

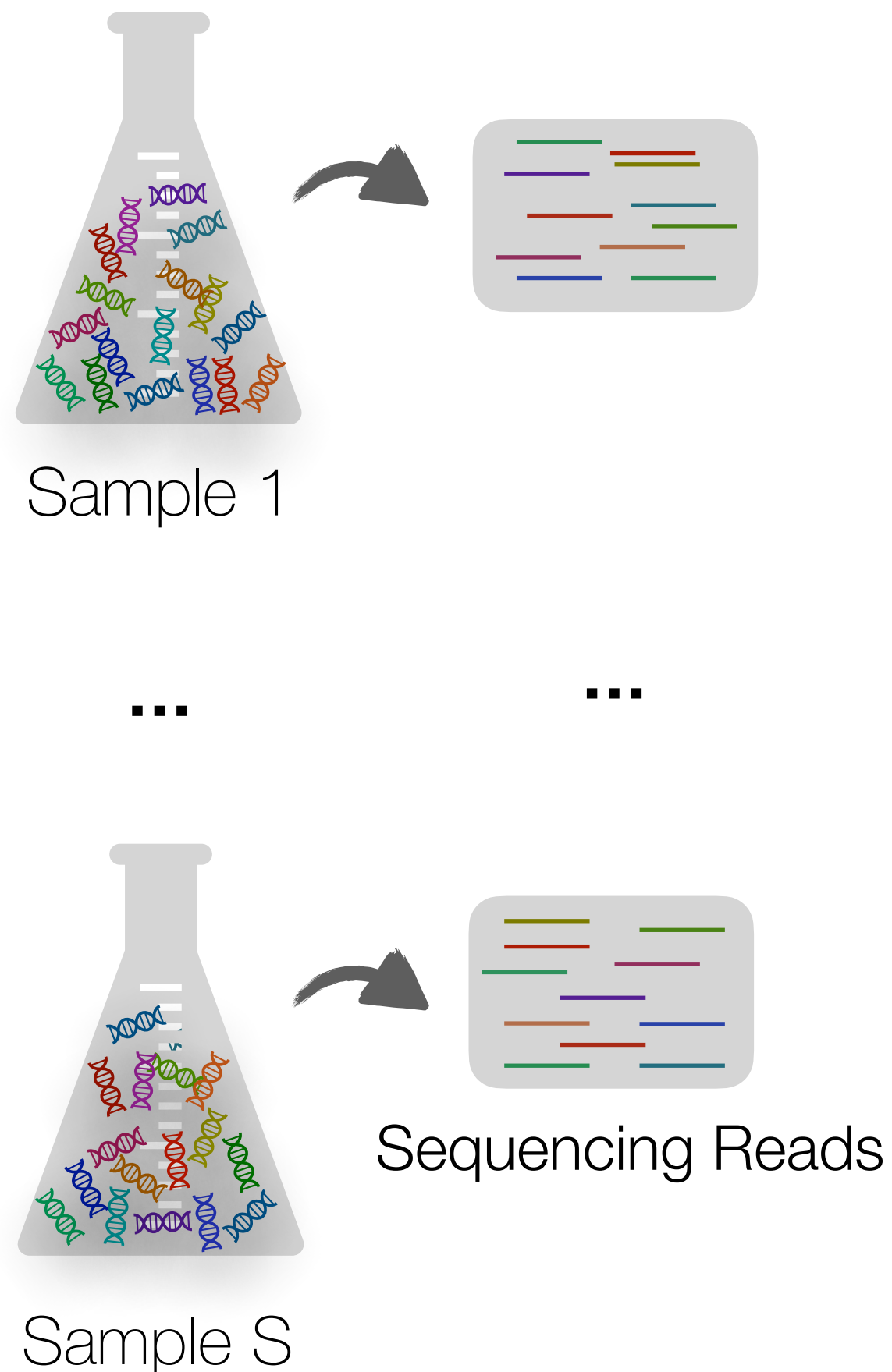
**Estimating read-genome distances  
with homologous  $k$ -mer matches  
*(and genome-wide phylogenetic placement)***

