



# Simka: fast kmer-based method for estimating the similarity between numerous metagenomic datasets

Gaëtan Benoit, Pierre Peterlongo, Dominique Lavenier, Claire Lemaitre

## ► To cite this version:

Gaëtan Benoit, Pierre Peterlongo, Dominique Lavenier, Claire Lemaitre. Simka: fast kmer-based method for estimating the similarity between numerous metagenomic datasets. JOBIM 2015, Jul 2015, Clermont-Ferrand, France. , 10.1093/Bioinforma-cs/btu406 . hal-01180603

HAL Id: hal-01180603

<https://inria.hal.science/hal-01180603>

Submitted on 27 Jul 2015

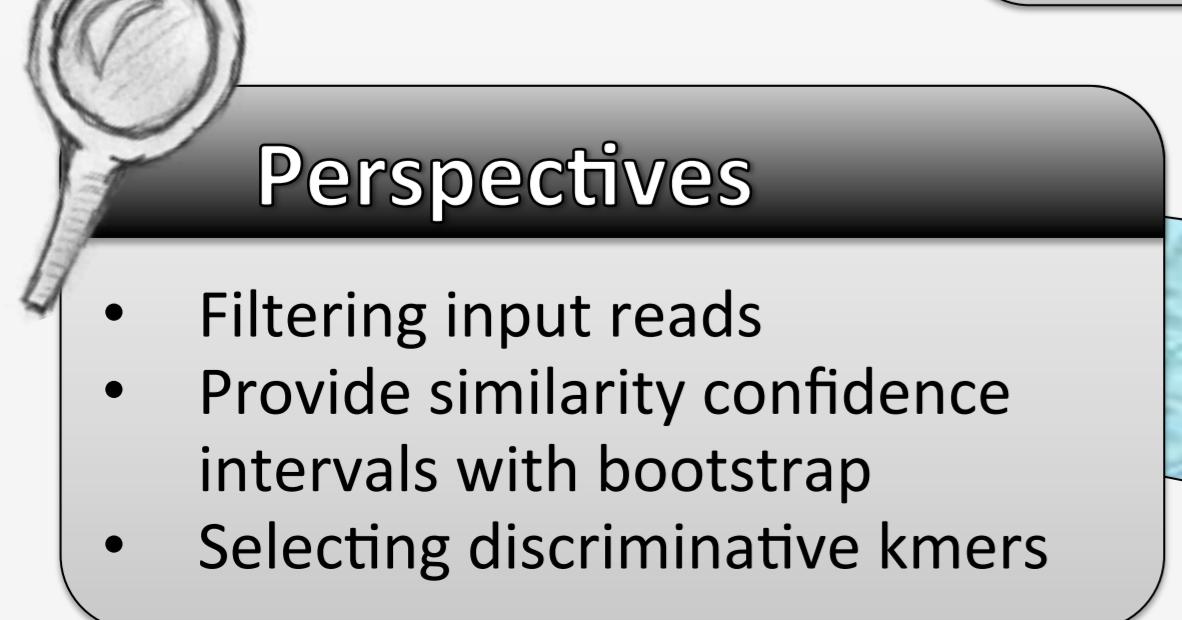
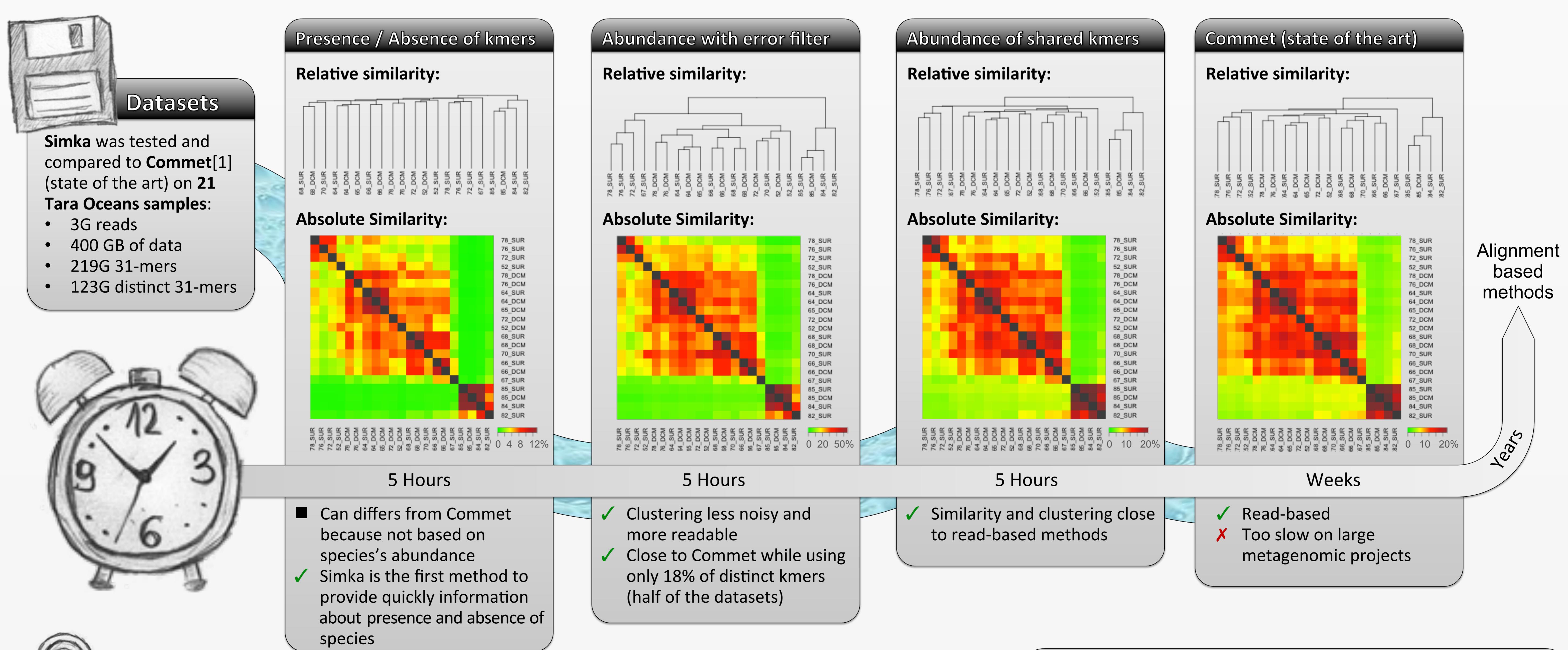
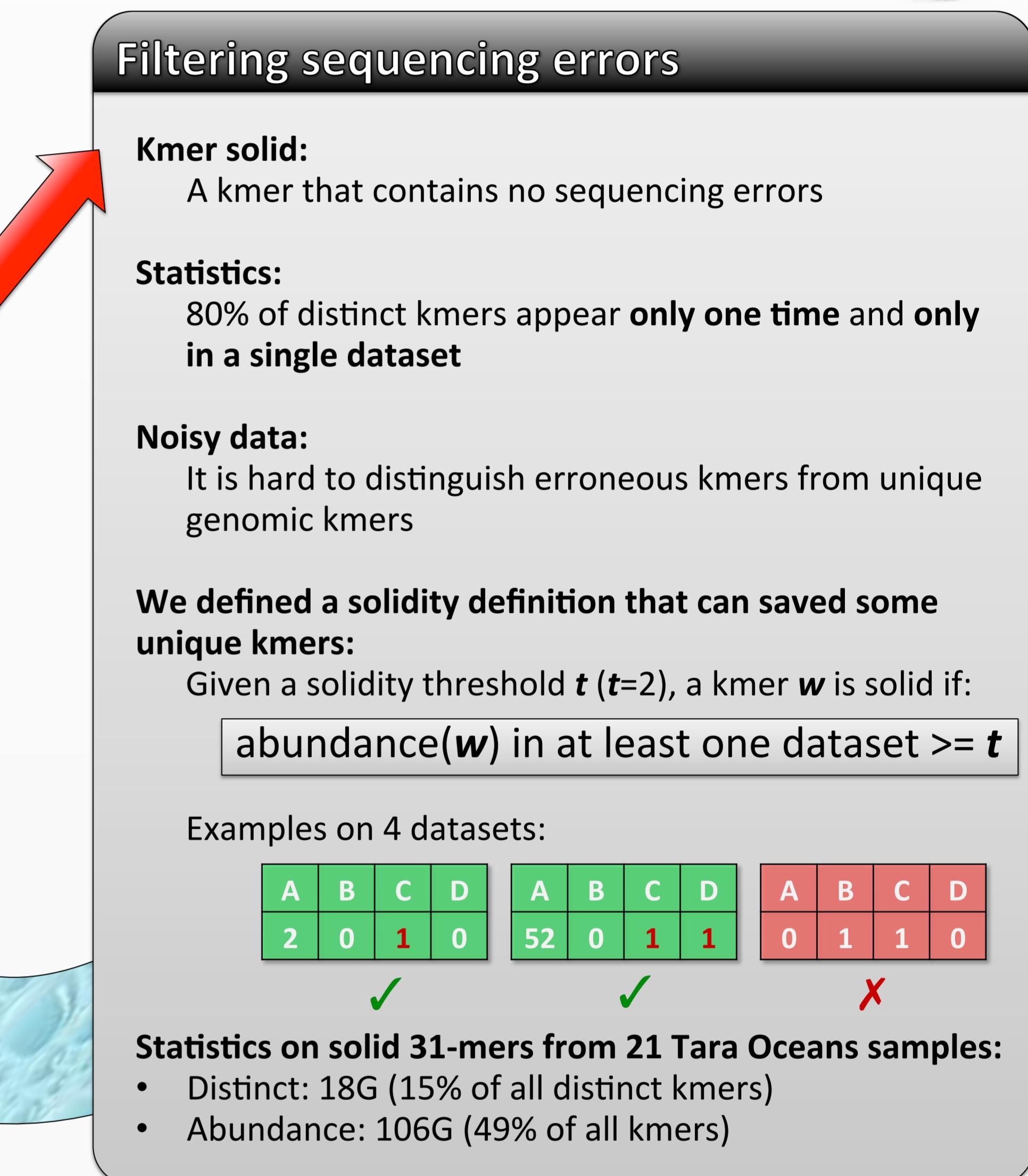
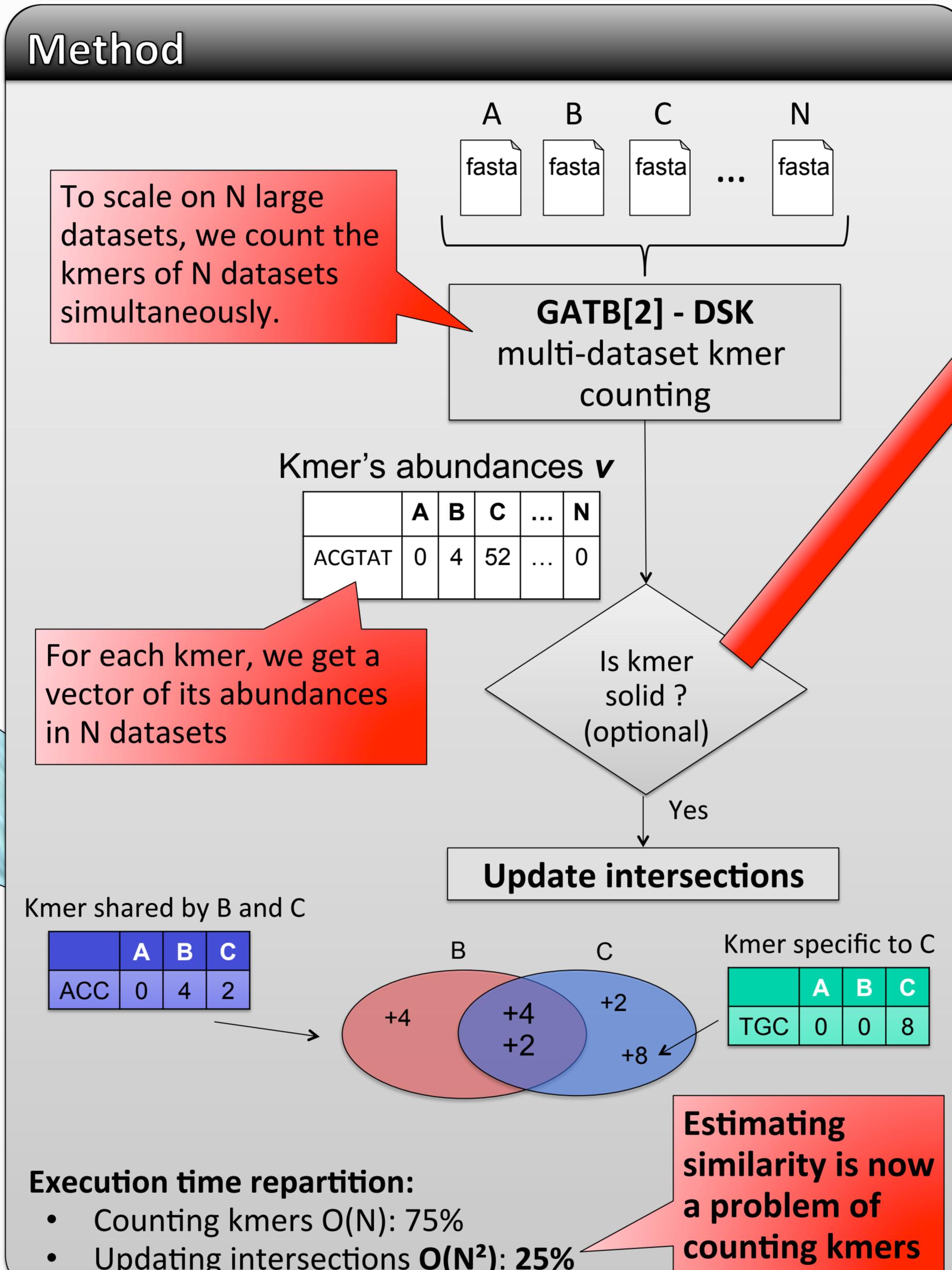
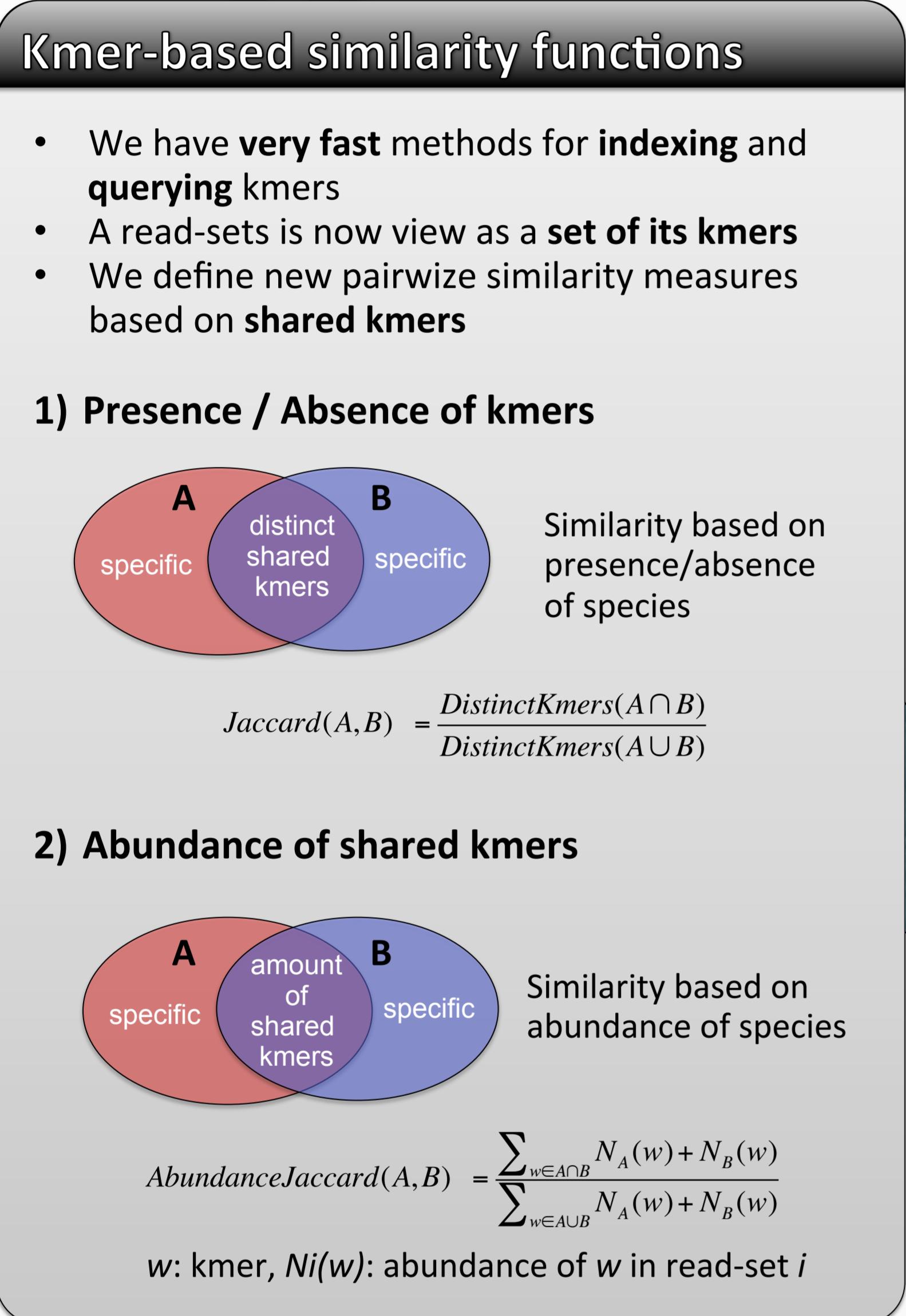
