

Evaluating Statistical Multiple Sequence Alignment in Comparison to Other Alignment Methods on Protein Data Sets

Michael Nute , Ehsan Saleh , Tandy Warnow 

Systematic Biology, Volume 68, Issue 3, May 2019, Pages 396–411, <https://doi.org/10.1093/sysbio/syy068>

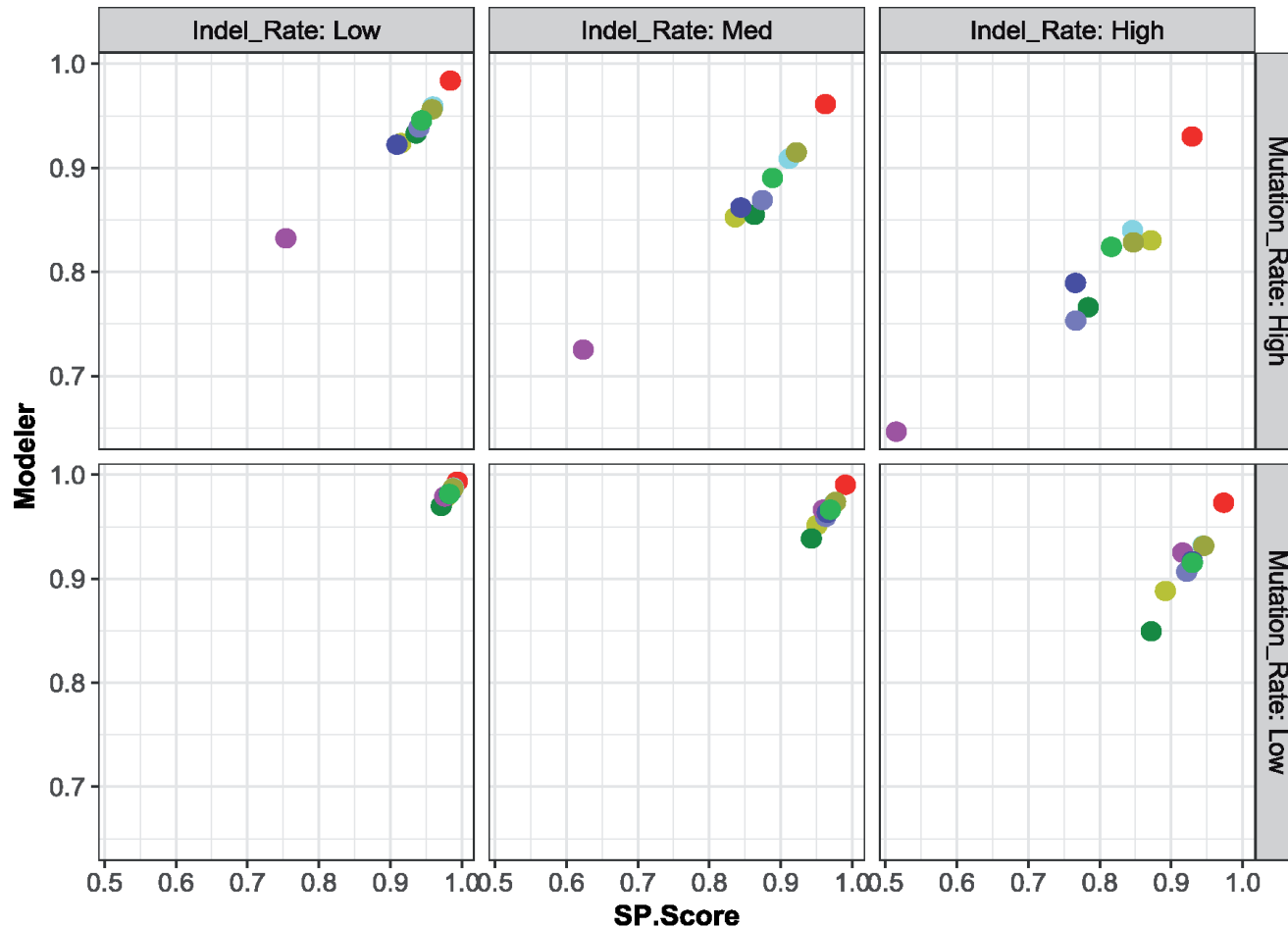
Published: 17 October 2018 **Article history** ▼

