# Aligning the unalignable: bacteriophage whole genome alignments

Sèverine Bérard, Annie Chateau, Nicolas Pompidor, Paul Guertin, Anne Bergeron and Krister M. Swenson
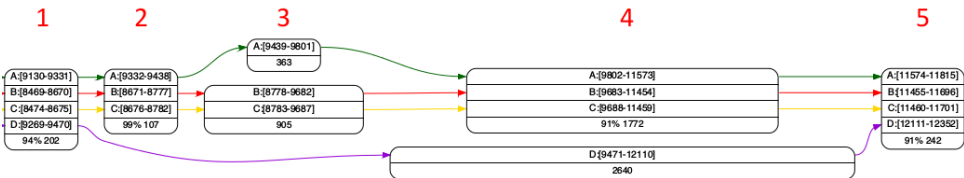
Colloque IPB 2015 Montpellier

# BACTERIOPHAGES GENOMES: A MOSAIC STRUCTURE

- Modular theory of phage genome organization [Botstein, 80]: biological functions are grouped into *modules* whose order is mostly conserved along the genomic sequence

- Each module has *variants* that perform the same functions, possibly encoded by dissimilar sequences

# CLASSICAL MULTIPLE ALIGNMENTS



Central dogma: sequence similarity $\Rightarrow$ functional similarity

# Phages multiple alignments



- The alignments rely on phages genomes features:
  - functional collinearity
  - a low duplication rate
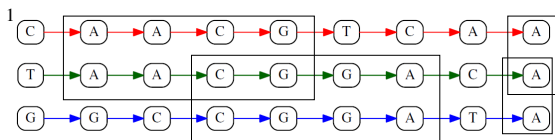  - the presence of long shared sequences

# THE ALPHA ALIGNER

Ideas

- to use exact matches (= anchor)

- to order them in a *graph* representing their succession in the genomes

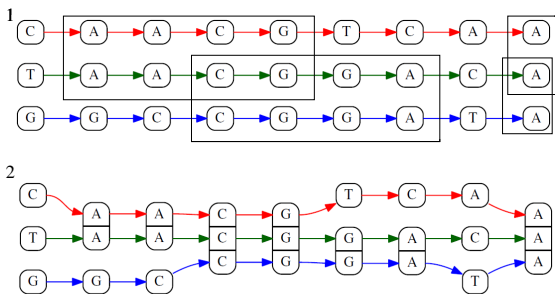- not to linearize the graph in an alignment
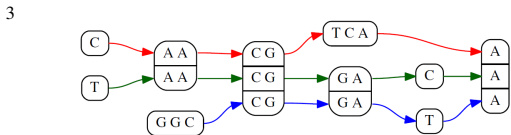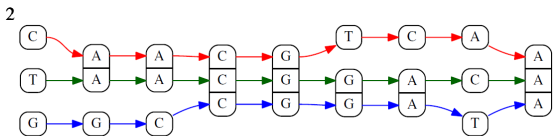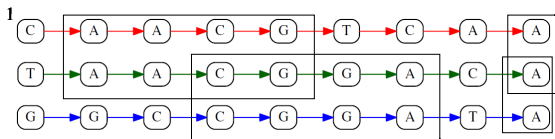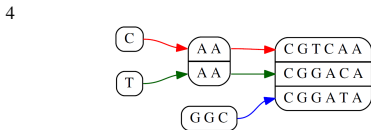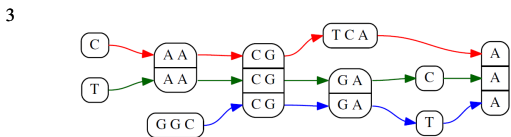
# THE ALPHA ALIGNER

Ideas

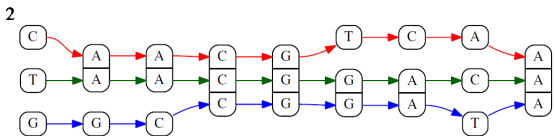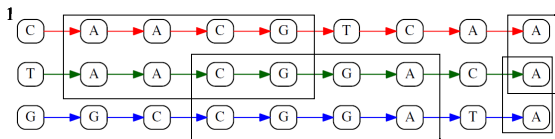- to use exact matches (= anchor)