Ali Sapci

8 April 2025

# Identifying metagenomic sequences and comparing samples

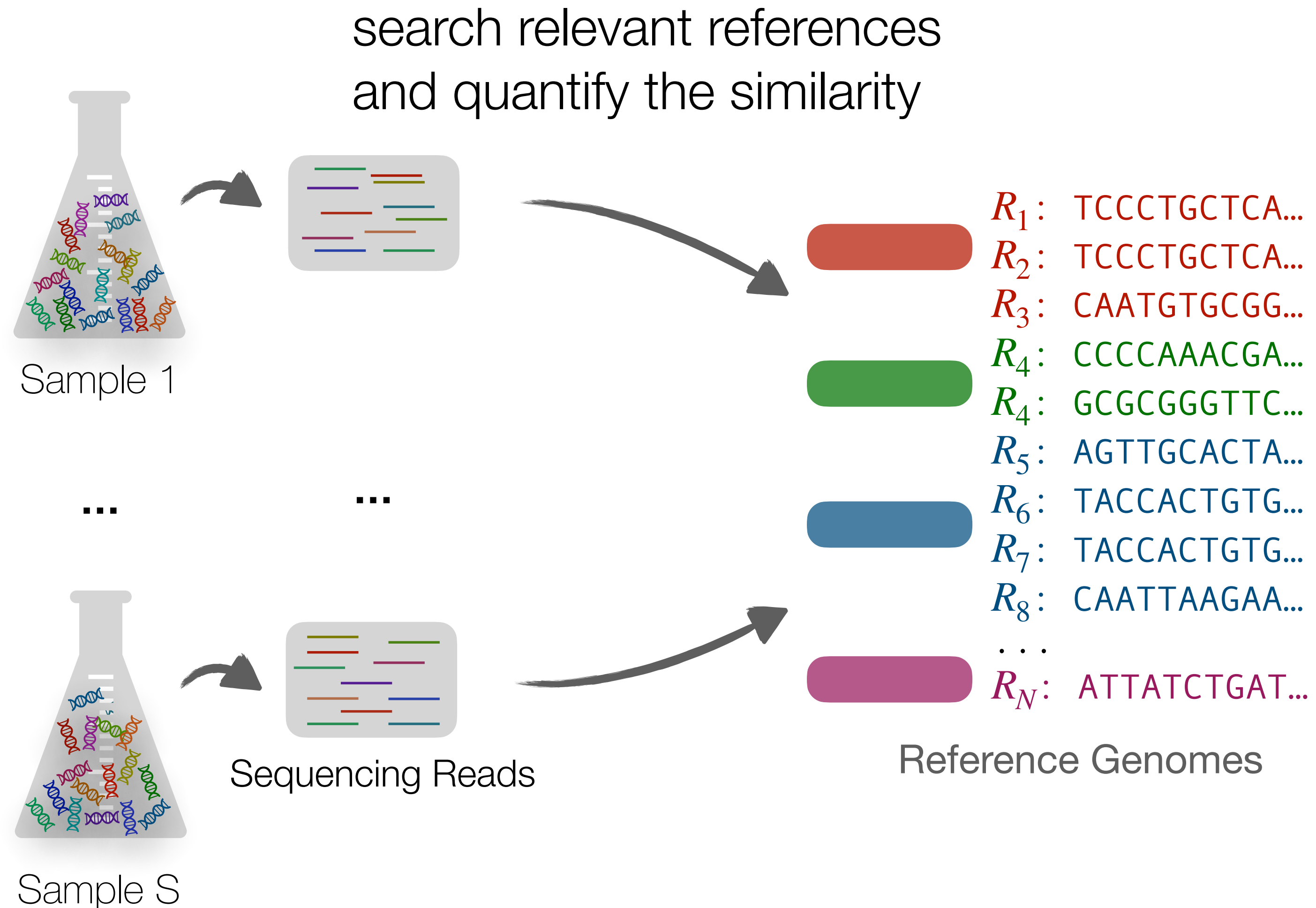
# Identifying metagenomic sequences and comparing samples



## High-level goals:

- Analyze & monitor present microbial communities
- Compare microbial composition of samples
- Measure diversity and detect novel sequences

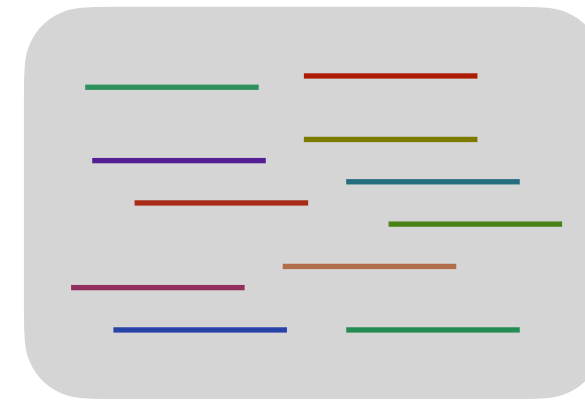
# Identifying metagenomic sequences and comparing samples



