

Towards an Automated Annotation of CRISPR-cas systems

Costa F¹, Alkhnbashi OS¹, Shah SA², Saunders SJ¹, Barrangou R⁴, Brouns SJJ⁵, Charpentier E⁶, Haft DH³, Horvath P⁷, Koonin EV³, Makarova KS³, Moineau S⁸, Mojica FJM⁹, Terns RM¹⁰, Terns MA¹⁰, van der Oost J⁵, White MF¹¹, Wolf YI³, Yakunin AF¹², Garrett RA², Backofen R^{1,13}



- Bioinformatics group, Department of Computer Science, University of Freiburg, Georges-Koehler-Allee 106, 79110 Freiburg, Germany
- Archaea Centre, Department of Biology, Copenhagen University, Ole Maaløes Vej 5, DK2200 Copenhagen N, Denmark
- National Center for Biotechnology Information, NLM, National Institutes of Health, Bethesda, MD 20894, USA
- Department of Food, Bioprocessing, and Nutrition Sciences, North Carolina State University, Raleigh, NC 27606, USA
- Laboratory of Microbiology, Wageningen University, Dreijenplein 10, 6703HB Wageningen, 16 Netherlands
- Department of Regulation in Infection Biology, Helmholtz Centre for Infection Research, D-38124 Braunschweig, Germany
- DuPont Nutrition and Health, BP10, Dangé-Saint-Romain 86220, France
- Département de biochimie, de microbiologie et de bio-informatique, Faculté des sciences et de génie, Groupe de recherche en écologie buccale, Félix d'Hérelle Reference Center for Bacterial Viruses, Faculté de médecine dentaire, Université Laval, Québec City, Québec, Canada
- Departamento de Fisiología, Genética y Microbiología, Universidad de Alicante. 03080-Alicante Spain
- Biochemistry and Molecular Biology, Genetics and Microbiology, University of Georgia, Davison Life Sciences Complex, Green Street, Athens, GA 30602
- Biomedical Sciences Research Complex, University of St Andrews, North Haugh, St Andrews, KY16 9TZ, UK
- Department of Chemical Engineering and Applied Chemistry, University of Toronto, Toronto, M5S 3E5, Canada
- BIOSS Centre for Biological Signaling Studies, Cluster of Excellence, University of Freiburg, Germany.

Automated annotation of CRISPR-cas loci in newly sequenced genomes is valuable both for general purposes of comparative genomics of archaea and bacteria, and for the progress of CRISPR research.