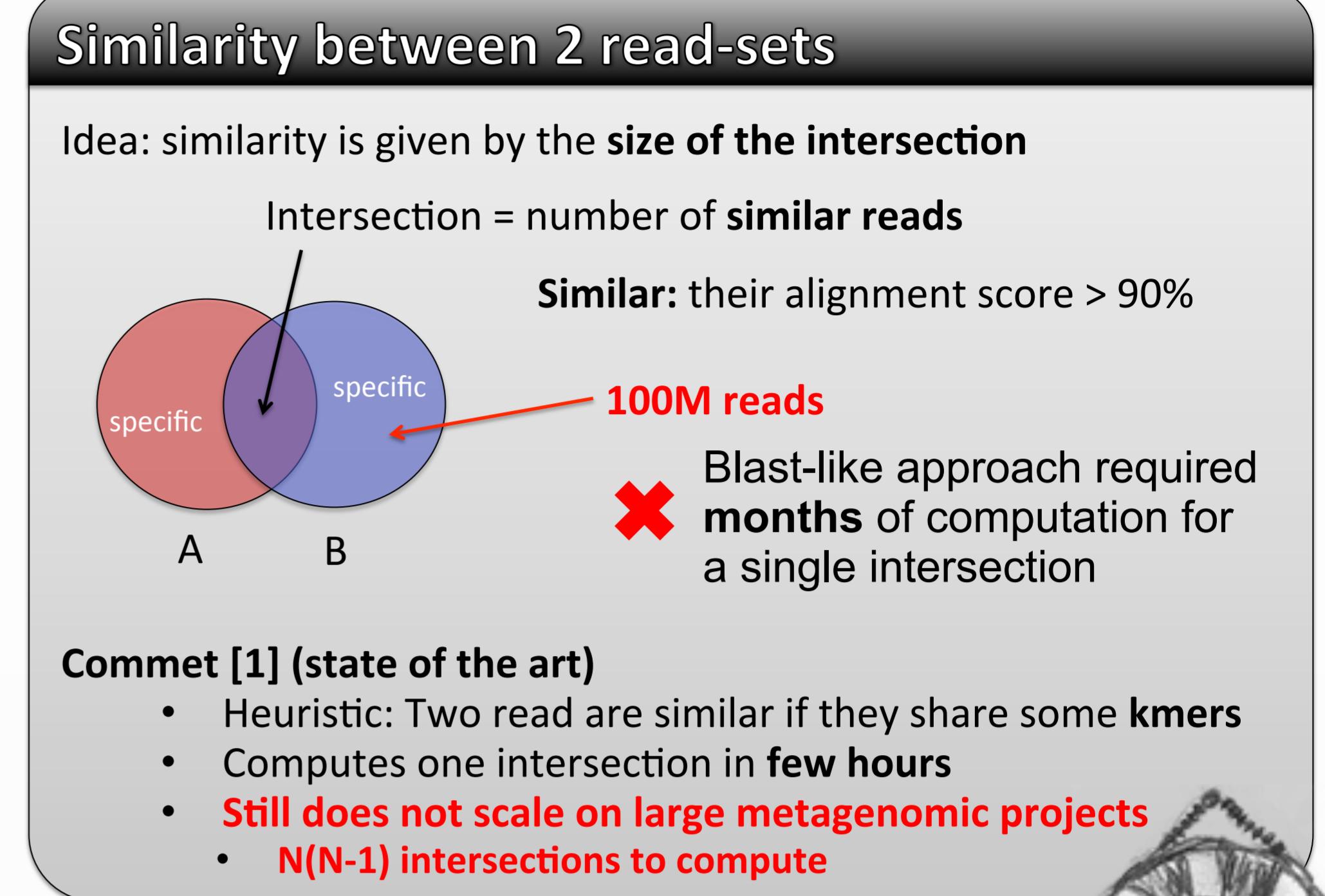
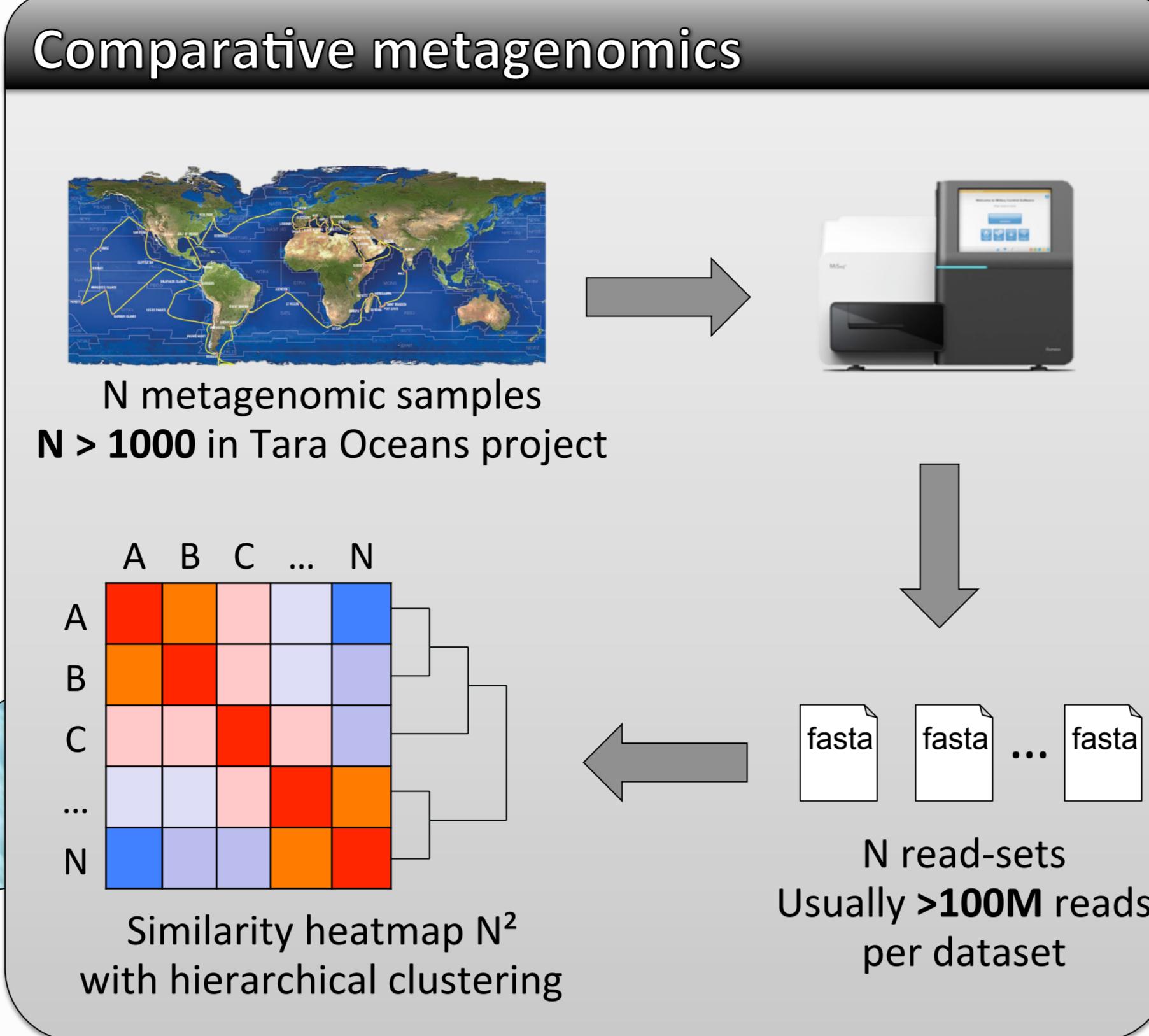
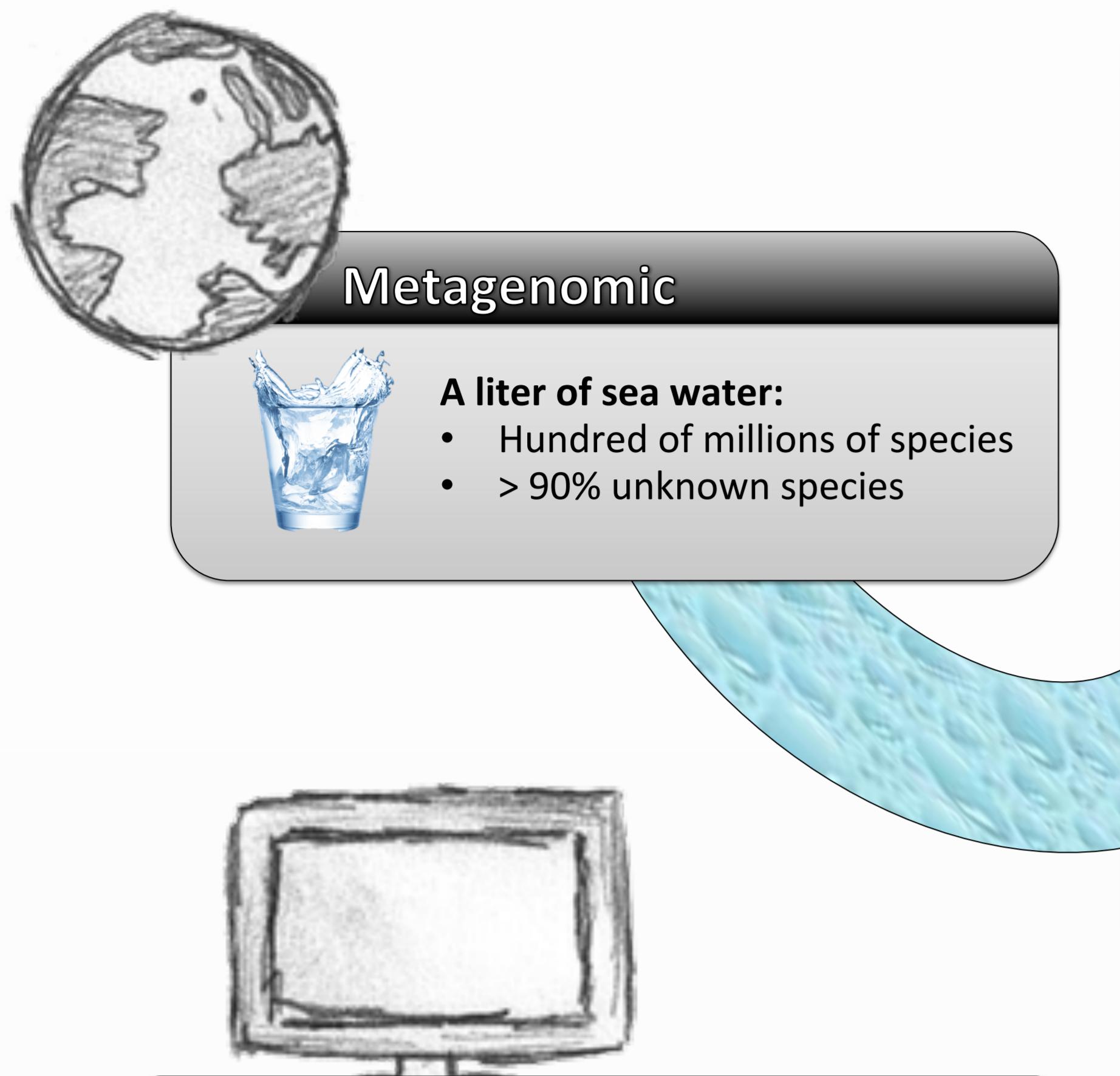
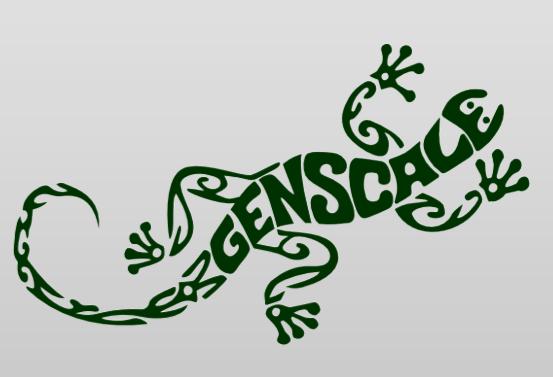
**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Simka: fast