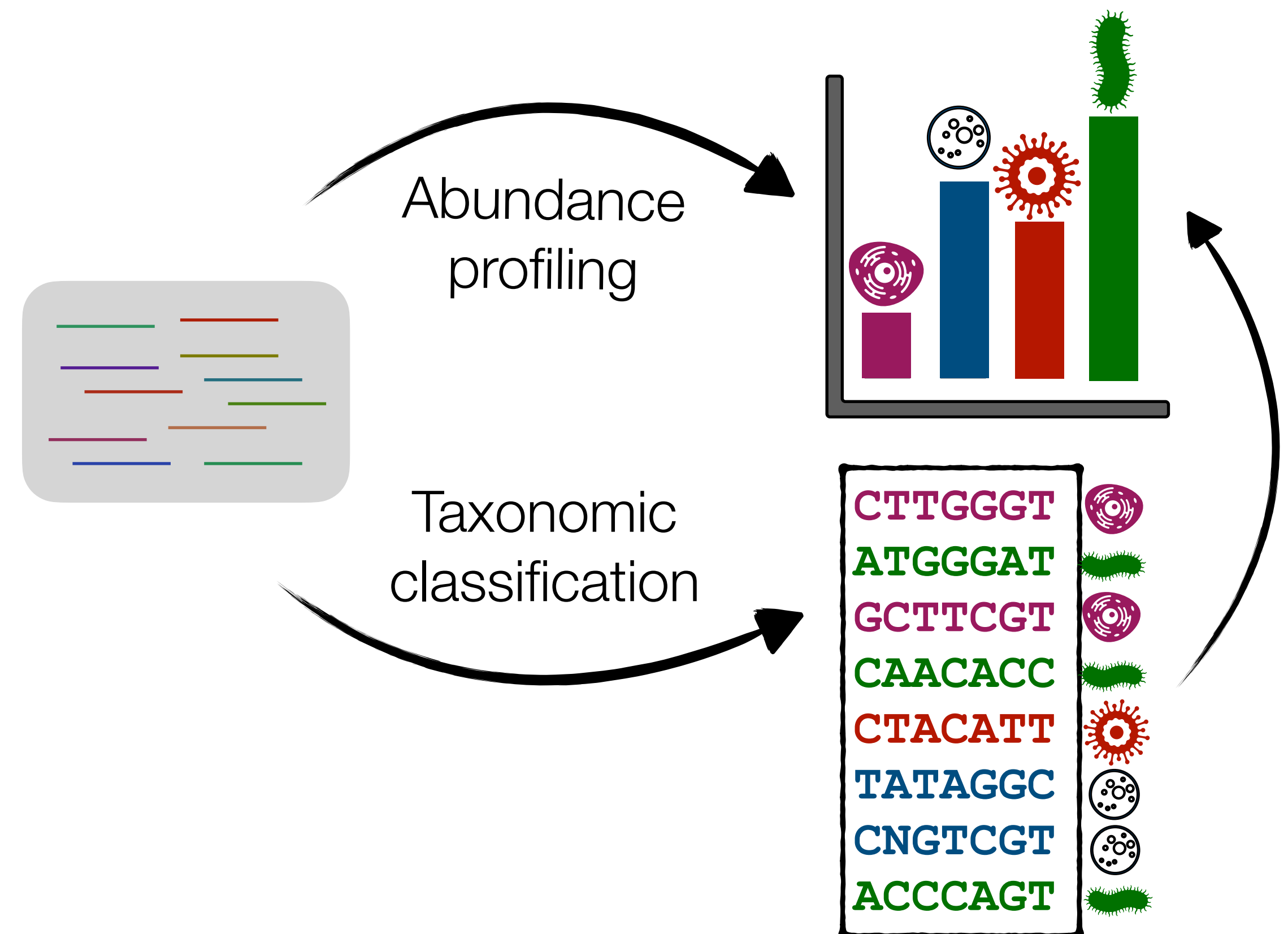
## High-level goals:

- Analyze & monitor present microbial communities
- Compare microbial composition of samples
- Measure diversity and detect novel sequences

