

# Pseudo-chromosome assembly of large and complex genomes using multiple references

Mikhail Kolmogorov<sup>1</sup>, Brian Raney<sup>2</sup>, Joel Armstrong<sup>2</sup>, Duncan Odom<sup>3,4</sup>, Paul Flicek<sup>5</sup>, David Thybert<sup>5,6</sup>, Benedict Paten<sup>2</sup> and Son Pham<sup>1</sup>

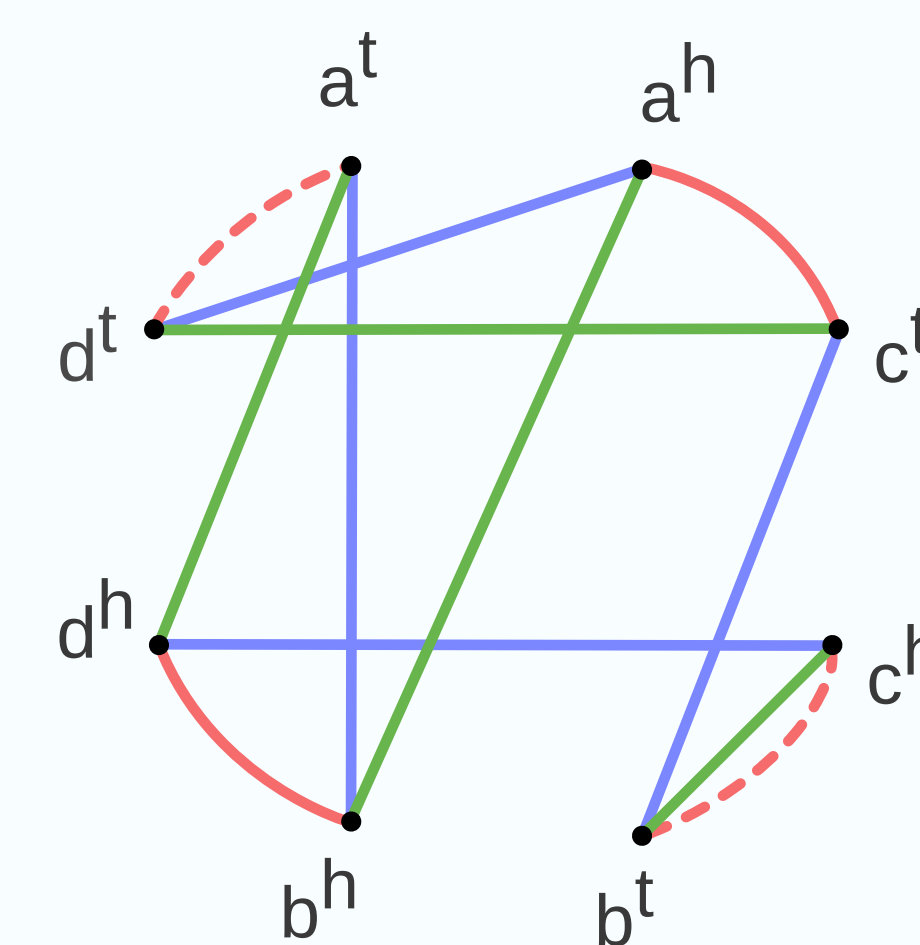
