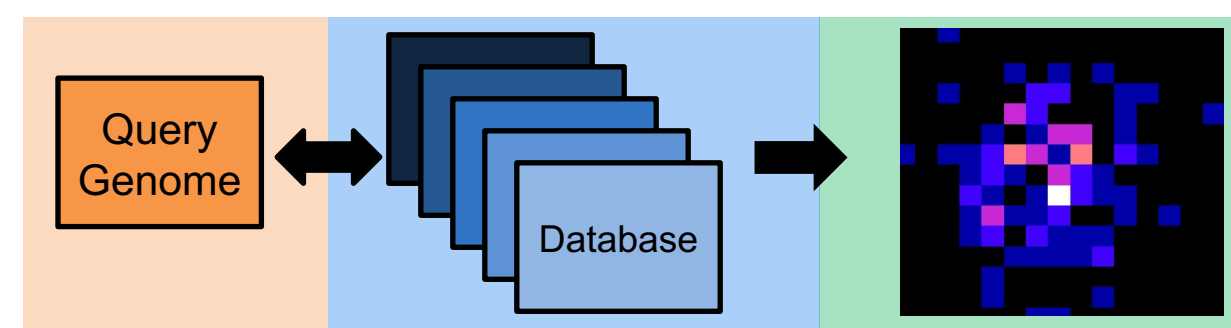


## BACKGROUND

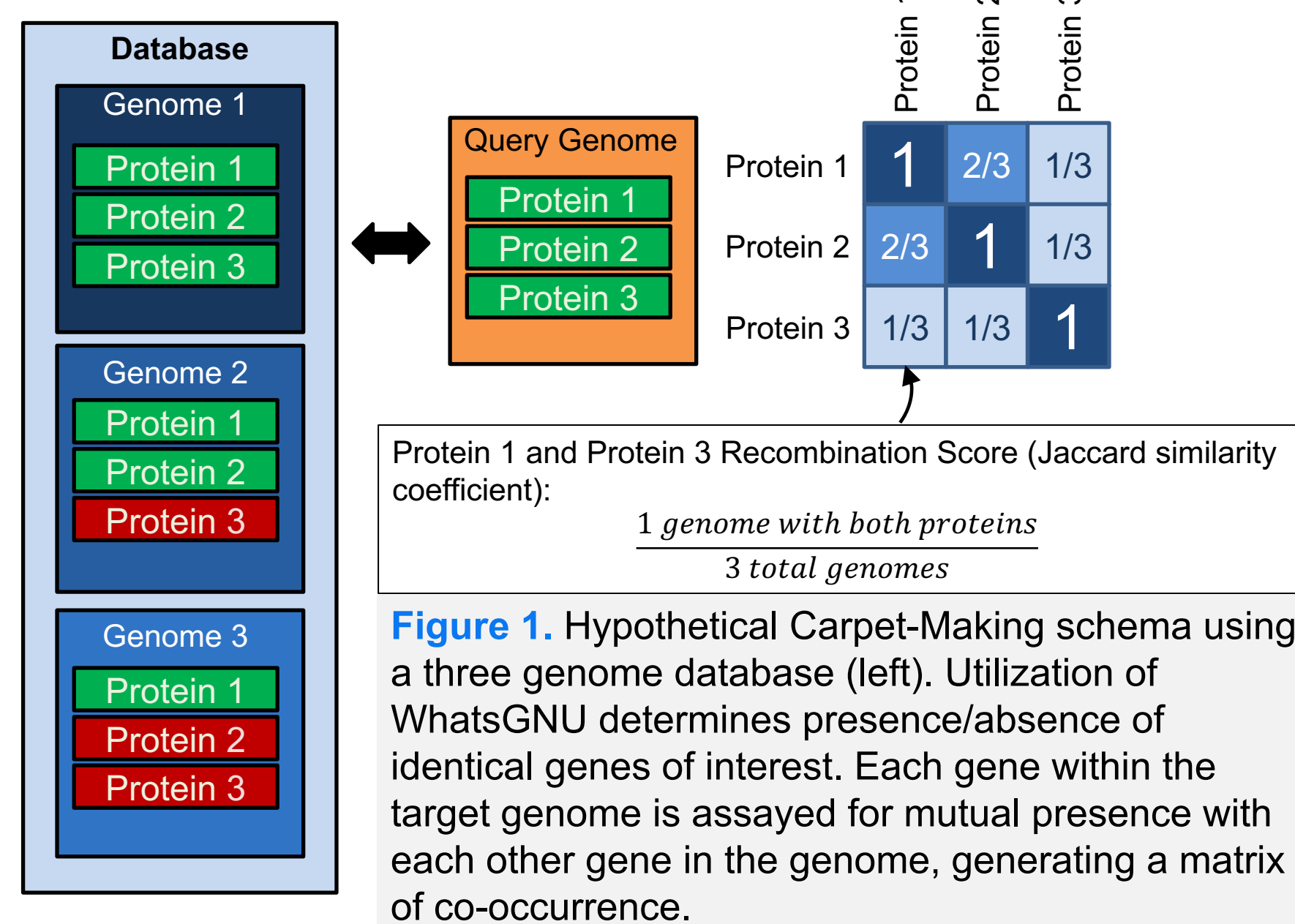
Genomic recombination is a generator of biological diversity and plays an important role in microbial adaptation to new environments, hosts, and niches. Recombination detection previously relied on genomic sequence alignment and phylogenetic or comparative techniques that are computationally expensive, especially with vast increases in available whole genome sequences. Here we present a database-driven technique that is alignment-free, does not rely on phylogeny or sequence similarity, and can be used rapidly on single genomes on a standard desktop. The tool, Redcarpet, combines the analytic output from our recently developed WhatsGNU algorithm<sup>1</sup> with a MinHash technique<sup>2</sup>. This operation is predicated on the idea that identical genes are more likely to appear in the same set of genomes if they share a common evolutionary history