# Poor man's solution: taxonomic profiling

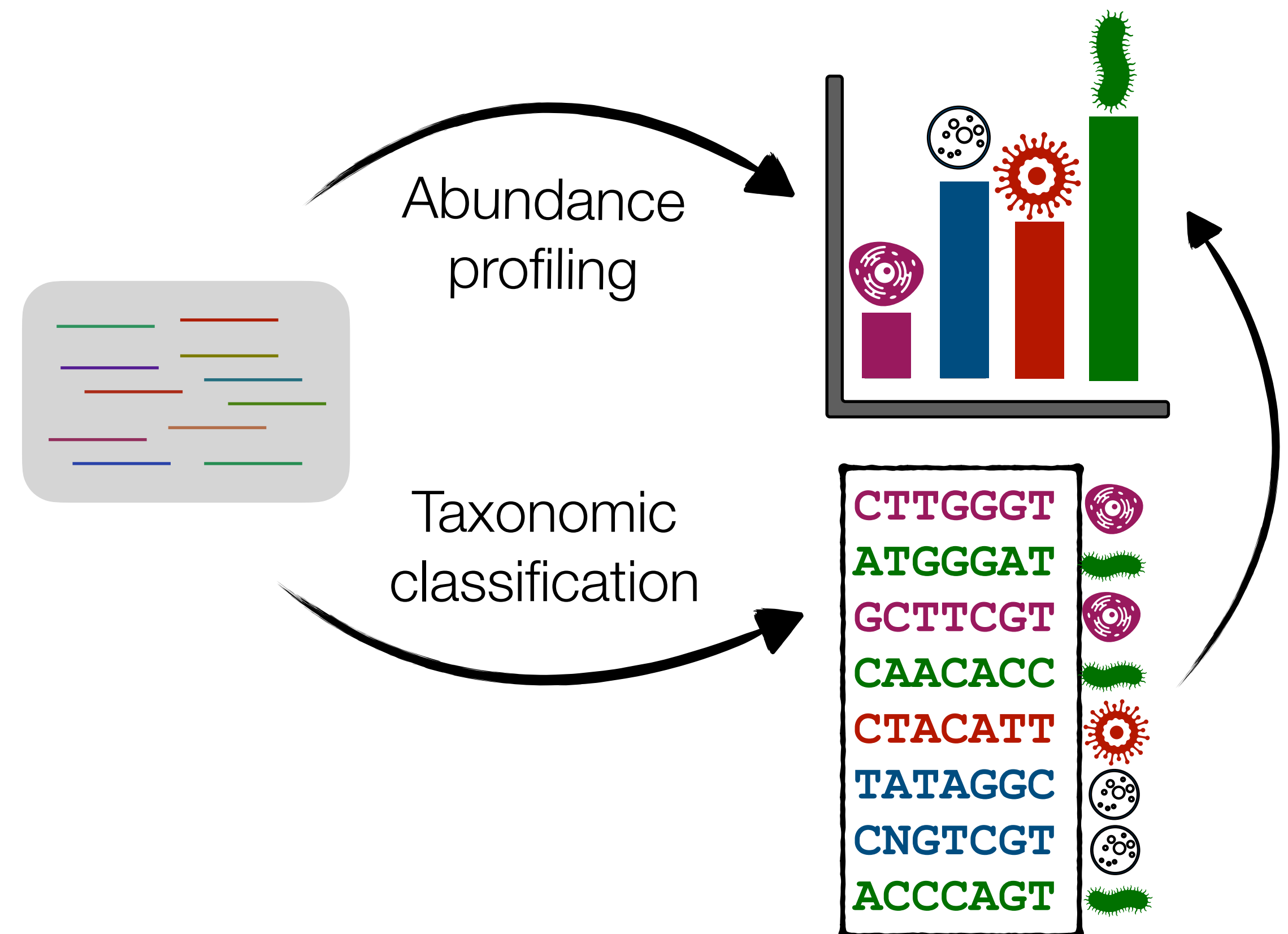


# Poor man's solution: taxonomic profiling



# Poor man's solution: taxonomic profiling

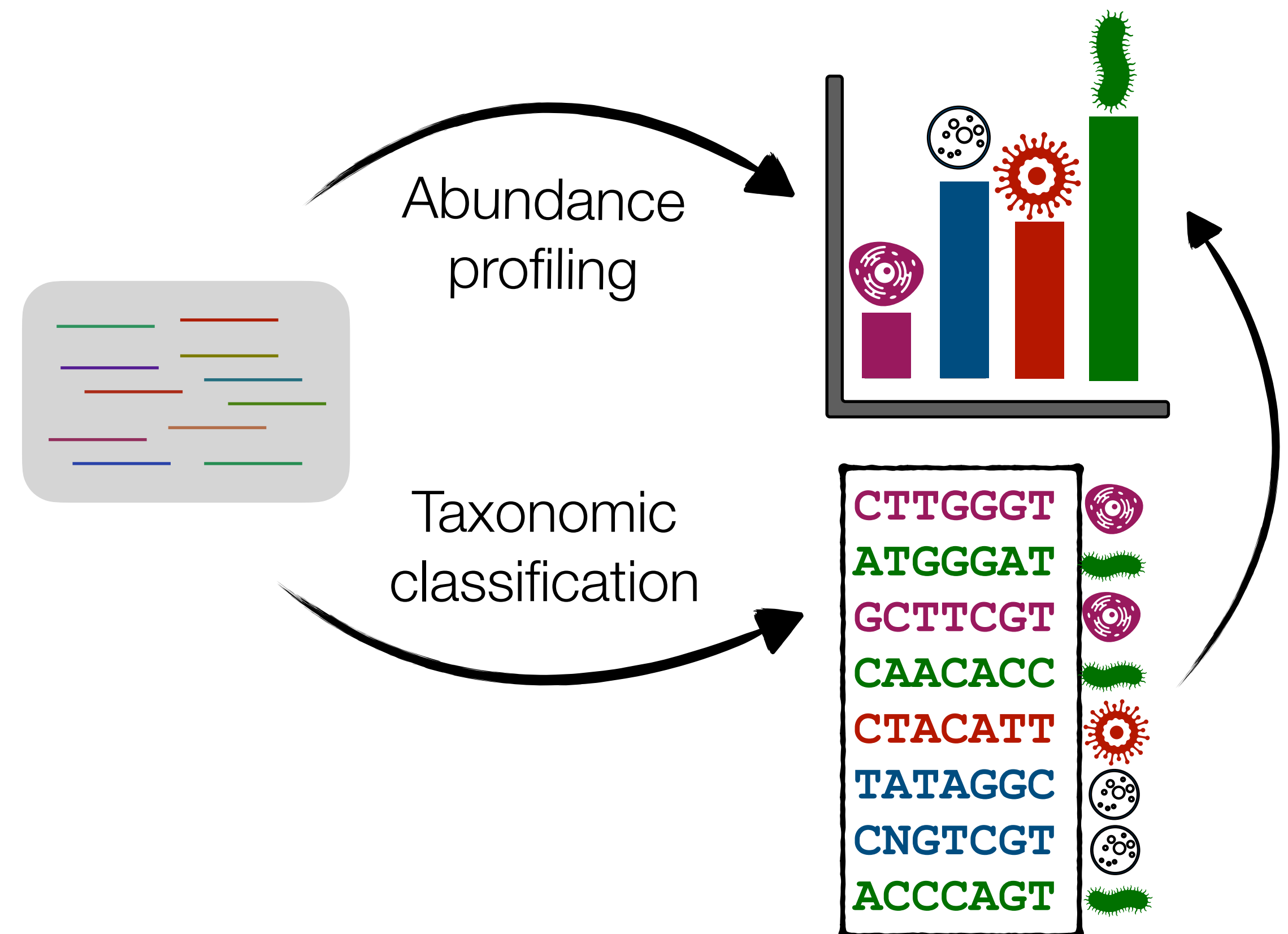
compare profile vectors  
(e.g., Bray-Curtis dissimilarity)