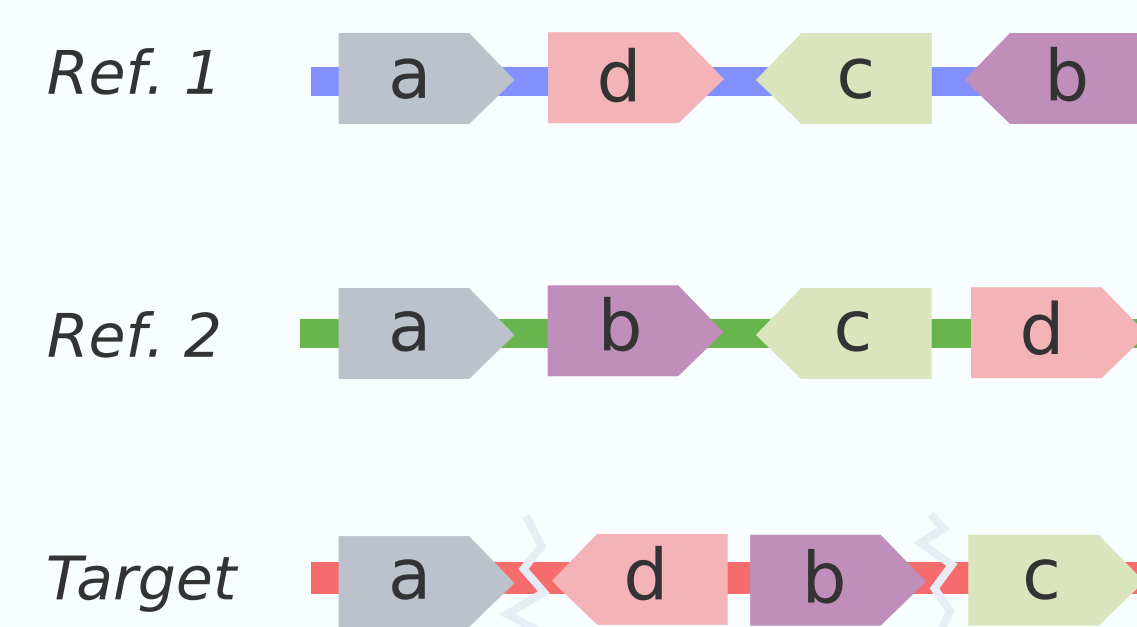
<sup>1</sup>University of California San Diego, USA <sup>2</sup>University of California Santa Cruz, USA, <sup>3</sup>University of Cambridge, UK, <sup>4</sup>Wellcome Trust Sanger Institute, UK, <sup>5</sup>European Bioinformatics Institute, UK, <sup>6</sup>The Genome Analysis Center, UK