## METHODOLOGY

Redcarpet takes in a **single query genome**, and for each encoded protein, determines **the set of genomes in a database** that contain an exact protein sequence match. It then **computes the Jaccard similarity coefficient between genome sets** for all pairwise protein comparisons in the genome.

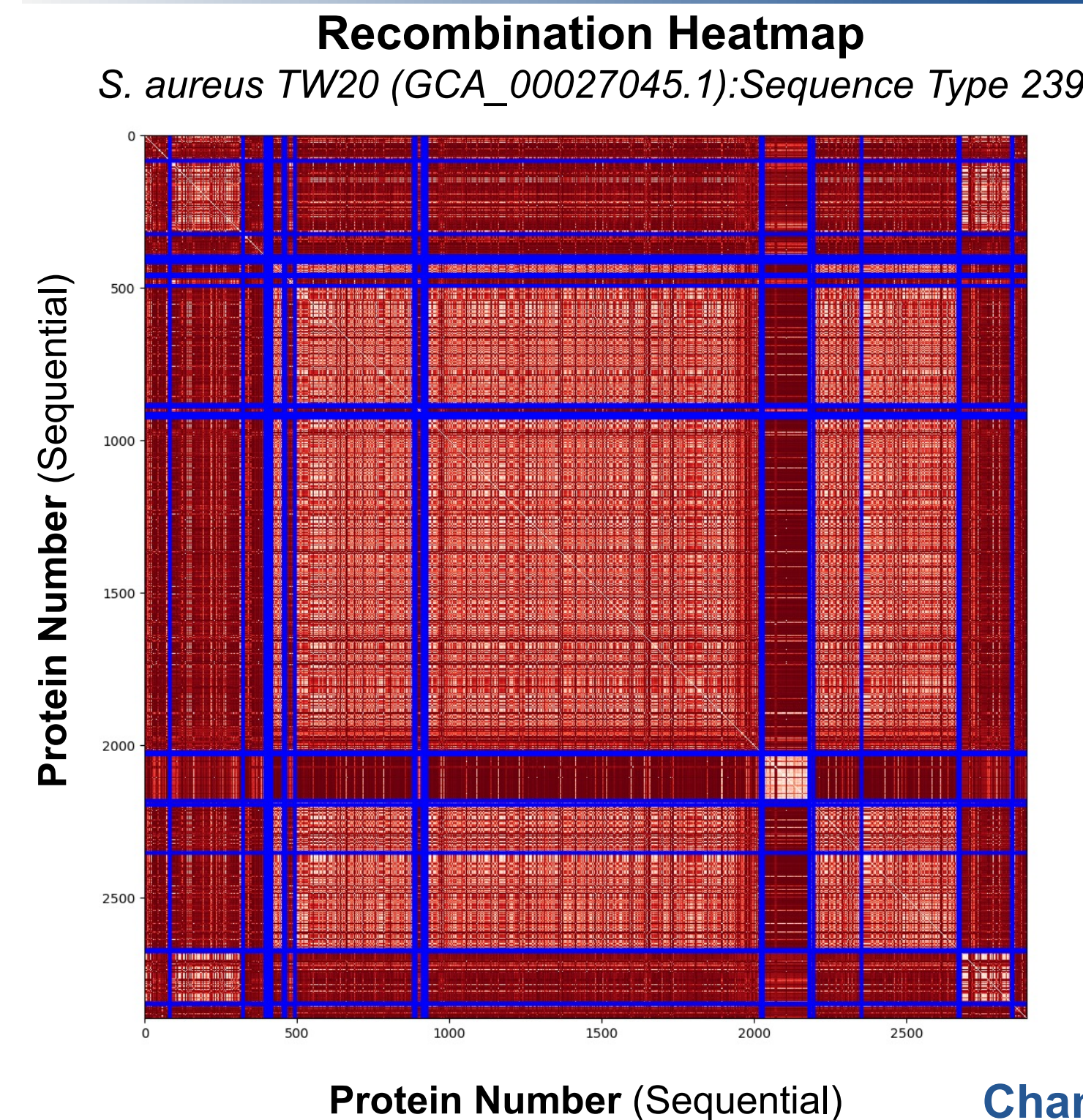


## CALCULATION THEORY



## RESULTS

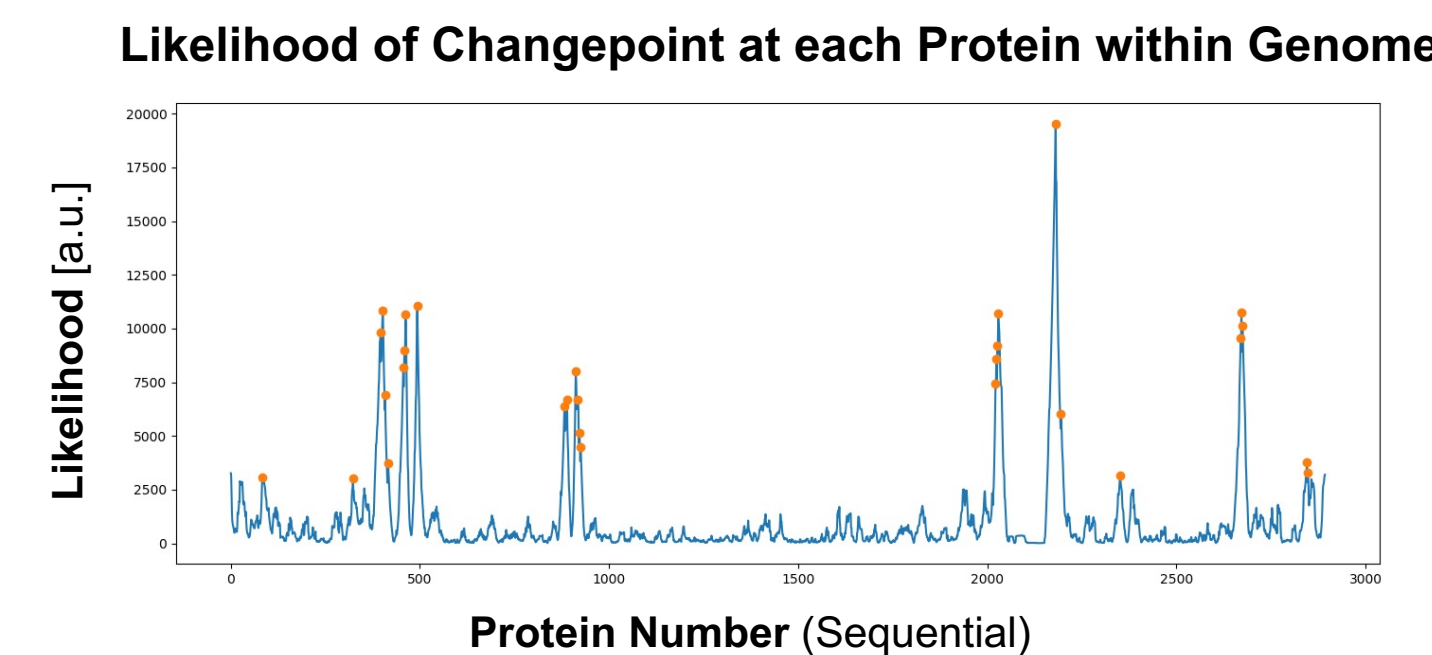
### Recombination



The output of Redcarpet is a pairwise, genome-set similarity matrix that can be visualized as a 2-D heatmap ordered by the gene location on the chromosome. The heatmap provides a visual tool for identifying recombination tracts. Beyond recombination detection, this tool can be further expanded to explore both (1) the phylogenetic history of intra-genomic recombinant regions and (2) the likely genomic regions of recombination.

**Figure 2.** Similarity Matrix (Redcarpet) for *Staphylococcus aureus* genome (GCA\_000027045), within ST239 that is known for large-scale evolutionary recombination. Use of 10,350 publicly-accessible NCBI *S. aureus* genomes were utilized as the comparative database in the generation of this heatmap. This analysis reveals the potential locations of recombination ("squares" within the heatmap), with blue lines added to depict the edge of the squares (i.e., the changepoints), as calculated within Figure 3.

### Changepoint Analysis



The recombinant changepoints, defined as likely location of recombination were calculated by comparing the protein values within a set number (~30) of proteins away from each protein in the genome. The local maxima of the generated curve (Fig. 3), approximated the likely region of recombination.

**Figure 3.** Regional Logarithmic-scaled parameter of changepoint likelihood for the same genome as studied within Fig. 2 surveying +/- 30 proteins determines likelihood of recombination: this can be utilized to analytically predict changepoints (orange dots) within the similarity matrix.

### Evolution and Recombinant Ancestry