# Poor man's solution: taxonomic profiling

- **Fast and scalable methods** based on *k*-mer search
- **Limited:**
  - ▶ has low resolution
  - ▶ often ambiguous (e.g., HGT)
  - ▶ omits within-group diversity
  - ▶ no notion of *distance*, novelty?

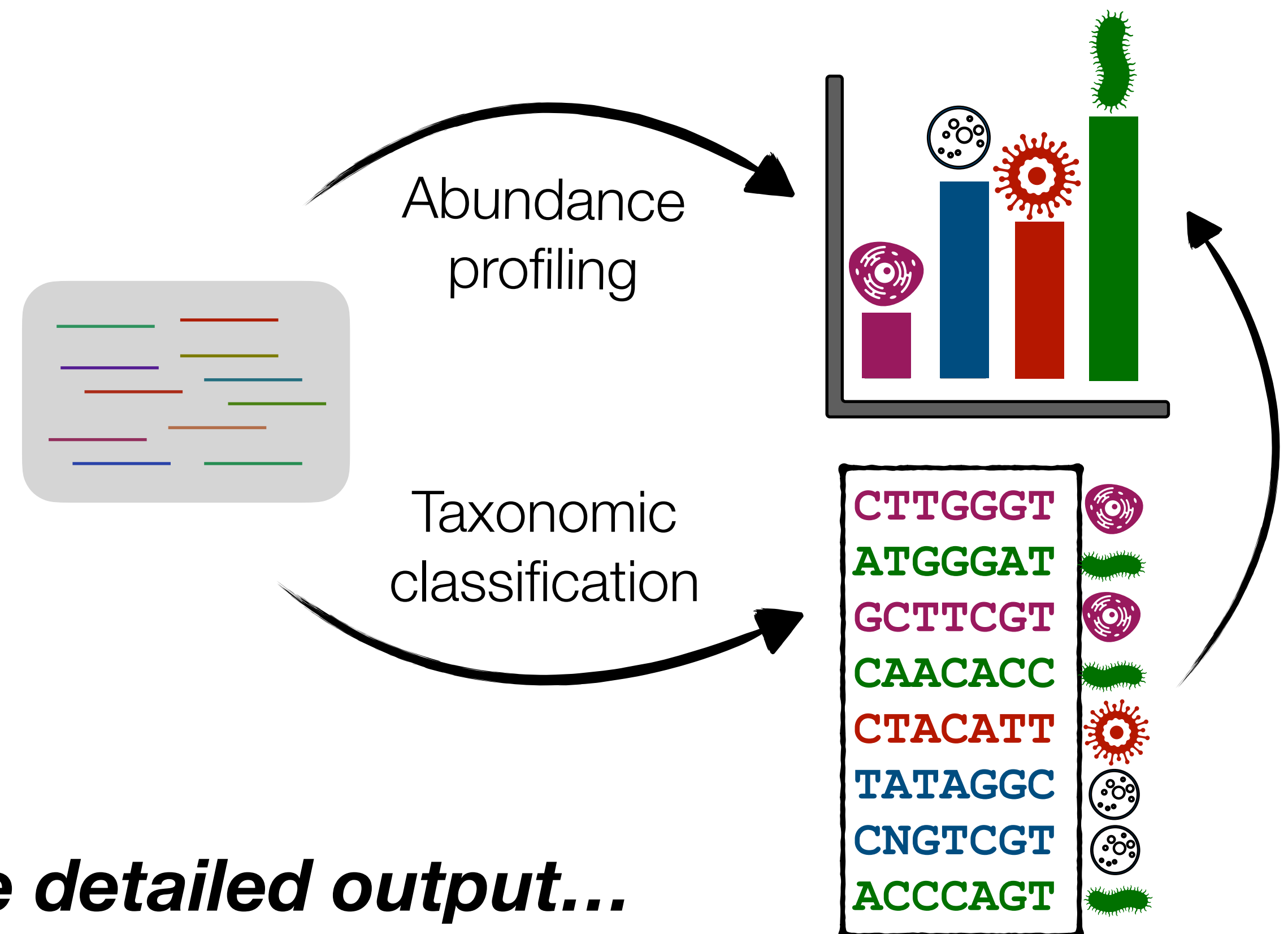
compare profile vectors  
(e.g., Bray-Curtis dissimilarity)



# Poor man's solution: taxonomic profiling

- **Fast and scalable methods** based on *k*-mer search
- **Limited:**
  - ▶ has low resolution
  - ▶ often ambiguous (e.g., HGT)
  - ▶ omits within-group diversity
  - ▶ no notion of *distance*, novelty?

compare profile vectors  
(e.g., Bray-Curtis dissimilarity)