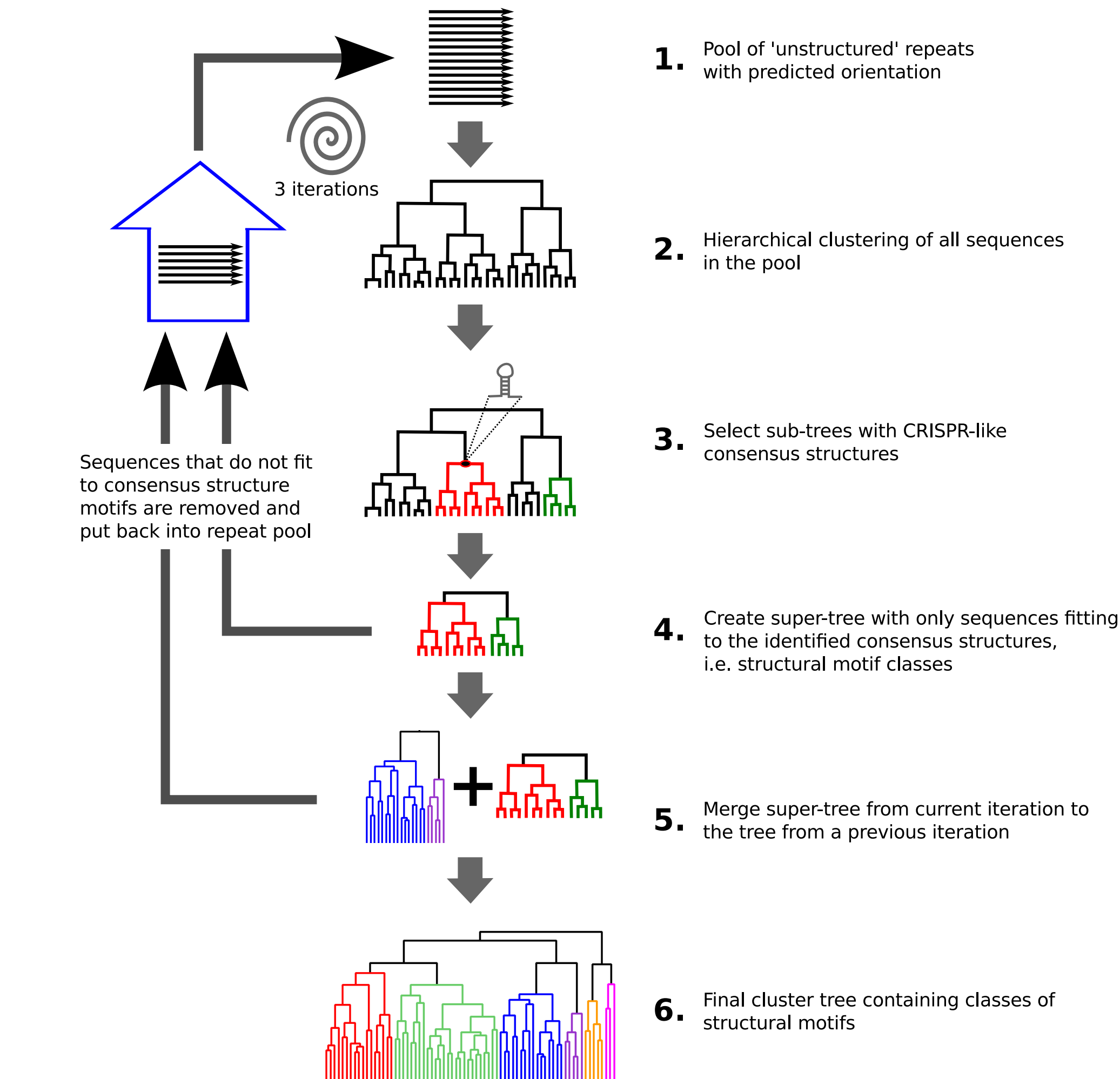
For a comprehensive analysis, several tasks have to be performed, among which the most important are:

- 1) the correct prediction of **repeat orientation** [CRISPRstrand]
- 2) a characterization of the repeats in terms of sequence and structure to infer **repeat evolution** [CRISPRmap]
- 3) annotation of the associated **loci subtype** according to the composition of cas genes.

