

A PROBABILISTIC FRAMEWORK FOR SEQUENCE DISTANCES & GENOME-WIDE (PHYLOGENETIC) PATTERN MATCHING

Motivation

Estimating alignment-free distances between sequences of any length:

- Scalable • Accurate and sensitive (e.g., >15%) • Modeling uncertainty

Applications of such a framework:

Metagenomics:

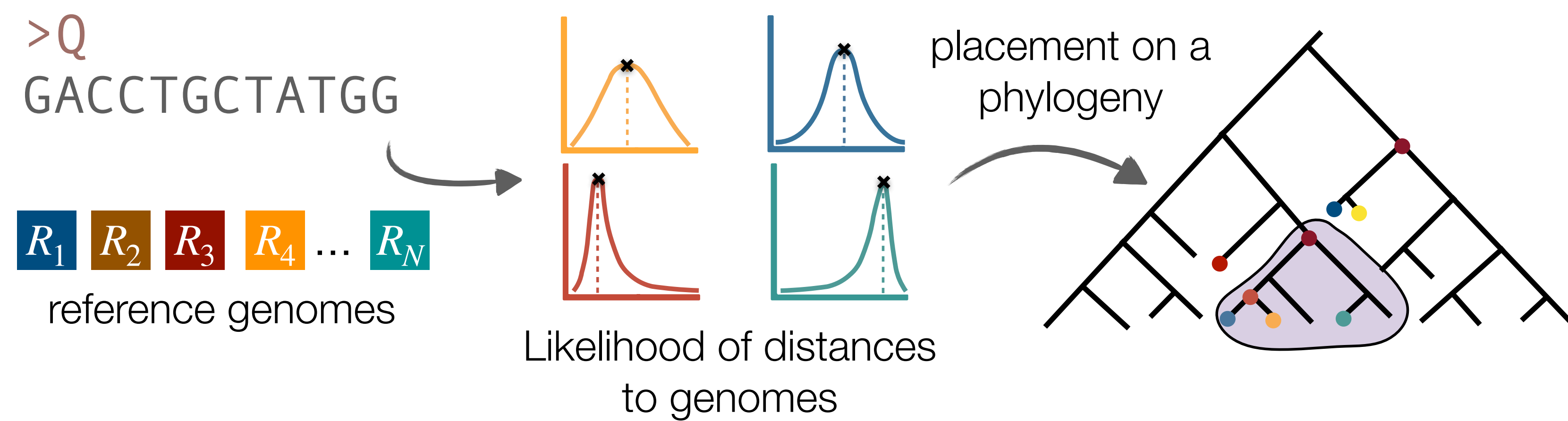
- Read-to-genome distances
- Phylogenetic placement

Comparative genomics:

- Horizontal gene transfer
- Ultra-conserved elements

Problem: alignment-free calculation of distances

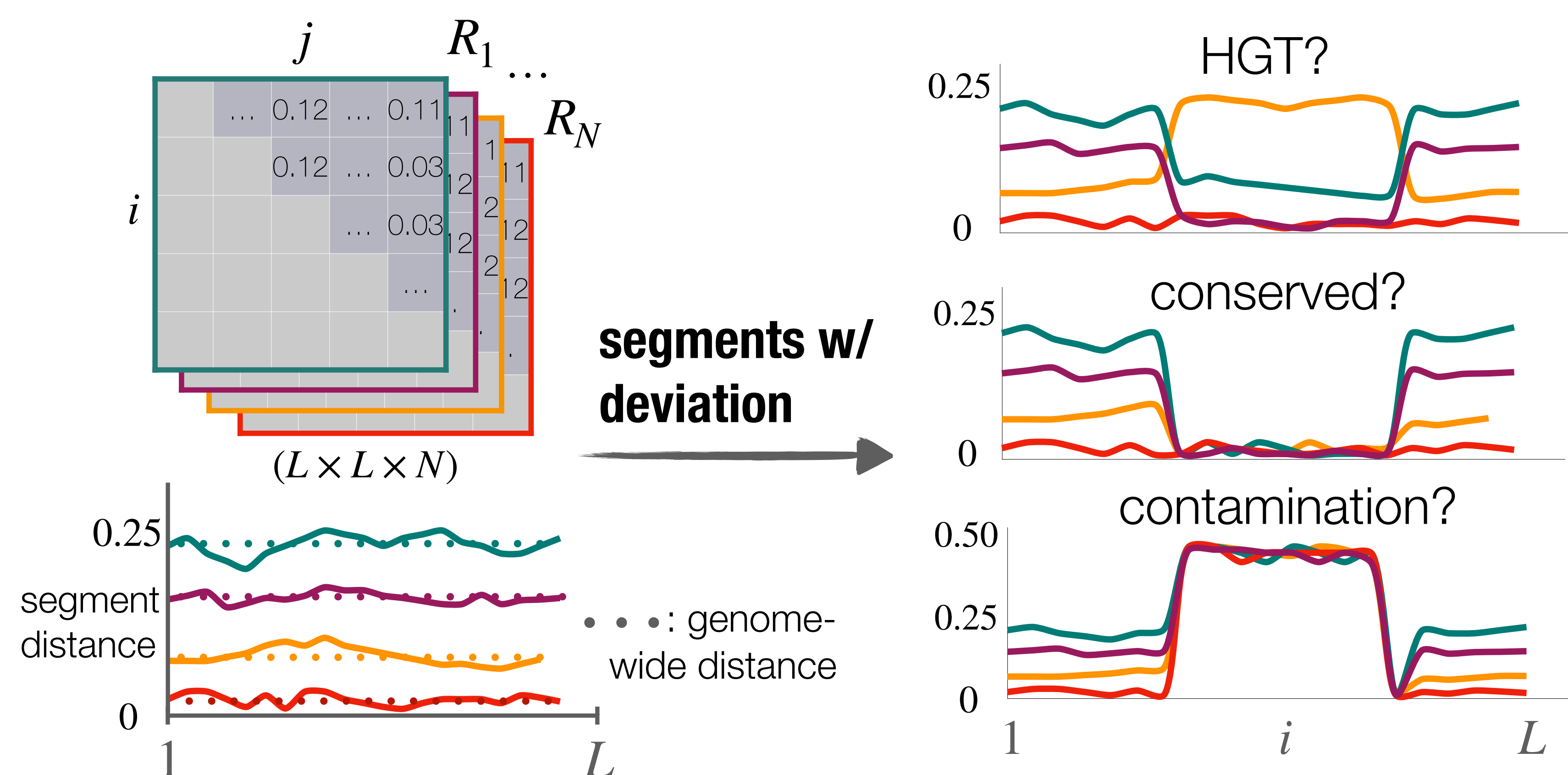
Modeling the distance between Q (e.g., read, contig, interval) and \mathcal{R} :



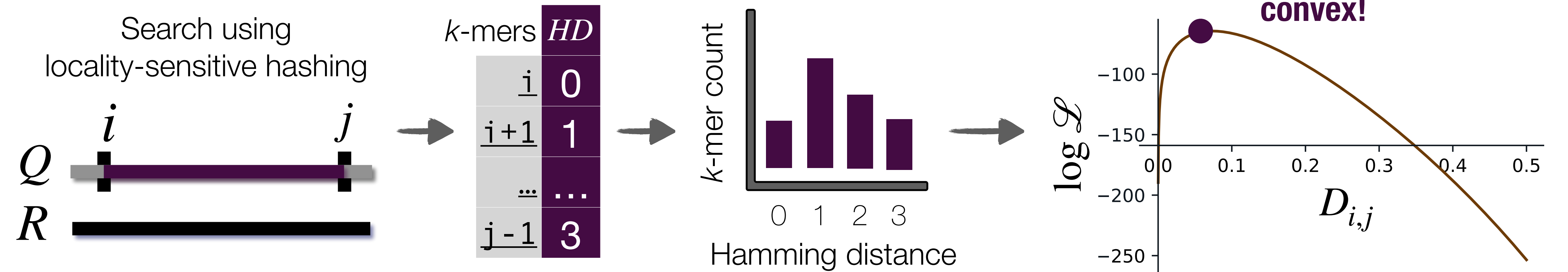
For a query Q and a reference $R \in \mathcal{R}$

- Compute $\mathcal{L}(D_{ij}; Q, R)$ using k -mers, where $D_{ij} = d(Q[i:j], R)$.
- Estimate $D = d(Q, R)$ and place Q on a phylogeny over \mathcal{R} .
- For all (i, j) and a threshold Δ , decide (in constant time) if $D_{ij} < \Delta$ (or $> \Delta$)
- Given τ , find all “maximal intervals” with $D_{ij} < \Delta$ and $j - i \geq \tau$
 $d(Q[i':j'], R) \geq \Delta$ for $i' \leq i \leq j \leq j'$, $(i, j) \neq (i', j')$

Detecting deviations of “local distances” from genome-wide distances:



Finding homologous k -mers and pseudo-likelihood formulation



Two key events:

For $x \in Q[i:j]$ and $y \in R$;

- They are d substitutions away:

$$P(\text{HD}(x, y) = d) = D^d (1 - D)^{(k-d)} \binom{k}{d}$$

- LSH values collide (true positive rate):

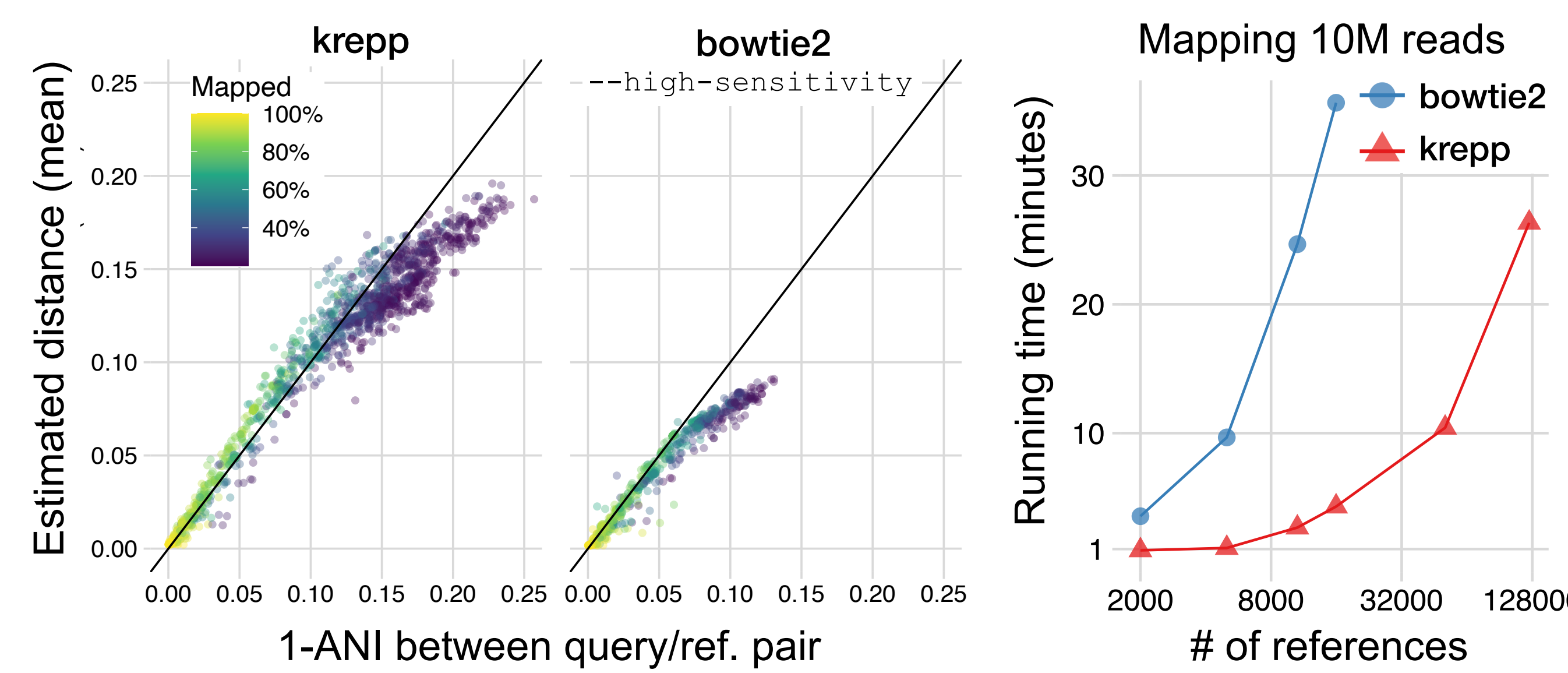
$$P(\text{LSH}(x) = \text{LSH}(y)) = \binom{k-h}{d} / \binom{k}{d}$$

$$\text{misses: } u = (j - i) - \sum_{d=0}^{\delta} v_d \quad \text{matches: } \mathbf{v} = [v_0, v_1, \dots, v_{\delta}]$$

$$\mathcal{L}(D; k, h, \delta, u, \mathbf{v}) = P_{\text{miss}}(D; k, h, \delta)^u \prod_{d=0}^{\delta} P_{\text{match}}(D; d, k, h)^{v_d}$$

We define probability of observing u misses & v_d HD = d matches based on **i**) subsampling rate **ii**) number of mutations **iii**) LSH match!