Compared BALi-Phy to other MSA methods on both 1192 protein biological benchmarks and 120 simulated datasets

Bali-Phy was run for 32 independent runs, each for 48 hours, to enable it to converge

All datasets at most 27 sequences

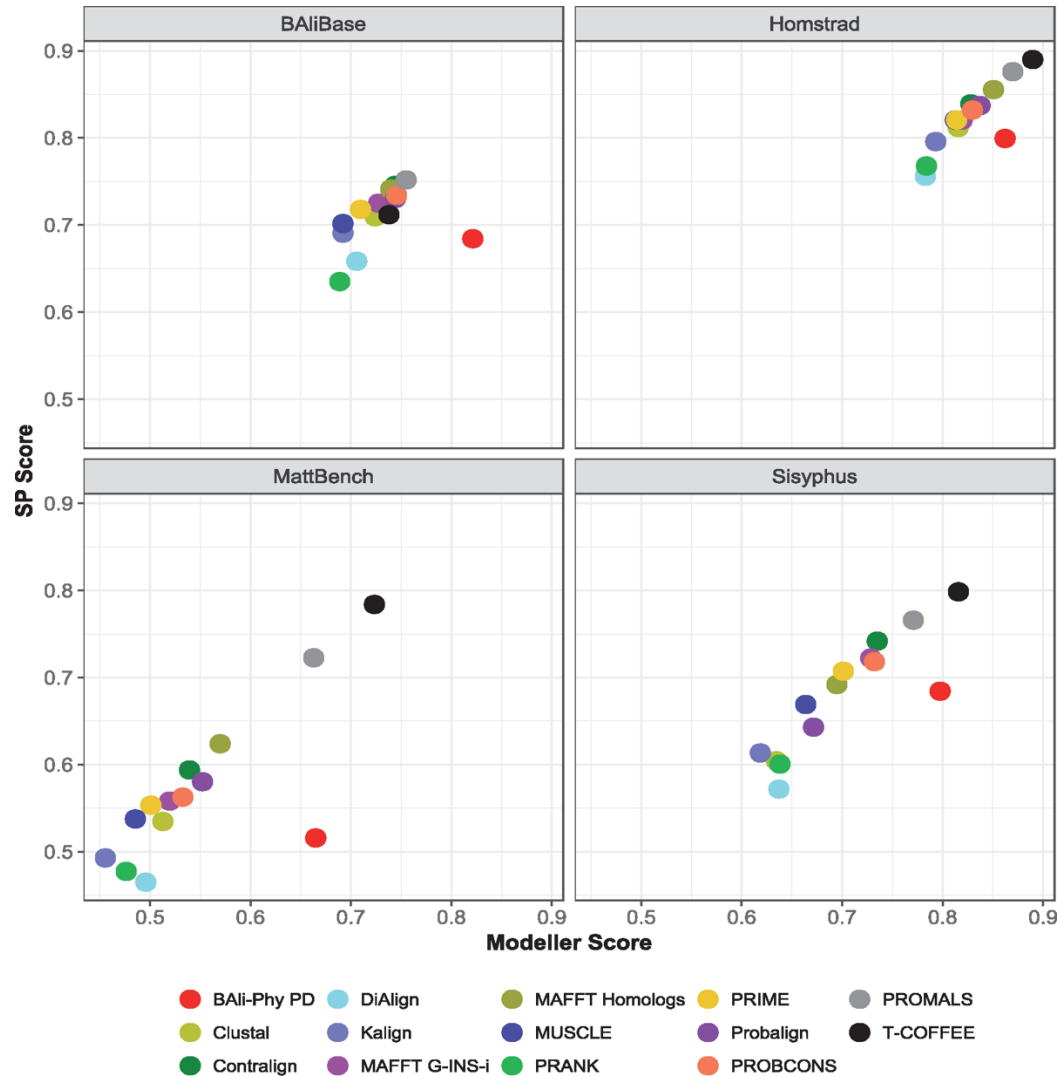
BAlI-Phy is best on small **simulated protein** datasets!