## Abstract

- Assembly of mammalian-scale genomes into complete chromosomes is challenging
- To address this, we developed Ragout, a reference-assisted assembly tool for large and complex genomes
- Taking as input an NGS assembly and multiple related references, Ragout infers their evolutionary relationship and builds the final assembly of the target genome using a genome rearrangement approach
- Using Ragout we assembled two mice genomes (*M. Caroli* and *M. Pahari*) with complicated chromosome-scale rearrangements into sets of high-quality pseudo-chromosomes
- Chromosome coloring confirms most the rearrangements that Ragout has detected

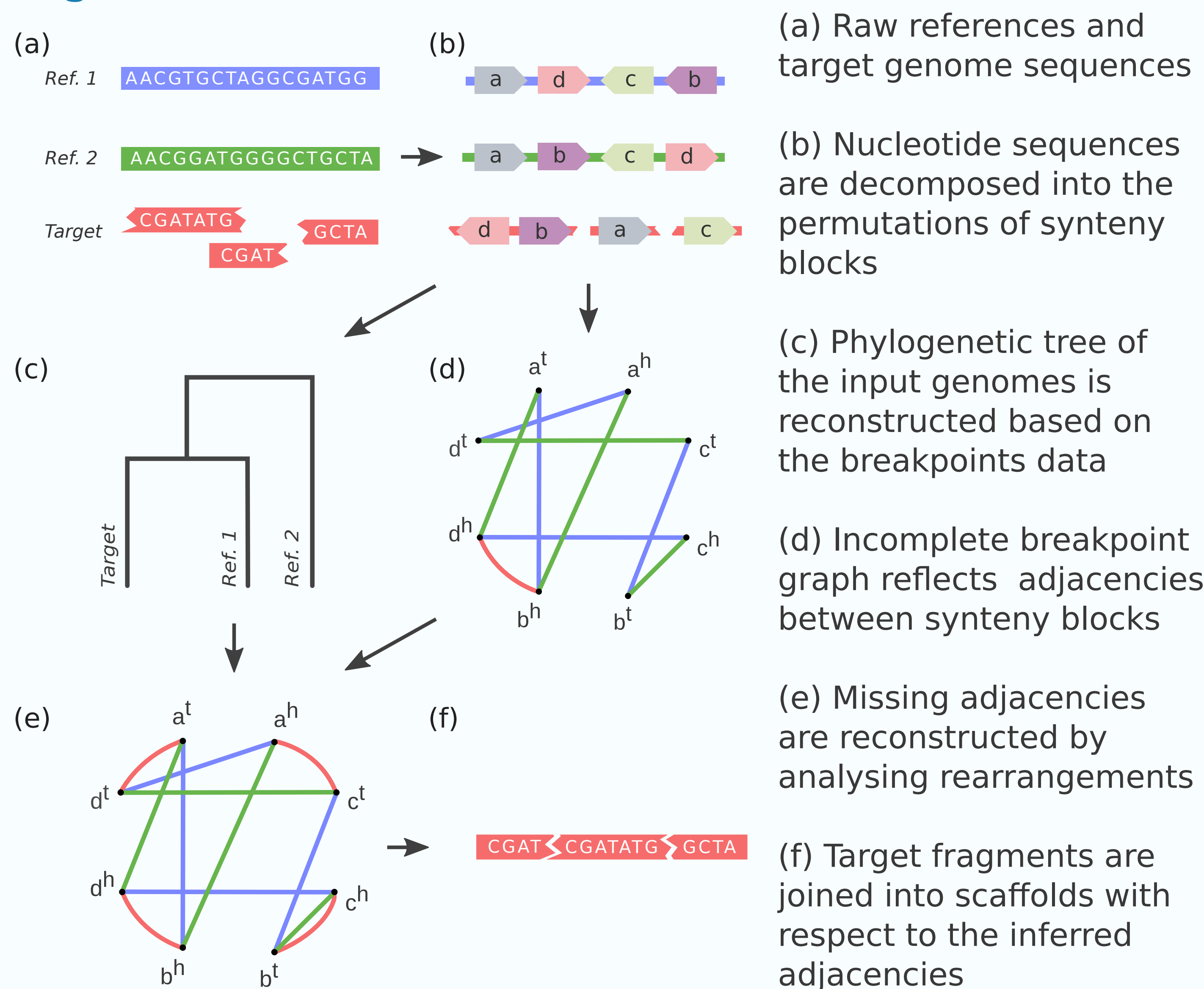


## Incomplete Breakpoint Graph Analysis



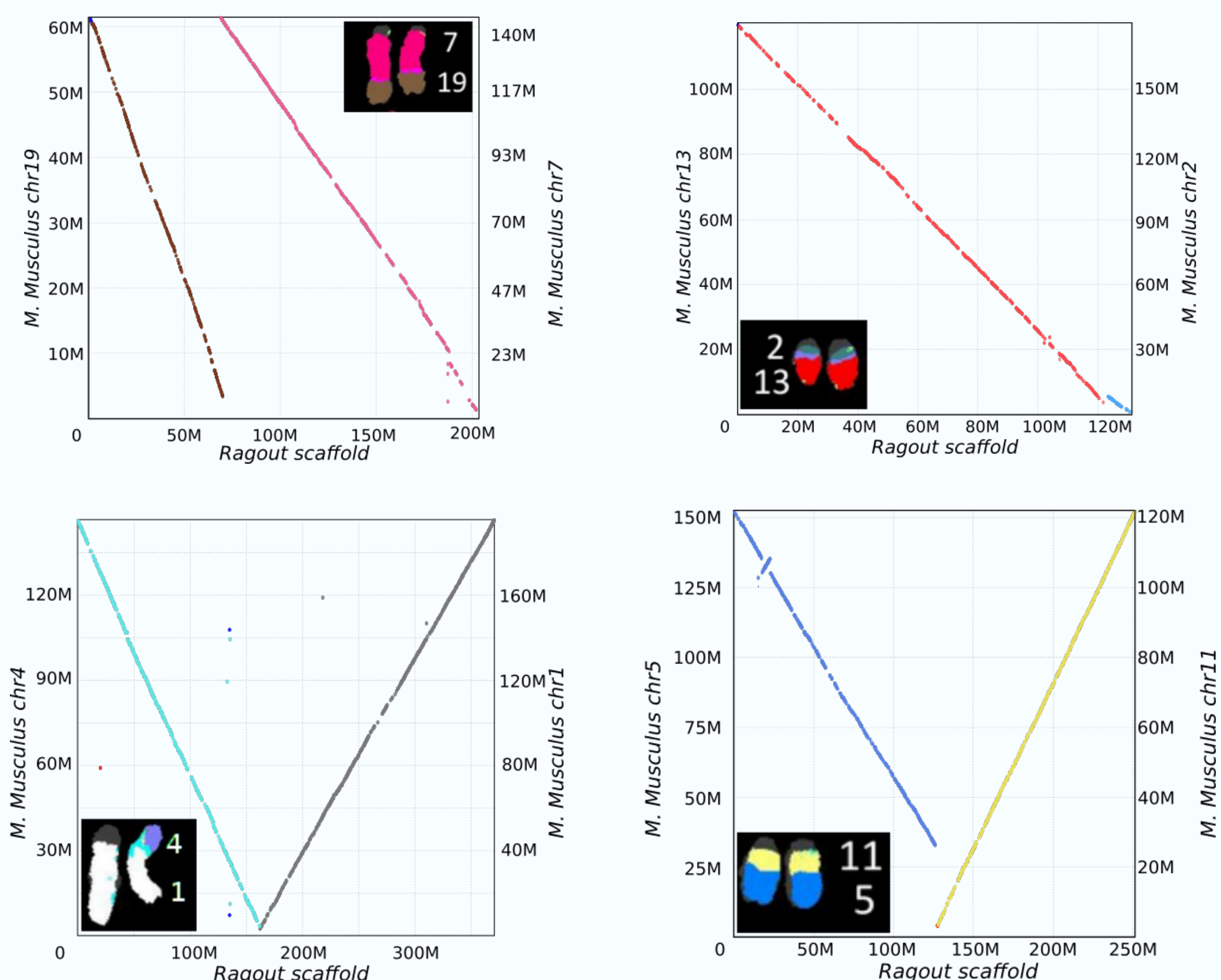
- Breakpoint graphs reflect adjacencies between syntenic blocks in different genomes
- If all genomes were complete, the edges of each color will define a perfect matching on the graph
- As the target genome is fragmented, some adjacencies of red color are missing
- Ragout recovers the missing adjacencies so as to minimize the weighted number of rearrangements between the genomes
- These adjacencies are then used to merge the target fragments into scaffolds