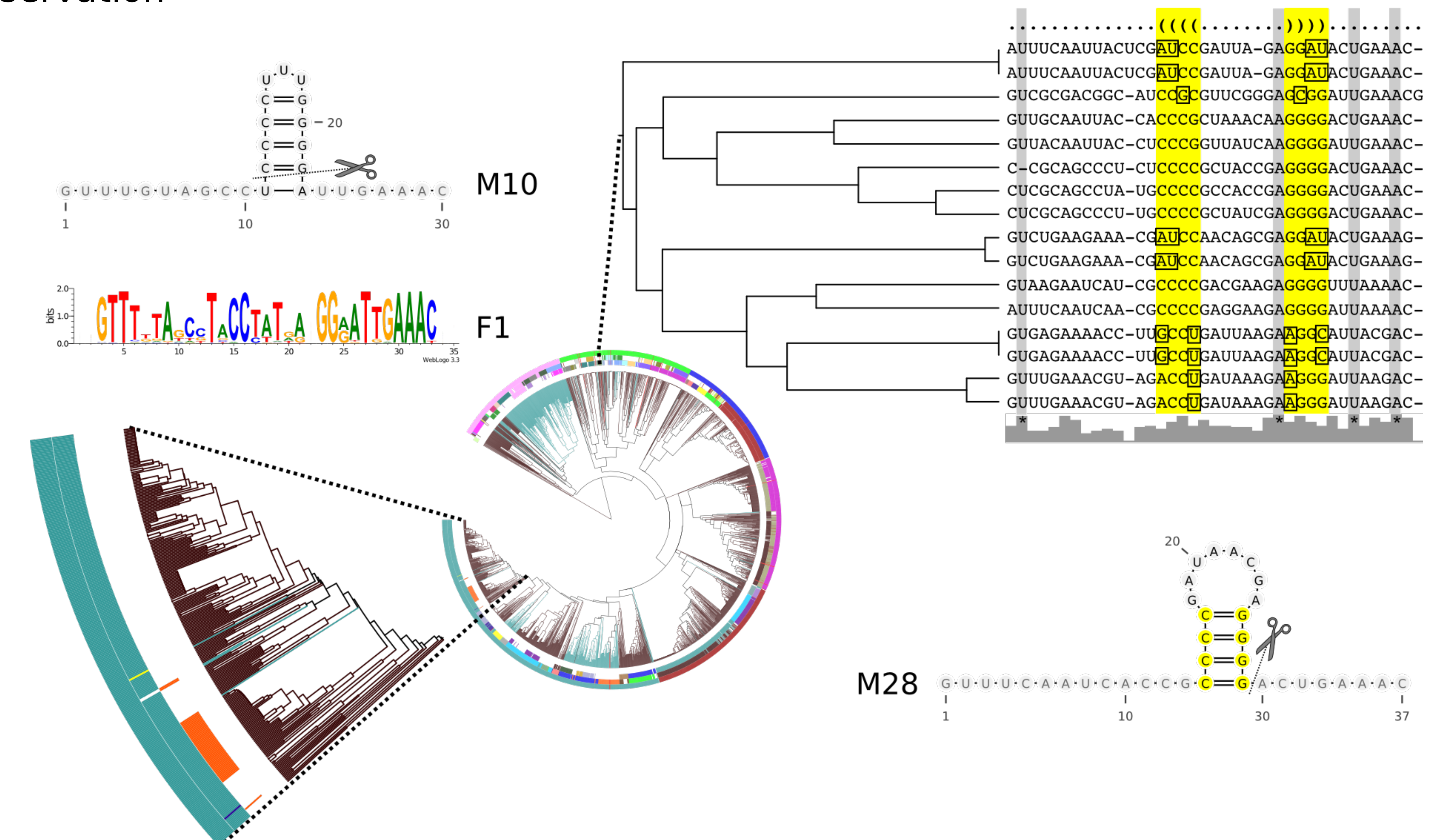
For the annotation of loci subtypes we present a solution that relies on a *novel similarity notion* for the interference modules. Loci annotation is achieved by nearest neighbor classification, which yields highly consistent results with respect to the current subtype classification.

CRISPRmap: an automated classification of repeat conservation

Basic procedure for identifying conserved hairpin motifs in all repeats.



We predicted CRISPR loci in all publically available genome sequences using the CRT tool and CRISPRfinder. In total, we analysed 4500 CRISPR loci in **4899 genomes**: 4590 bacteria, 309 archaea