Estimating read-to-genome distances



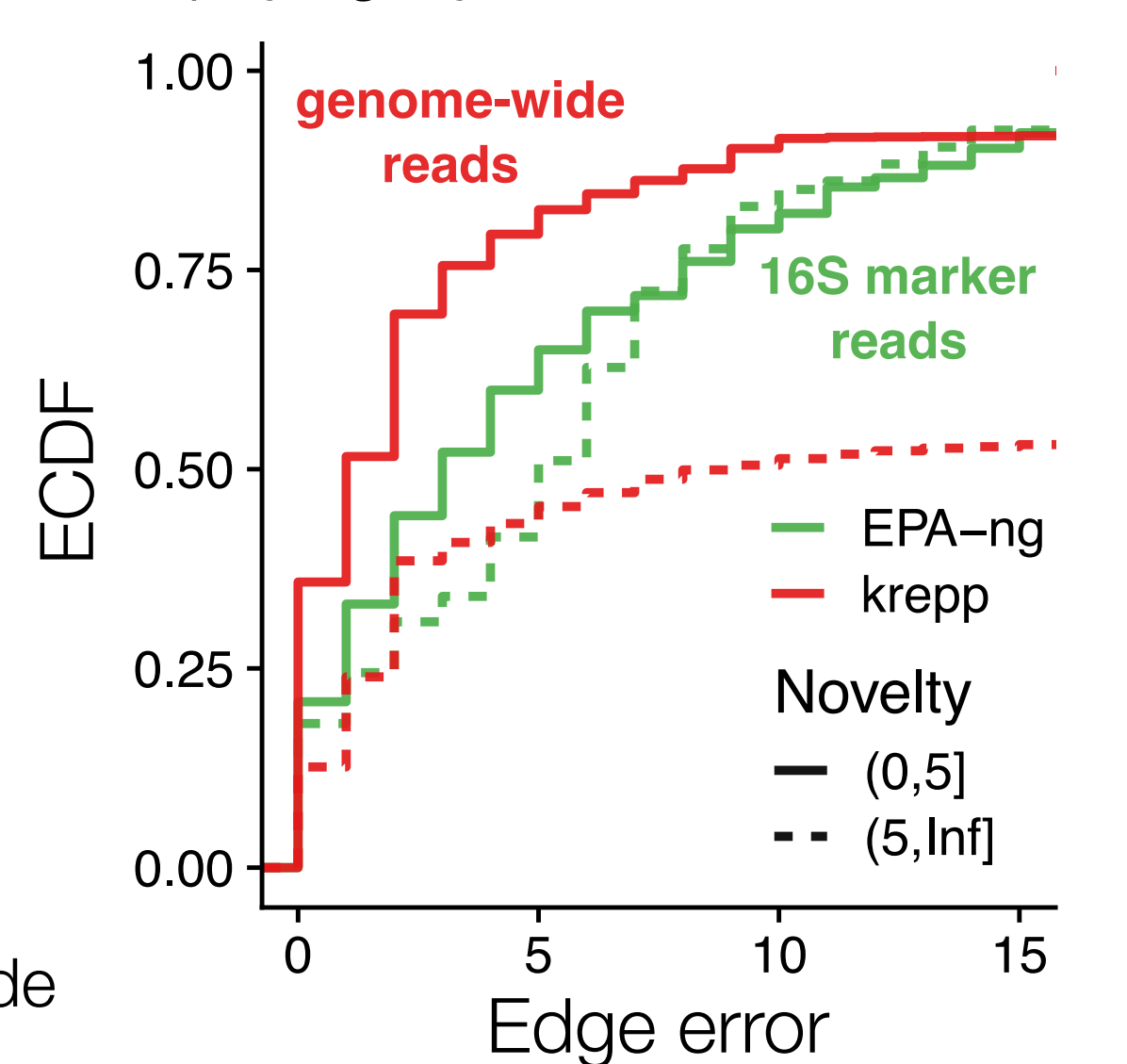
Phylogenetic placement

ancestral distances & likelihood-ratio test for placement:

$$\lambda_{LR} = \frac{\mathcal{L}_*(D; k, h, \delta, u_*, \mathbf{v}_*)}{\mathcal{L}_*(D^*; k, h, \delta, u_*, \mathbf{v}_*)}$$