- to order them in a *graph* representing their succession in the genomes

- not to linearize the graph in an alignment

$\Rightarrow$ we work on graphs, representing <span style="color:red">partial order</span>, to capture mosaicism of sets of genomes

# ALPHA'S CHARACTERISTICS

- One parameter $m$: minimal length of matches. Allow to adjust the alignment view

# Alpha's characteristics

- One parameter $m$: minimal length of matches. Allow to adjust the alignment view

- Time complexity in $\mathcal{O}(\alpha(M_l)M_l + nN_l)$
  (where $M_l$ is the total length of matches, $\alpha$ is the inverse of the Ackermann function, $n$ is the number of genomes and $N_l$ is the total length of genomes)
  = reasonable calculation time

# Alpha's characteristics

- One parameter $m$: minimal length of matches. Allow to adjust the alignment view

- Time complexity in $\mathcal{O}(\alpha(M_l)M_l + nN_l)$
  (where $M_l$ is the total length of matches, $\alpha$ is the inverse of the Ackermann function, $n$ is the number of genomes and $N_l$ is the total length of genomes)
  = reasonable calculation time

- Functional collinearity hypothesis: Alpha detects large rearrangements events that contradict this hypothesis and does not perform the alignment in such a case

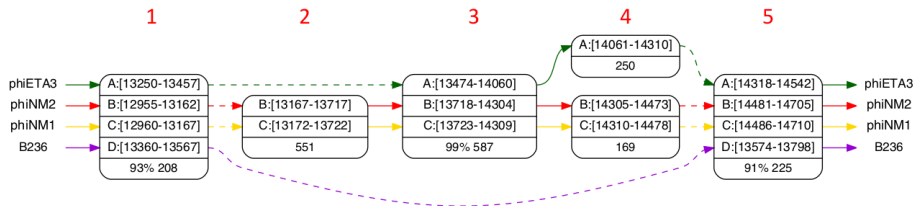# Highlighting rearrangements (*e.g.* deletions)



**Figure 4  A local alignment of four S. aureus bacteriophages.** In column 3, three phages are in a 587 bp exact alignment, with 99 % identical columns. A major deletion in phage B236 spans columns 2, 3, and 4, and the corresponding arrow is dotted to reflect the fact that some basepairs are not shown, 9 bp in this case.

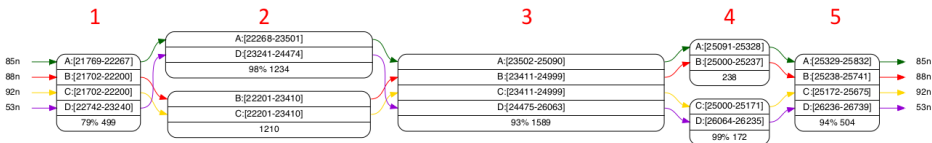# HIGHLIGHTING REARRANGEMENTS (*e.g.* RECOMBINATION)