*We need a more detailed output...*

**What alternatives do we have?**

# **What alternatives do we have?**

- **Sample-wide containment analysis (using MinHash)?**

# What alternatives do we have?

- **Sample-wide containment analysis (using MinHash)?**
  - higher resolution compared to taxonomic profiles (+)
  - no assignments for individual sequences (-)
  - no distances btw. queries and references (-)

# What alternatives do we have?

- **Sample-wide containment analysis (using MinHash)?**
  - higher resolution compared to taxonomic profiles (+)
  - no assignments for individual sequences (-)
  - no distances btw. queries and references (-)
- **Using marker genes?**

# What alternatives do we have?

- **Sample-wide containment analysis (using MinHash)?**
  - higher resolution compared to taxonomic profiles (+)
  - no assignments for individual sequences (-)
  - no distances btw. queries and references (-)
- **Using marker genes?**
  - aligning to a MSA is possible & and distances (+)
  - limited and potentially biased (-)
  - inability to capture novel sequences and queries (-)

# What alternatives do we have?

- **Sample-wide containment analysis (using MinHash)?**
  - higher resolution compared to taxonomic profiles (+)
  - no assignments for individual sequences (-)
  - no distances btw. queries and references (-)
- **Using marker genes?**
  - aligning to a MSA is possible & and distances (+)
  - limited and potentially biased (-)
  - inability to capture novel sequences and queries (-)
- **Aligning reads to all references?**

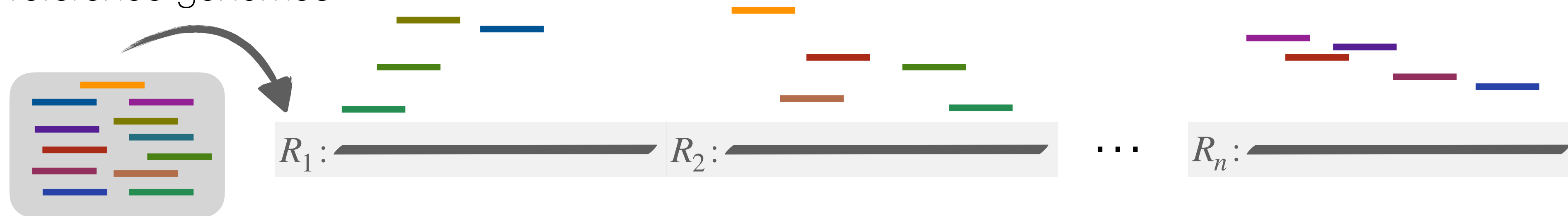
# Aligning sequences of unknown origins

align reads to  
reference genomes



# Aligning sequences of unknown origins

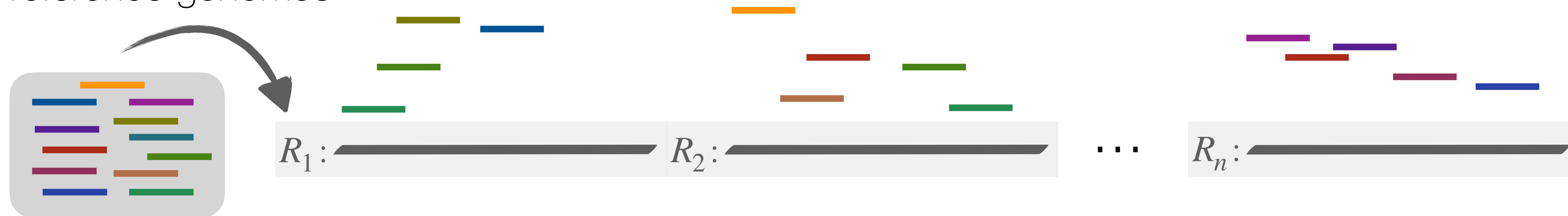
align reads to  
reference genomes



(+) quantifying similarity — as detailed as it gets (in fact redundant)

# Aligning sequences of unknown origins

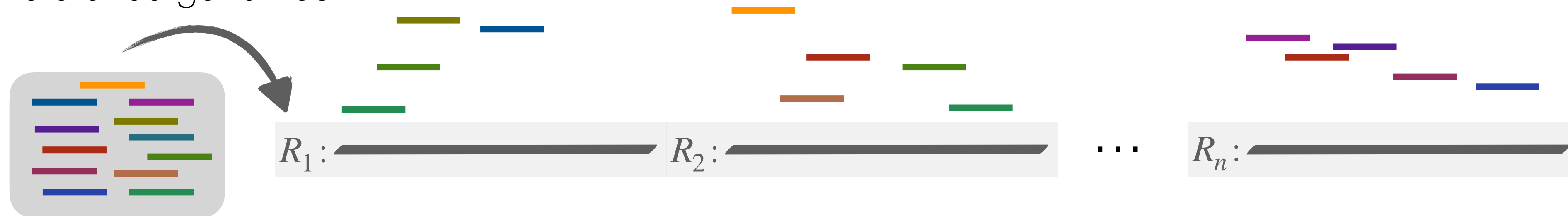
align reads to  
reference genomes



- (+) quantifying similarity — as detailed as it gets (in fact redundant)
- (+) placing sequences on a phylogeny — not new but limited