kmer-based method for estimating the similarity between numerous metagenomic datasets

Gaëtan Benoit<sup>1</sup>, Pierre Peterlongo<sup>1</sup>, Dominique Lavenier<sup>1</sup>, Claire Lemaitre<sup>1</sup>

<sup>1</sup> Inria/IRISA GenScale, Campus de Beaulieu, 35042 Rennes cedex, France.  
gaetan.benoit@inria.fr, claire.lemaitre@inria.fr, pierre.peterlongo@inria.fr



## Simka

- ★ New similarity functions based on **shared kmers**
- ★ Based on **abundance** and **presence/absence** of species
- ★ Results close to read-based methods
- ★ Fast and low memory thanks to the **GATB library [2]**

## References

- [1] COMMET: comparing and combining multiple metagenomic datasets  
N. Maillet, G. Collet, T. Vannier, D. Lavenier, P. Peterlongo  
IEEE BIBM, 2014
- [2] GATB: Genome Assembly & Analysis Tool Box  
E. Drezen, G. Rizk, R. Chikhi, C. Delteil, C. Lemaitre, P. Peterlongo, D. Lavenier  
10.1093/Bioinformatics/btu406, 2014  
<https://gatb.inria.fr/>

Funding: ANR Hydrogen, ANR-14-CE23-0001