**Figure 5 Modules and variants.** Local alignment of four *S. aureus* bacteriophages clearly showing modules and variants. Notice that sequences with the same variant in the second column are switched in the fourth.
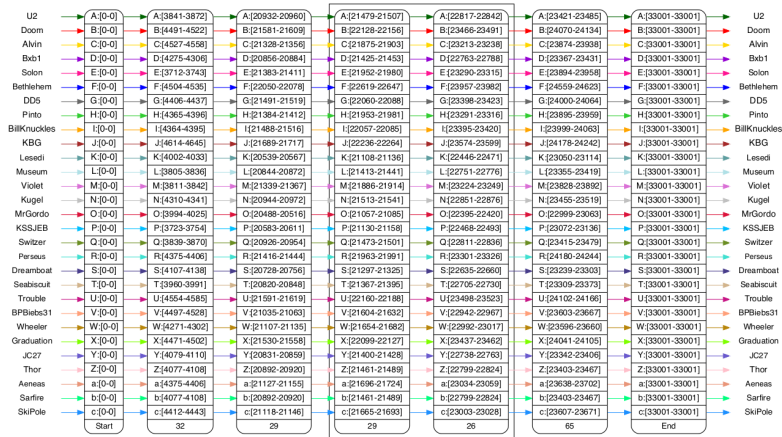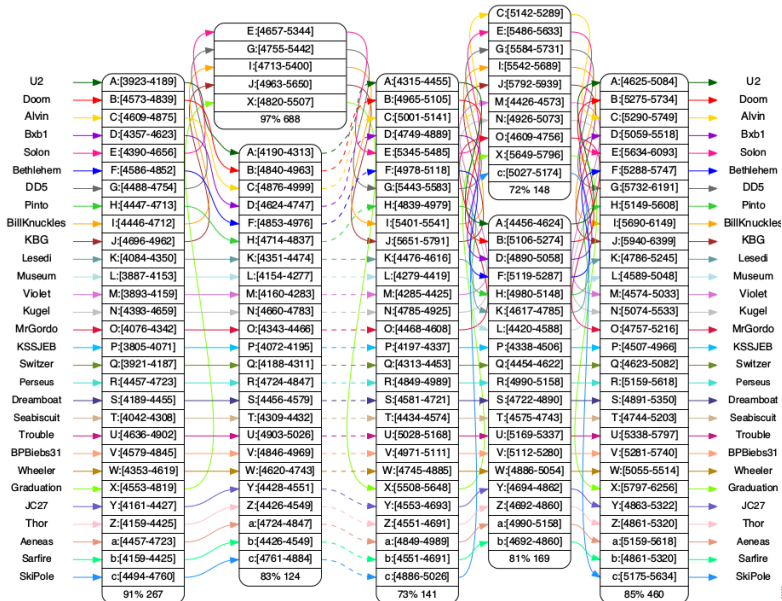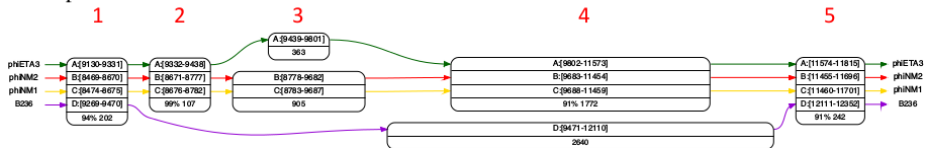
# Anchors visualization, backbone



**Figure 6 The anchor view.** Aligning a set of 29 *Mycobacterium* bacteriophages using $m = 175$: the aligner can display the sequence of anchors, or backbone, allowing the user to zoom quickly to a specific region. Boxes group pairs of anchors bounding a gapless alignment.