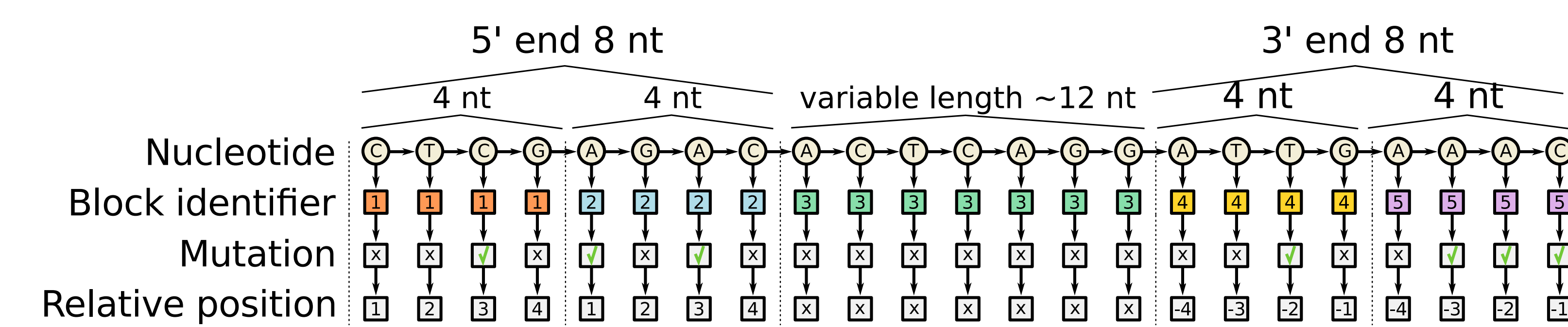


Highlighting the advantage of independent clustering approaches. (A) CRISPRs in the largest sequence family, F1, are mostly unstructured; however, for 50 CRISPRs also a conserved structure motif, M10, was identified. This indicates that subsets of conserved families can be structured. F1 contains the conserved 50 tag, marked with the magenta box. (B) Structure motif M28 shows no sequence conservation, but a conserved structure (base pairs are highlighted in yellow). The many compensatory base pairs are marked in the alignment with squares. This structure has been verified via mutational analyses. Potential cleavage sites are indicated as observed in the literature

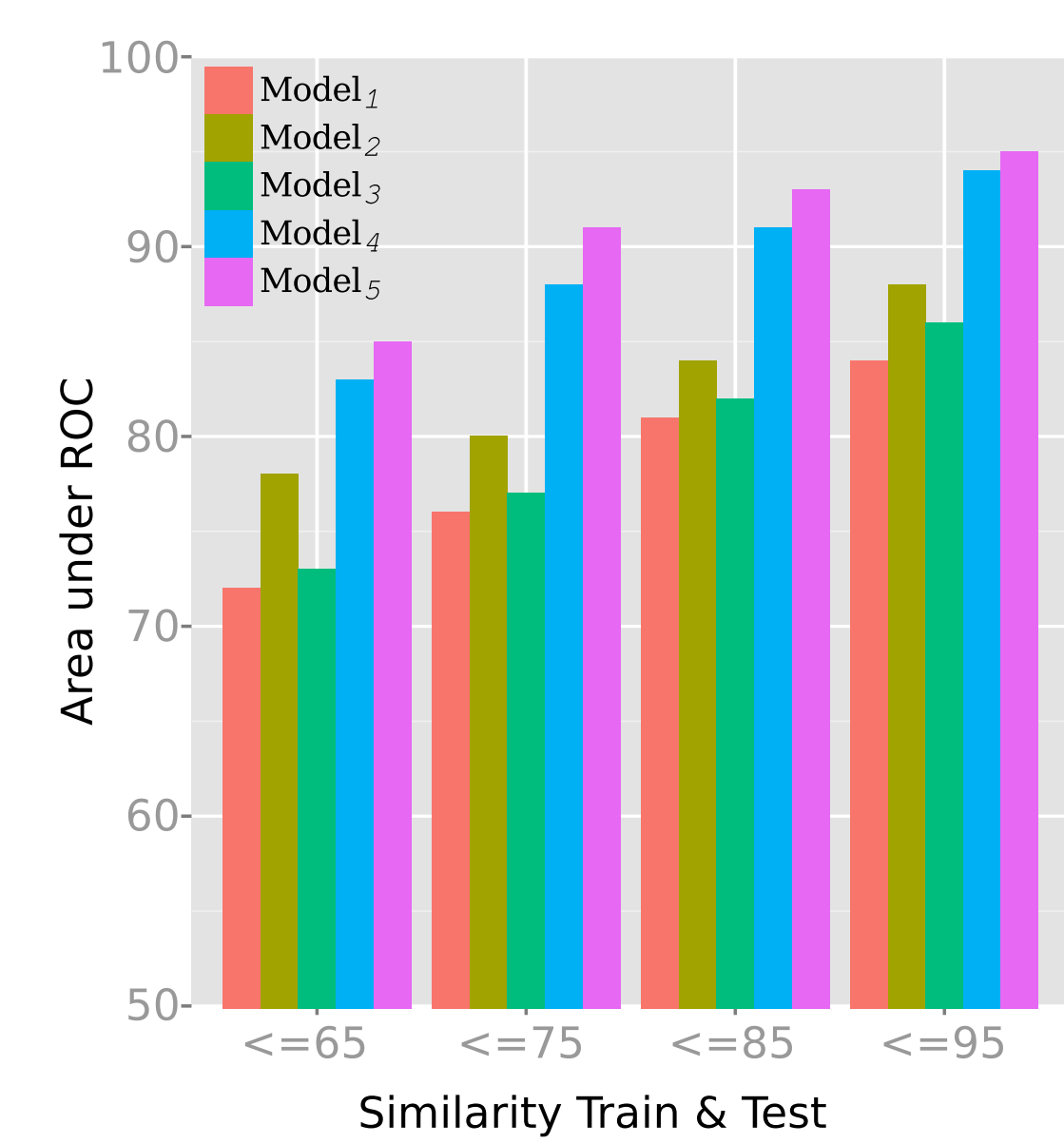
CRISPRstrand: predicting repeat orientations at CRISPR loci