BAlI-Phy is best!

● BAlI-Phy ● ContrAlign ● MUSCLE ● PRIME ● PROBCONS
● Clustal ● MAFFT G-INS-i ● PRANK ● Probalign

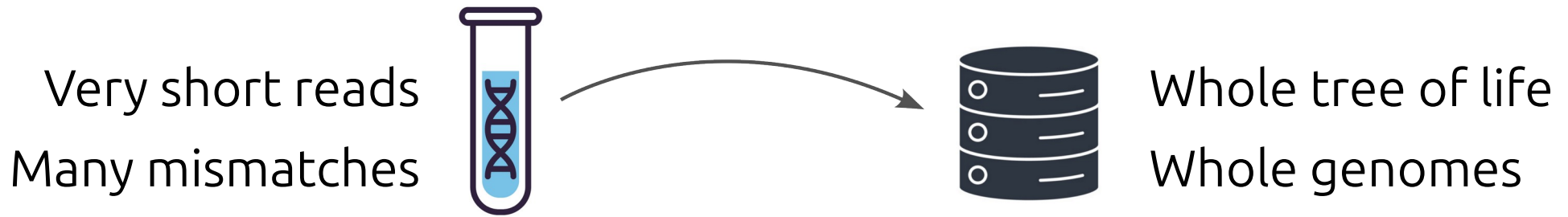
BAlI-Phy not so great on on 1192 small **biological protein** datasets



T-Coffee and PROMALS
are best!

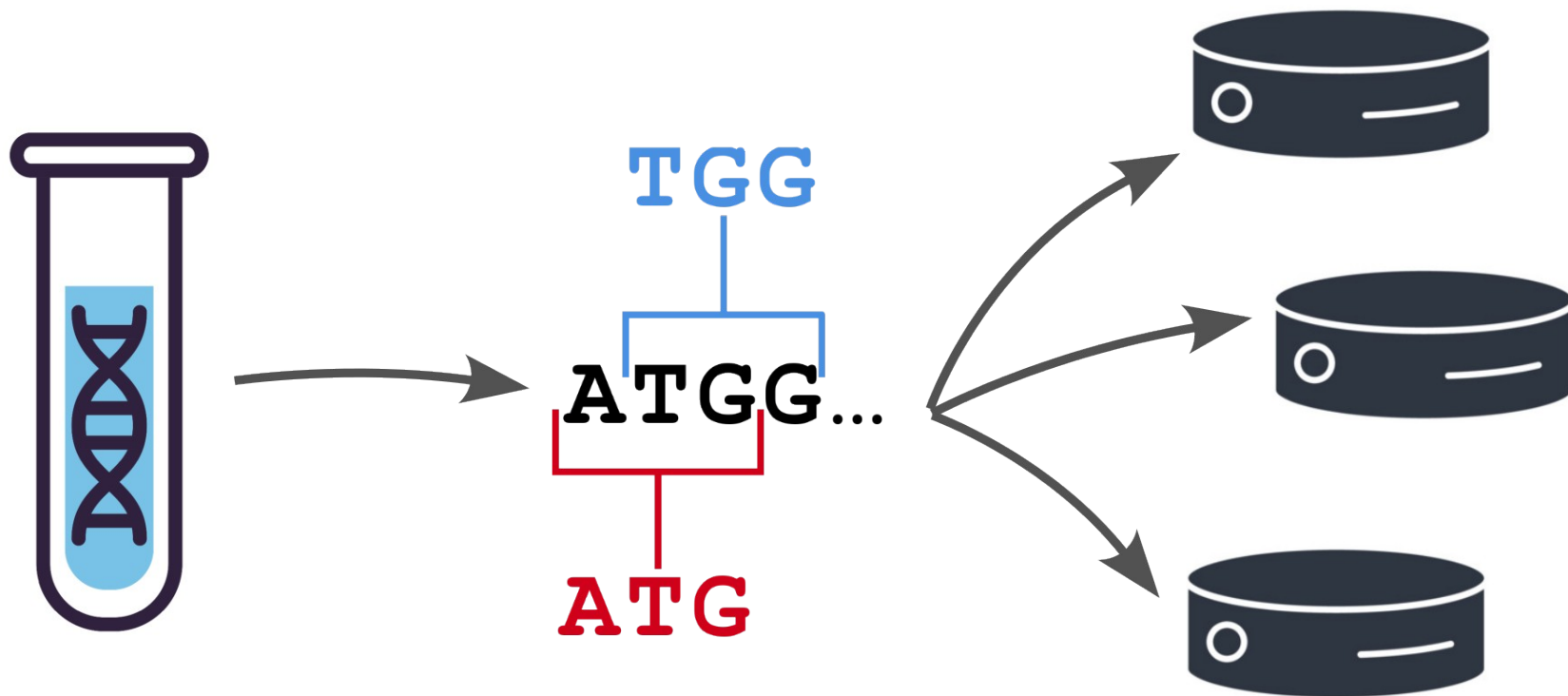
BAlI-Phy good for
Modeler score, but not
so good for SP-Score
(e.g., MAFFT better)