## Algorithm Overview

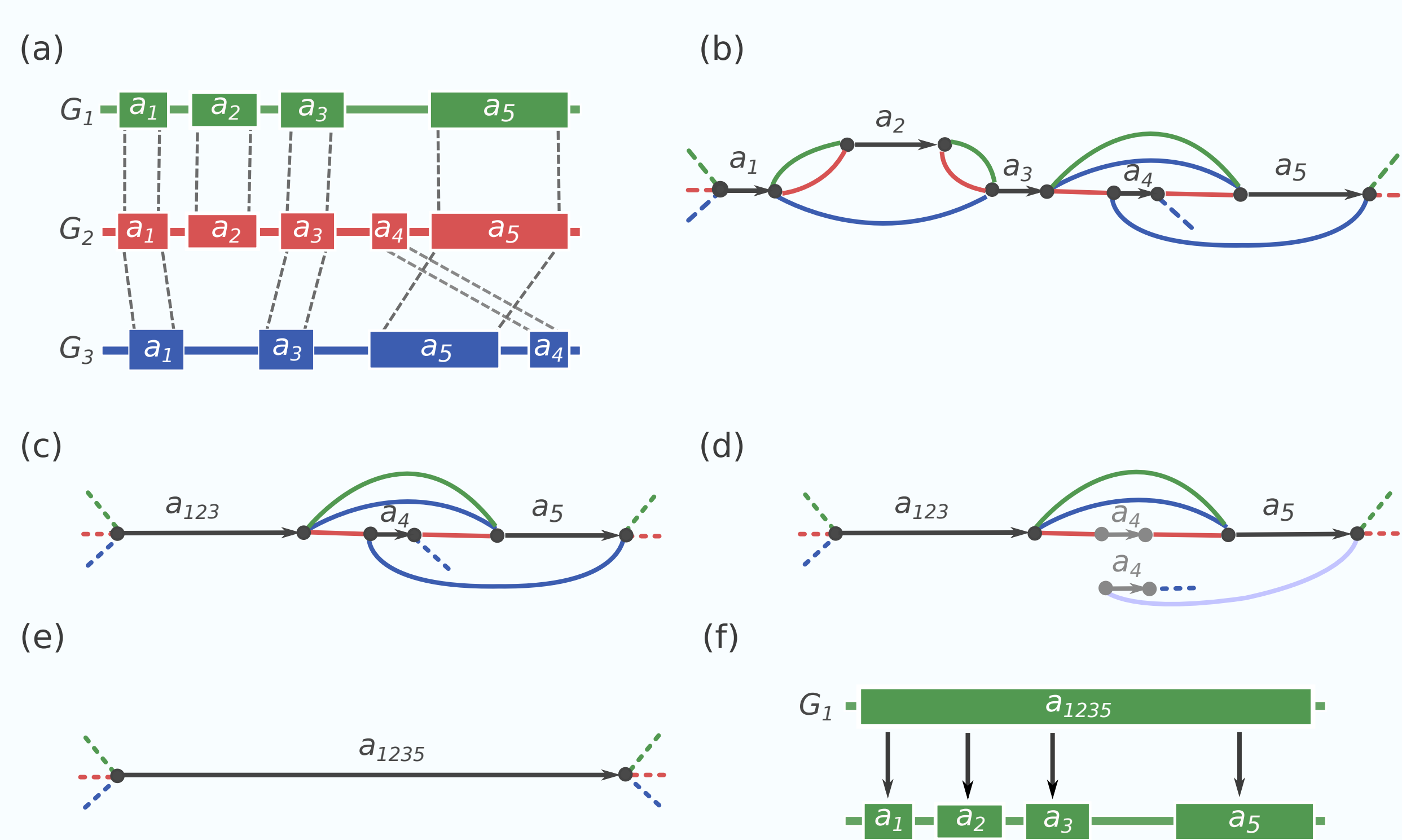


## Results

- We assembled two genomes from Murinae family: *Mus. Caroli* and *Mus. Pahari* using *Mus. Musculus* and *Rattus Norvegicus* as references. The assemblies contains 20 and 23 pseudo-chromosomes, respectively with at most 2% of unlocalized sequence.
- *M. Caroli* shows 5% sequence diversity from *Mus. Musculus* and has the same karyotype (which was confirmed by Ragout). However, we have detected a large inversion in chr17.
- *M. Pahari* has 10% sequence diversity and contain many inter-chromosomal rearrangements. Ragout has detected four of them, which were also confirmed by chromosome coloring.
- Some of the rearrangements remain undetected, as the corresponding breakpoints are missing from the NGS data. However, they could be recovered using the aid of chromosomal maps.



## Syntenic Blocks



Syntenic blocks help to separate small sequence variations from large-scale rearrangements. (a) An alignment between three genomes with complicated sub-structure. (b) A-Brujn graph representation of the alignment. Small sequence variations correspond to bubbles, while rearrangements form more complicated structures. (c) A bubble is removed during the graph simplification, forming a larger syntenic block  $a_{123}$  (d) Masking smaller block  $a_4$  allows to make the graph structure even simpler. (e) After removing another bubble, the whole alignment is represented as a large syntenic block. (f) A hierarchical representation of a syntenic block.

## Availability & Contacts

- Ragout is an easy to use package, written in Python/C+. It is freely available at <http://fenderglass.github.io/Ragout>
- Email: fenderglass@gmail.com