Although existing bioinformatics tools can recognize CRISPR loci by their characteristic repeat-spacer architecture, they generally output CRISPR arrays of ambiguous orientation and thus do not determine the strand from which crRNAs are processed.

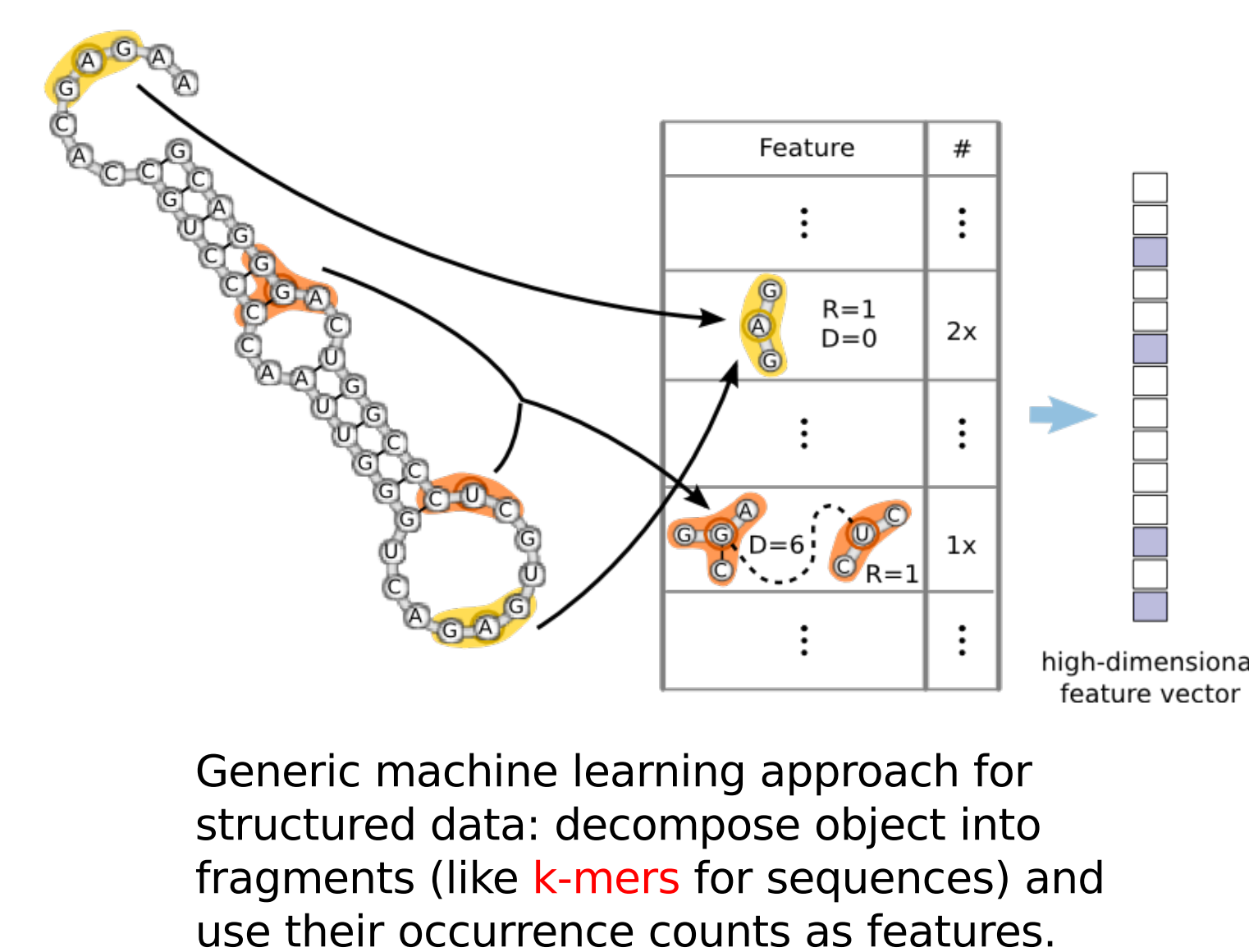
Knowledge of the **correct orientation is crucial** for many tasks, including the classification of CRISPR conservation, the detection of leader