# Zooming from anchor view

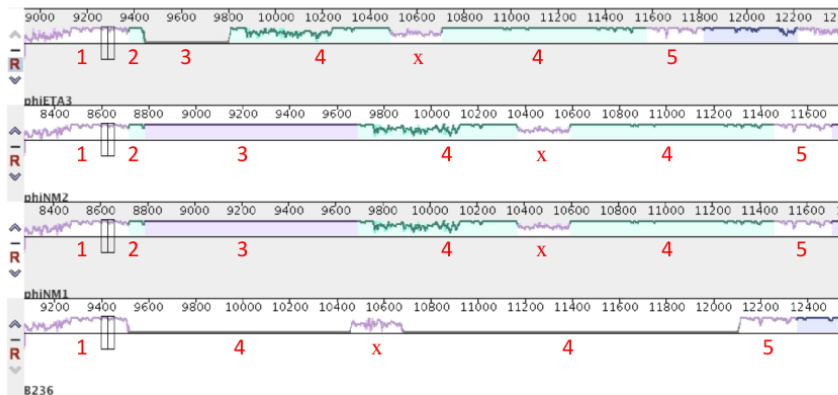# COMPARISON WITH MAUVE

# Comparison with Mauve



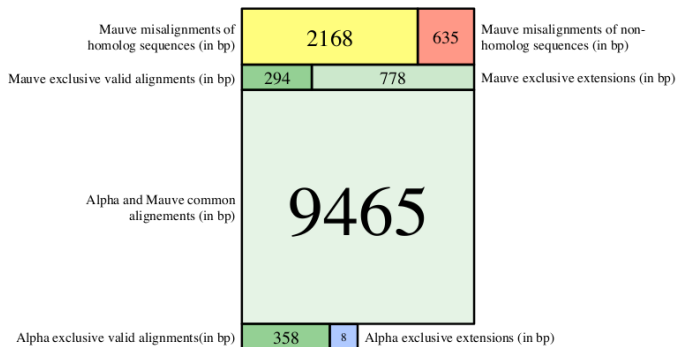**Figure 9 Alignments comparison.** Diagram representing the repartition of 13706 bp of 70 block alignments by Alpha or Mauve. The top 5 rectangles are from alignments produced by Mauve, and the bottom 3 are from Alpha alignments. Green and blue rectangles represent valid alignments, while yellow and red are misalignments.

# Conclusion and perspectives

- Conclusion

  - ▸ Alpha produces biologically meaningful alignments

# Conclusion and perspectives

- Conclusion

  - ▶ Alpha produces biologically meaningful alignments

  - ▶ Implemented in an interactive aligner

## Conclusion and perspectives

- Conclusion

  ▸ Alpha produces biologically meaningful alignments

  ▸ Implemented in an interactive aligner

  ▸ The model is almost parameter free

# Conclusion and perspectives

- Conclusion

  ▶ Alpha produces biologically meaningful alignments

  ▶ Implemented in an interactive aligner

  ▶ The model is almost parameter free

- Perspectives

  ▶ Transfer of annotations

# Conclusion and perspectives

- Conclusion

  - Alpha produces biologically meaningful alignments

  - Implemented in an interactive aligner

  - The model is almost parameter free

- Perspectives

  - Transfer of annotations

  - Extend our mathematical model to deal with rearrangements

## Thank you for your attention

- References:
  - ▶ Article: "Aligning the unalignable: bacteriophage whole genome alignments", *BMC Bioinformatics*, to appear

  - ▶ Tool: Alpha
    http://www.lirmm.fr/∼swenson/alpha/alpha.htm

# Any questions ?
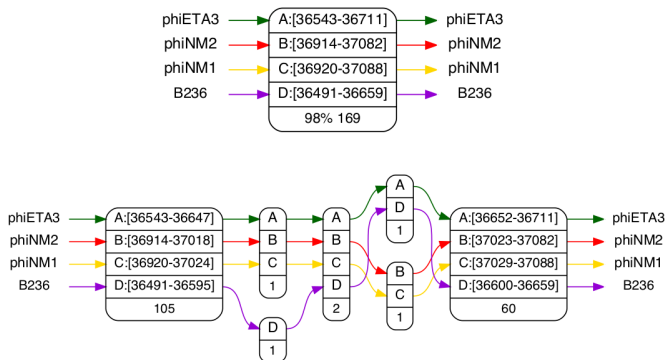
# Alpha produces gapless alignment



**Figure 3 Contracted vertices.** A example of a contracted vertex (top) and its expanded graph (bottom) in the alignment of four *S. aureus* bacteriophages. The length of the vertice is 169 nucleotides, and the percentage of identity is 98 %. In the expanded graph, there are two single nucleotide mutations that account for the 98% identity.

## Quality of Alpha gapless alignments

- Comparison with 3 other alignment tools : Clustal, MUSCLE and T-Coffee

- 491 alignments compared over 3 datasets

- Only 8 regions of disagreement. For 7 regions, Alpha alignments are validated by translating in a.a. (tblastx). For the last region, the right solution is unknown.