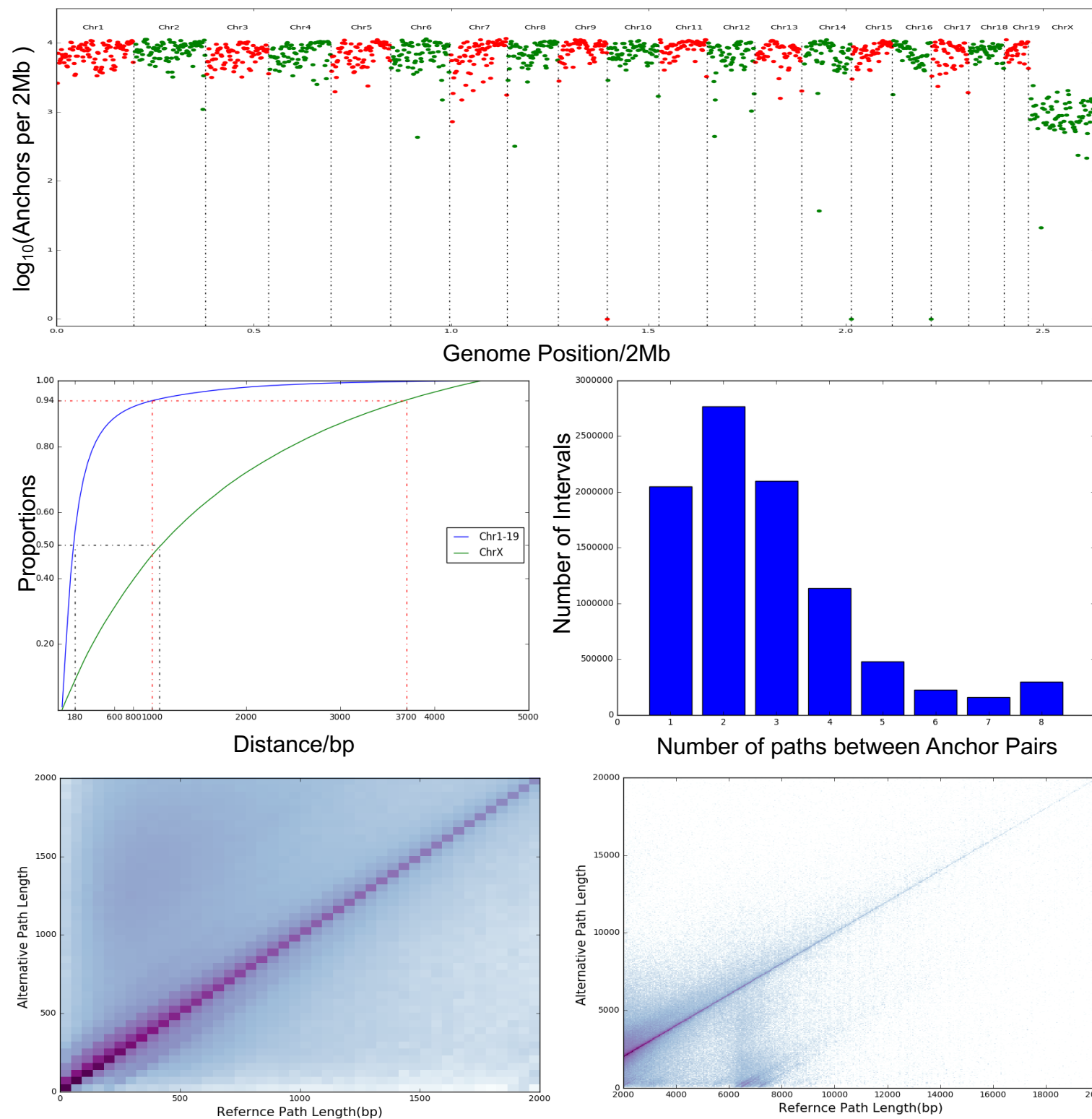
**Genomic Contig:** a Path from a Source to a Sink node