regions, the identification of target sites (protospacers) on invading genetic elements and the characterization of protospacer-adjacent motifs.



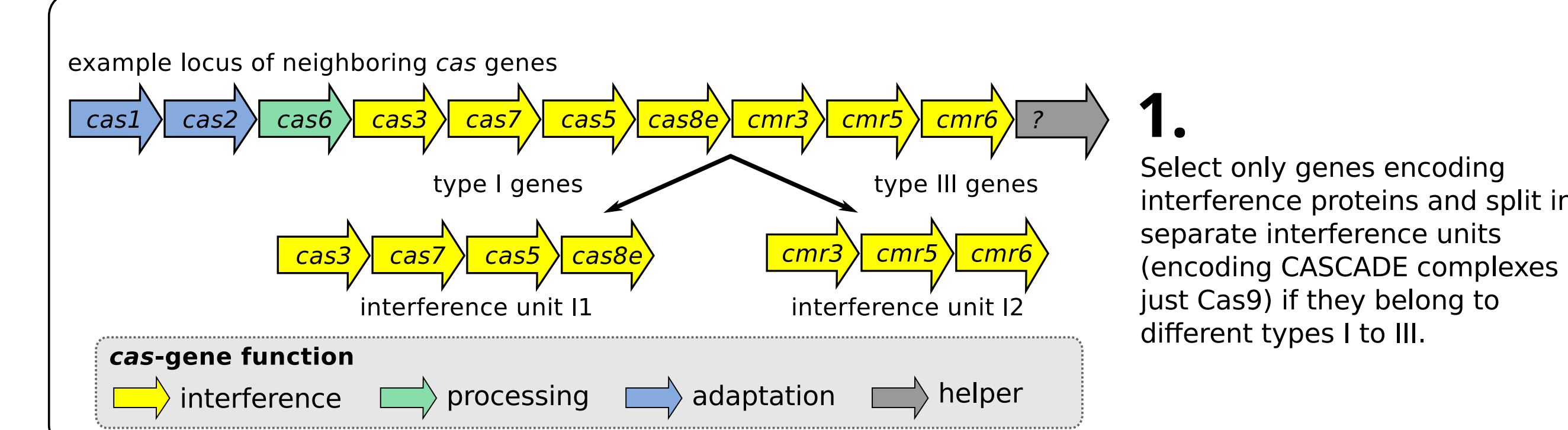
Graph encoding the consensus repeat sequence. The consensus nucleotide information is represented as a path graph, and additional information is modelled as a chain of additional vertices: block identifiers, mutation and relative position in the block.