Taxonomic assignment and mapping of ancient environmental DNA at scale



We want mapping
– not just taxonomic assignment –
to estimate damage

Pre-select reference shards via k-mers,
then map reads to selected shards

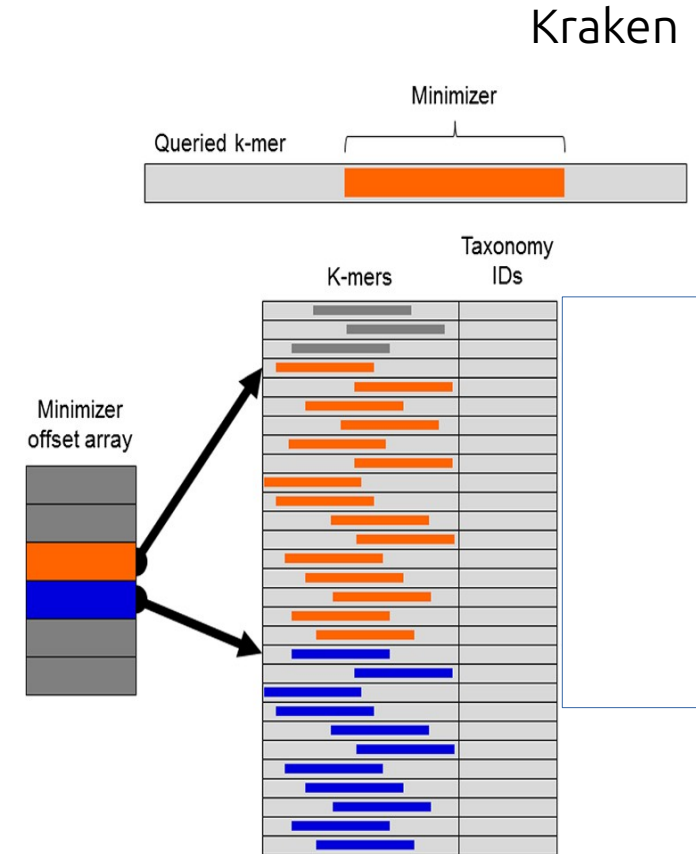


Problems with established approaches

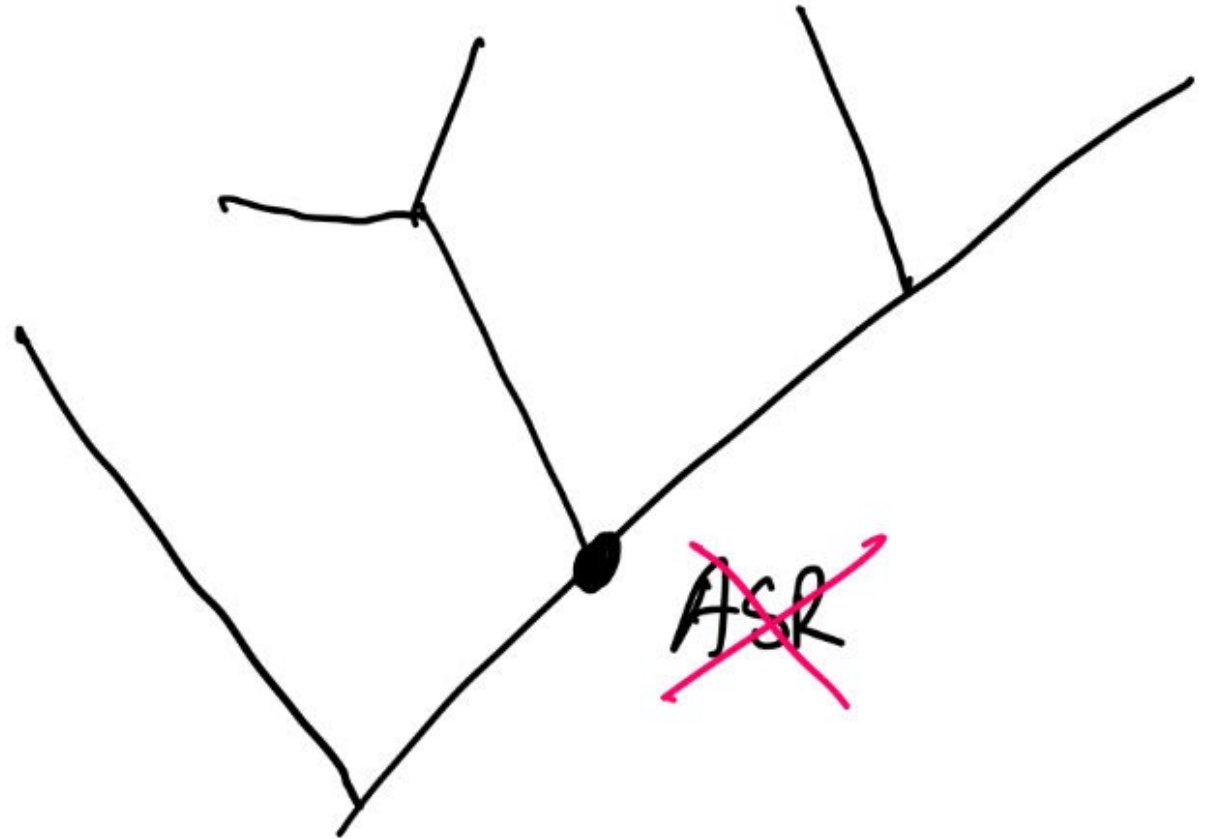
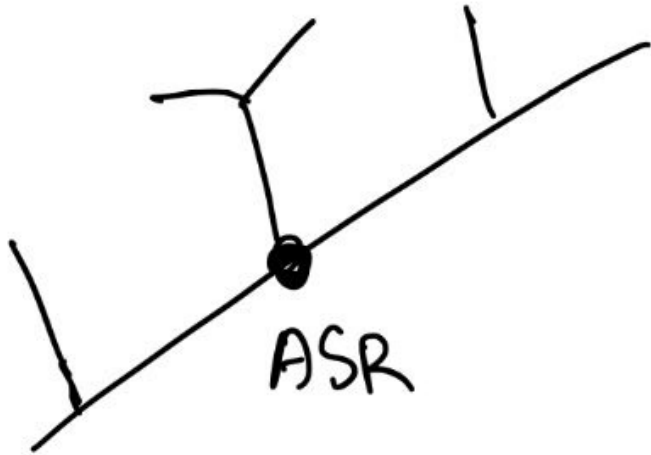
Storing the k-mer values:

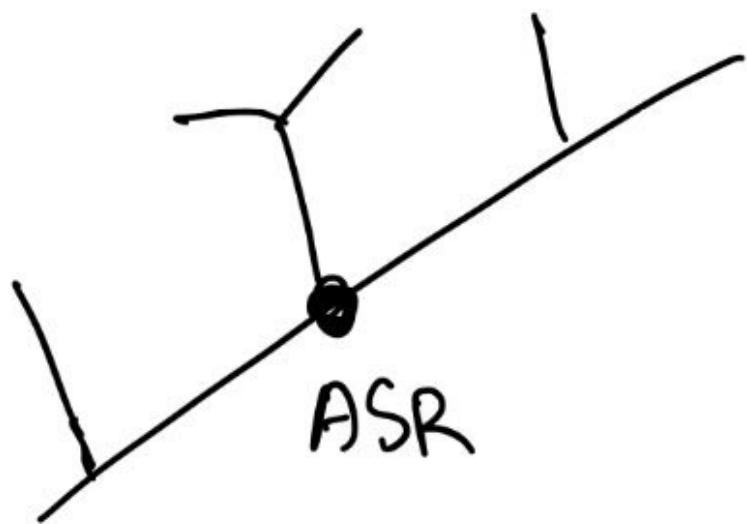
- Too much to keep all occurring k-mers
- Subset the kmers (e.g., minimizers)
- Loss of sensitivity

For aDNA, we need to keep more k-mers with shorter k (in the order of the mapping seed length)