## Collaborative Cross Graphical Genome

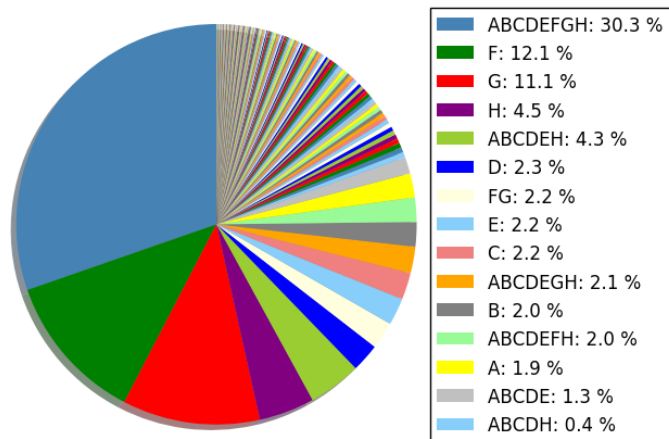
### CCGG Overview



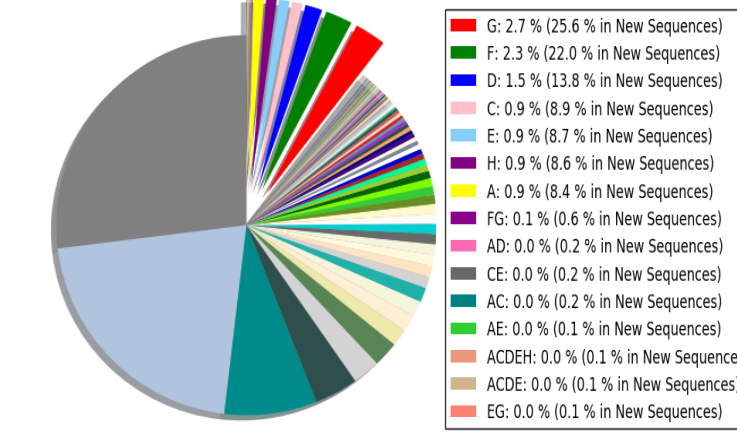
### Graph Properties



### Edge Sharing

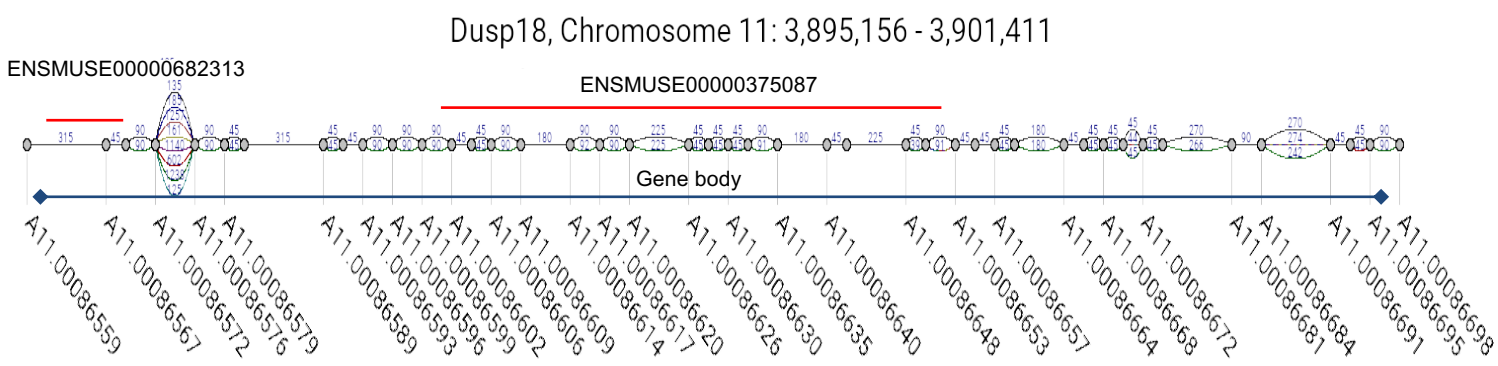


### Sequence Sharing

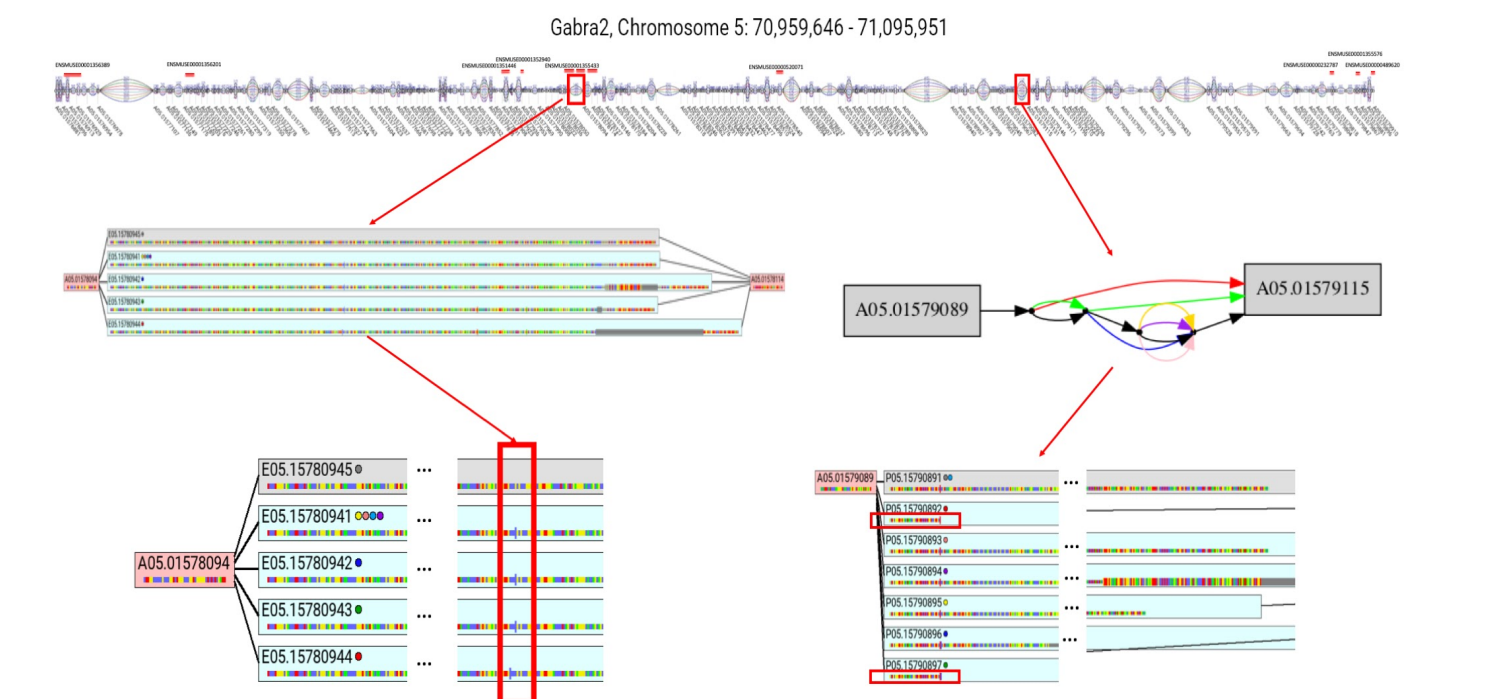