AUC ROC performance comparison of the five models that encode increasing amount of information about the CRISPR arrays.



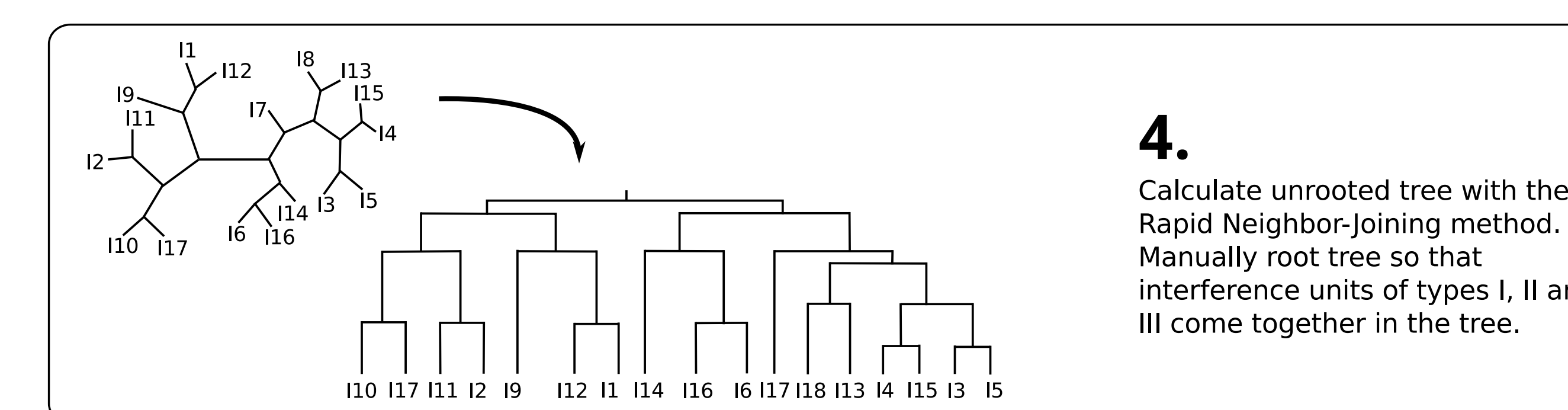
Clustering of CRISPR-Cas crRNA-effector modules by protein sequence similarity



$$symbit(a,b) = \frac{bitscore(a,b) + bitscore(b,a)}{2}$$

$$sim(a,b) = \frac{symbit(a,b)}{\sqrt{symbit(a,a) \times symbit(b,b)}}$$

$$s(Ix, Iy) = \frac{1}{|Ix||Iy|} \sum_{a \in Ix} \sum_{b \in Iy} sim(a,b)$$



Given the variability of the gene composition and architecture of the effector modules and the lack of a universal marker suitable for phylogenetic analysis, we developed a simple clustering approach based on **sequence similarity between the proteins** of these modules.

Cas loci can contain multiple modules of different types, we therefore separated them accordingly. Gene-to-gene similarities were calculated and aggregated into a locus-wide value. The resulting cluster dendrogram showed a high correlation with the subtype classification.

Using prior information on the association between sequence profiles and CRISPR-cas loci and the corresponding classification of the effector modules, we developed a **k-nearest neighbor classifier** that achieved **0.998 accuracy**.