place as sister to the largest indistinguishable clade

a phylogeny with 10,000 leaves



Distance-based segmentation & genome-wide pattern matching

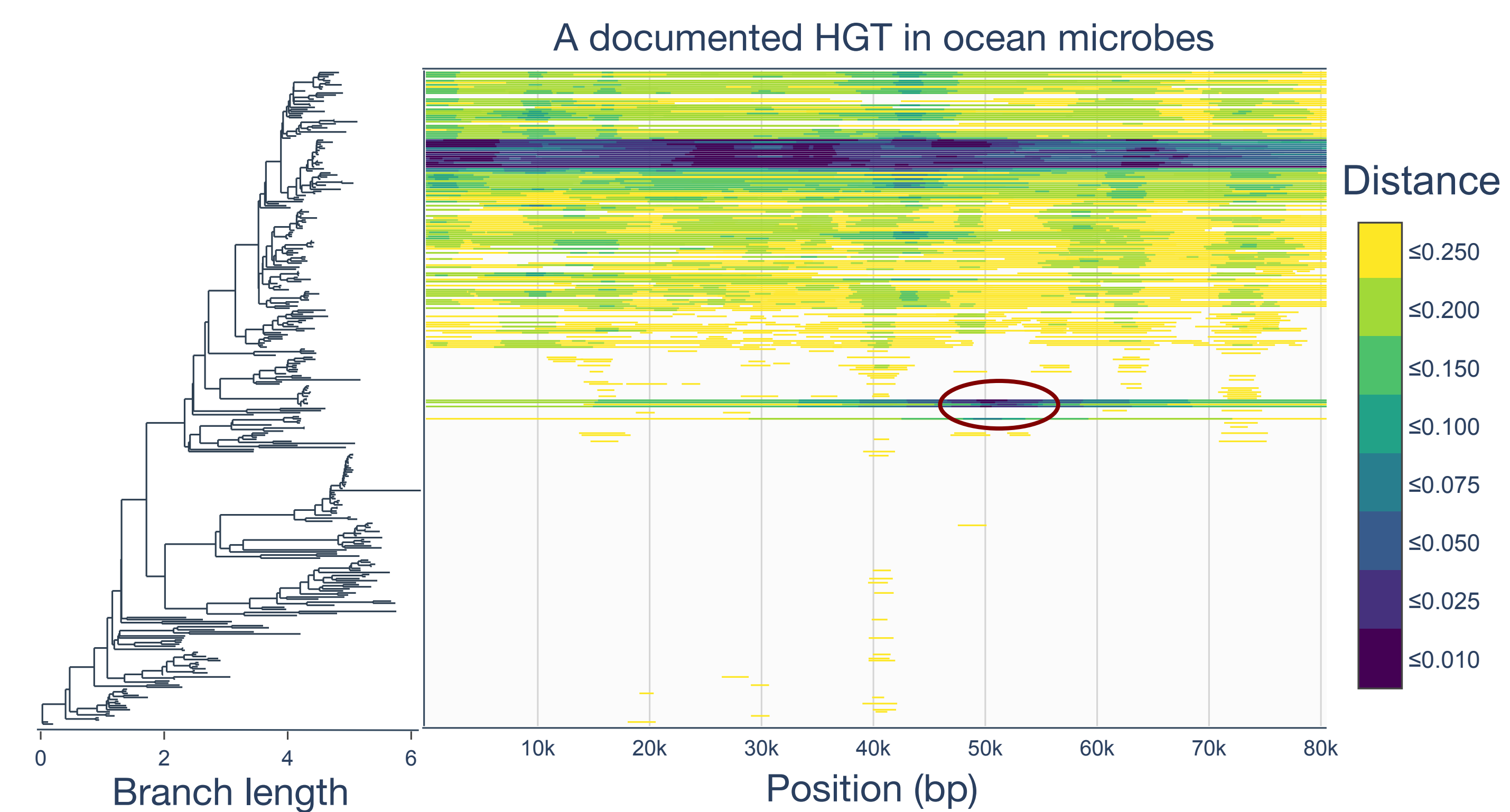
$$\text{MLE } \hat{D} = \Delta \rightarrow \ell'(\Delta) = 0$$

$$\text{If } \Delta < \hat{D} \rightarrow \ell'(\Delta) > 0$$

$$\text{If } \Delta > \hat{D} \rightarrow \ell'(\Delta) < 0$$

S : prefix-sum array for $\ell'(\Delta)$

- $\hat{D}_{ij} < \Delta$ iff $s_i > s_j$ (constant time)
- All maximal intervals (linear time)



krepp & gdiff

krepp for read-to-genome distances and placement:



gdiff for genome-wide local distances and segmentation:



Soon...