## CCGG Graph Structure in Functional Regions

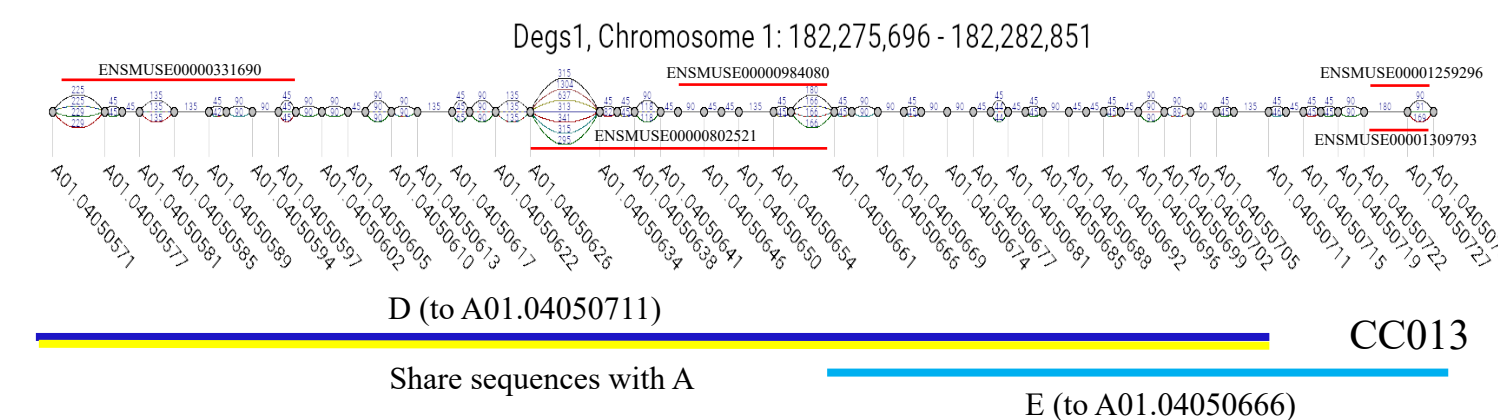
### Genomic Structure: identify conserved and highly variable regions



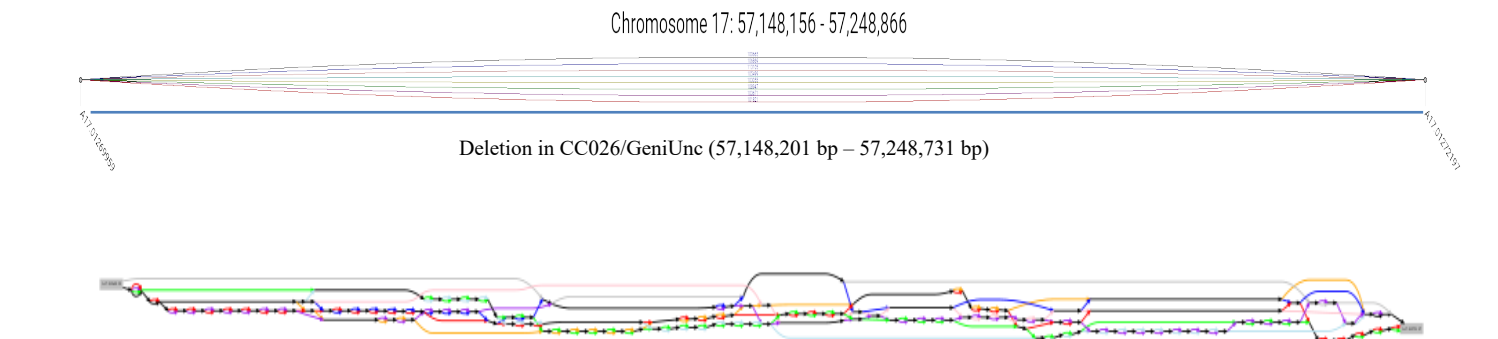
### Genomic Structure and Variants representation



### CC Recombination in more than 2000 genes



### Genomic Structure of Long Gaps



- Reference Pangenome for the Collaborative Cross
- Graph reveals sequence commonality and diversity
- Functional variants and recombination analysis in CCGG
- CCGG establish a better standard for genome annotation and coordinates
- Downstream tool chains: aligners, variant calling...