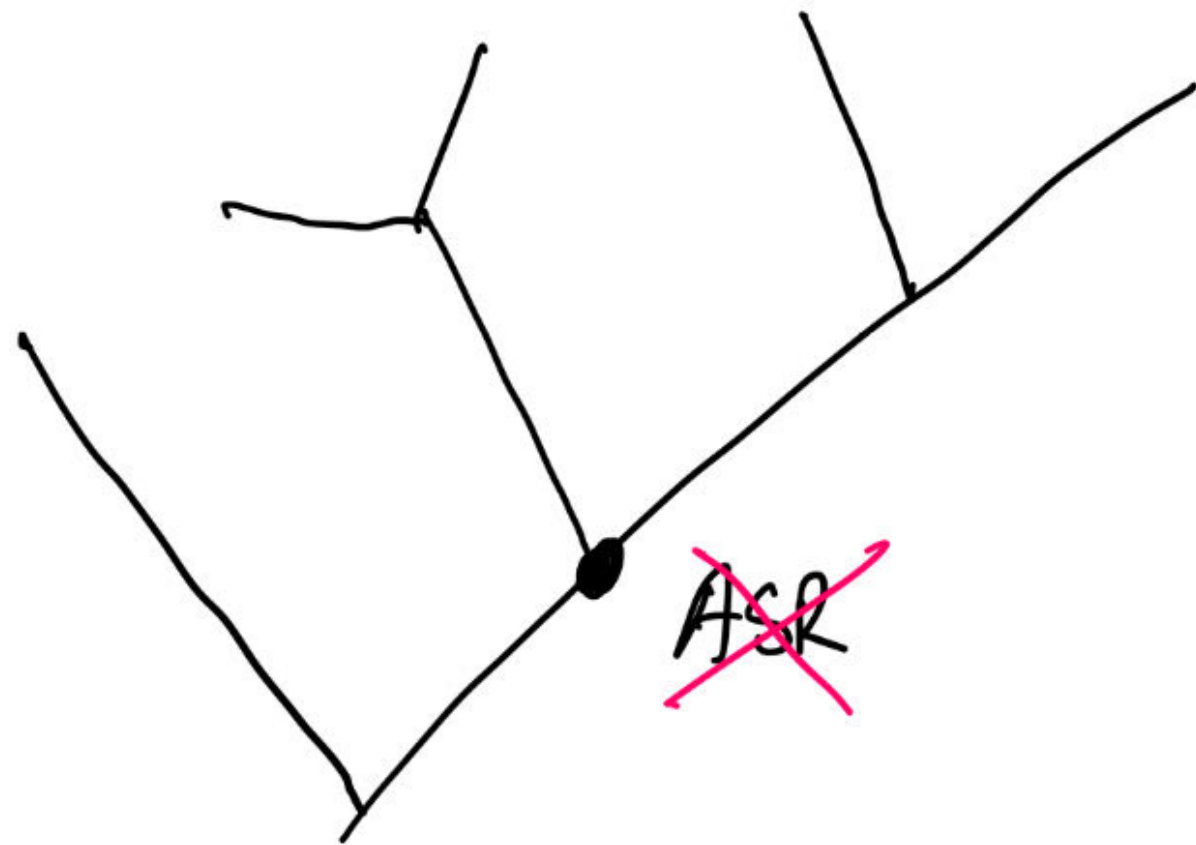


Ancestral sequence reconstruction beyond the twilight zone?



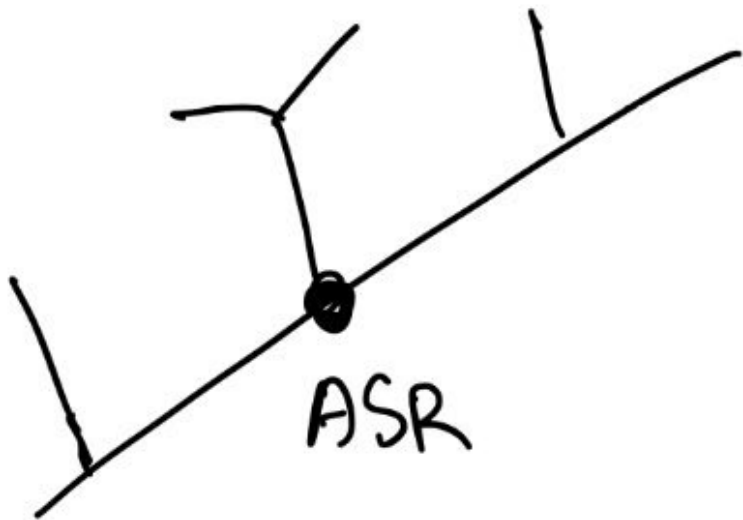


Structure



Sequence

MSEAKF...

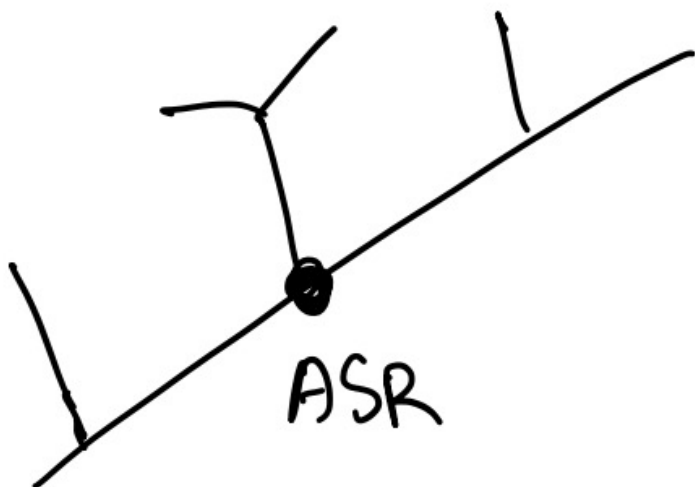


Structure
(3Di)

AKDVRC...
(3Di)

The problem:
going back to amino acid sequence from 3Di sequence

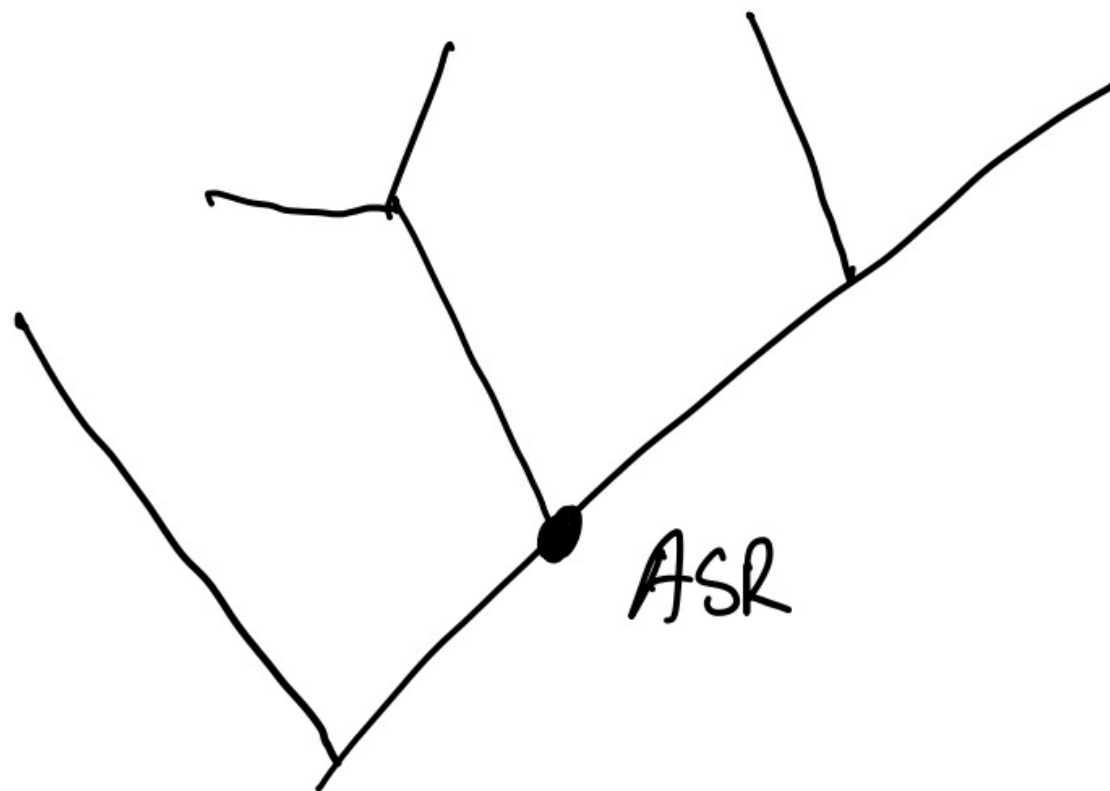
MSEAKF... \Leftarrow AKDVRC...
(3Di)



Certain 3Di sequences

Weighting?

Joint probability?



Uncertain AA sequences

- Let T be a tree with branch length on the node set V . Let sequences $S = S_1 \dots S_n$ evolve on T using some CTMC like JC or GTR (note: $S_i \in \{A, C, G, T\}^L$). Let $T_{i,j}$ be distances between S_i and S_j on T .
- Vanila: Find a function $\Phi(s) : \{A, C, G, T\}^L \rightarrow \mathbb{R}^d$ such that

$$\forall i, j : \lim_{L \rightarrow \infty} d(\Phi(S_i), \Phi(S_j)) \rightarrow T_{i,j}$$

for some measure of distance d such as ℓ_2 . What is the smallest d where this is possible?

- Extended: Instead of \mathbb{R}^d , we can look for embedding in another geometric space.