# Aligning sequences of unknown origins

align reads to  
reference genomes



- (+) quantifying similarity — as detailed as it gets (in fact redundant)
- (+) placing sequences on a phylogeny — not new but limited
- (-) not scalable for large  $n$ , even with efficient indexes

# Aligning sequences of unknown origins

align reads to  
reference genomes



- (+) quantifying similarity — as detailed as it gets (in fact redundant)
- (+) placing sequences on a phylogeny — not new but limited
- (-) not scalable for large  $n$ , even with efficient indexes
- (-) not good for higher distances & novel sequences (>15%)

**Our goal: the best of both worlds**

# **Our goal: the best of both worlds**

Given a query sequence:

# Our goal: the best of both worlds

Given a query sequence:

- estimate accurate distances for references
  - ▶ scale to modern references (>100K microbial genomes)
  - ▶ akin to alignment — but without solving the difficult problem

# Our goal: the best of both worlds

Given a query sequence:

- estimate accurate distances for references
  - ▶ scale to modern references (>100K microbial genomes)
  - ▶ akin to alignment — but without solving the difficult problem
- place it on a large reference phylogeny (>100,000)
  - ▶ always relies on marker genes — with one exception

# **Problem statement**

# Problem statement

Given:

- query sequence  $q$
- set of references  $\mathcal{R} = \{R_1, \dots, R_n\}$
- a backbone phylogeny  $T$

$>q$   
CCTGCTA...



## reference genomes

```
R1: TCCCTGCTCA...  
R2: TCCCTGCTAA...  
R3: CCCCTGGCAG...  
R4: ATTATCTGAT...  
...  
Rn: CCCCAAACAA...
```

# Problem statement

Given:

- query sequence  $q$
- set of references  $\mathcal{R} = \{R_1, \dots, R_n\}$
- a backbone phylogeny  $T$



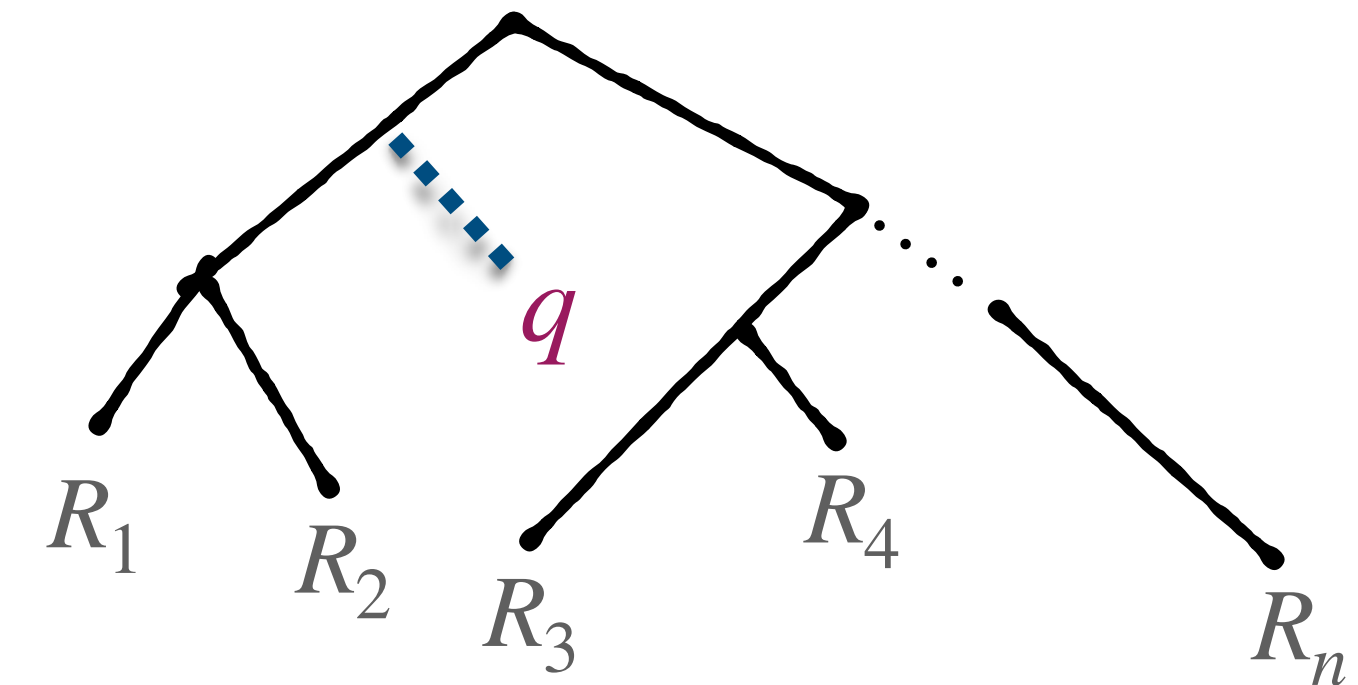


# Problem statement

Given:

- query sequence  $q$
- set of references  $\mathcal{R} = \{R_1, \dots, R_n\}$
- a backbone phylogeny  $T$

Phylogenetic placement:



Interpretation of the distances:

i) 
$$d(q, R) = \frac{\text{\# of mismatches}}{\text{length of } q}$$



ii) 
$$\mathbb{E}[d(q, R)] \approx 1 - \text{ANI}(Q, R)$$
  
 where  $Q$  is the source genome of  $q$

# **Outline of the method & sub-tasks we need to solve**

**krepp: k-mer-based read phylogenetic placement**

# Outline of the method & sub-tasks we need to solve

**krepp: k-mer-based read phylogenetic placement**

- I. Going beyond exact  $k$ -mer search & finding homologous  $k$ -mers

# Outline of the method & sub-tasks we need to solve

## **krepp: k-mer-based read phylogenetic placement**

- I. Going beyond exact  $k$ -mer search & finding homologous  $k$ -mers
- II. An efficient mapping between  $k$ -mers and reference genomes

# Outline of the method & sub-tasks we need to solve

## **krepp: k-mer-based read phylogenetic placement**

- I. Going beyond exact  $k$ -mer search & finding homologous  $k$ -mers
- II. An efficient mapping between  $k$ -mers and reference genomes
- ◆ III. Modeling Hamming distances of  $k$ -mer matches to estimate distances

# Outline of the method & sub-tasks we need to solve

## **krepp: k-mer-based read phylogenetic placement**

- I. Going beyond exact  $k$ -mer search & finding homologous  $k$ -mers
- II. An efficient mapping between  $k$ -mers and reference genomes
- ◆ III. Modeling Hamming distances of  $k$ -mer matches to estimate distances
- IV. Finding a placement on a reference tree that best explains distances

# **Computing Hamming distances w/ locality sensitive hashing**

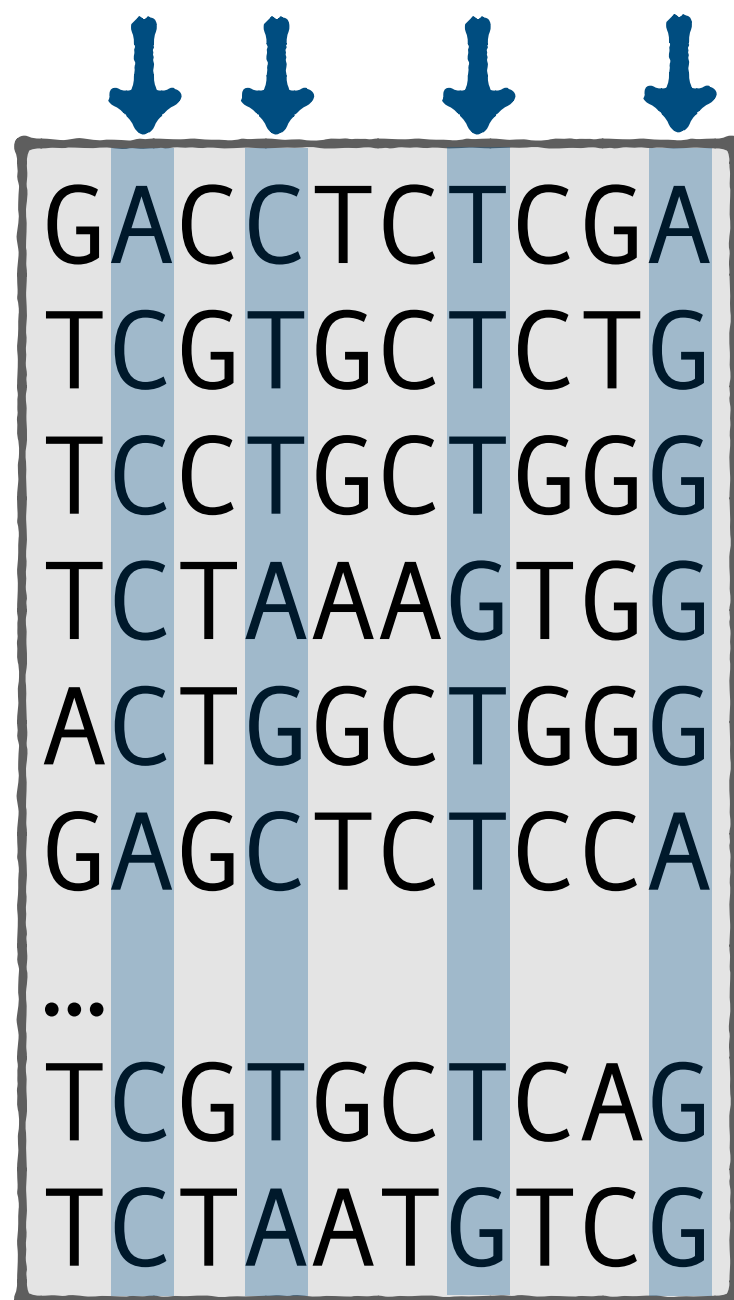
# Computing Hamming distances w/ locality sensitive hashing

```
GACCTCTCGA
TCGTGCTCTG
TCCTGCTGGG
TCTAAAGTGG
ACTGGCTGGG
GAGCTCTCCA
...
TCGTGCTCAG
TCTAATGTCG
```

reference  $k$ -mers

# Computing Hamming distances w/ locality sensitive hashing

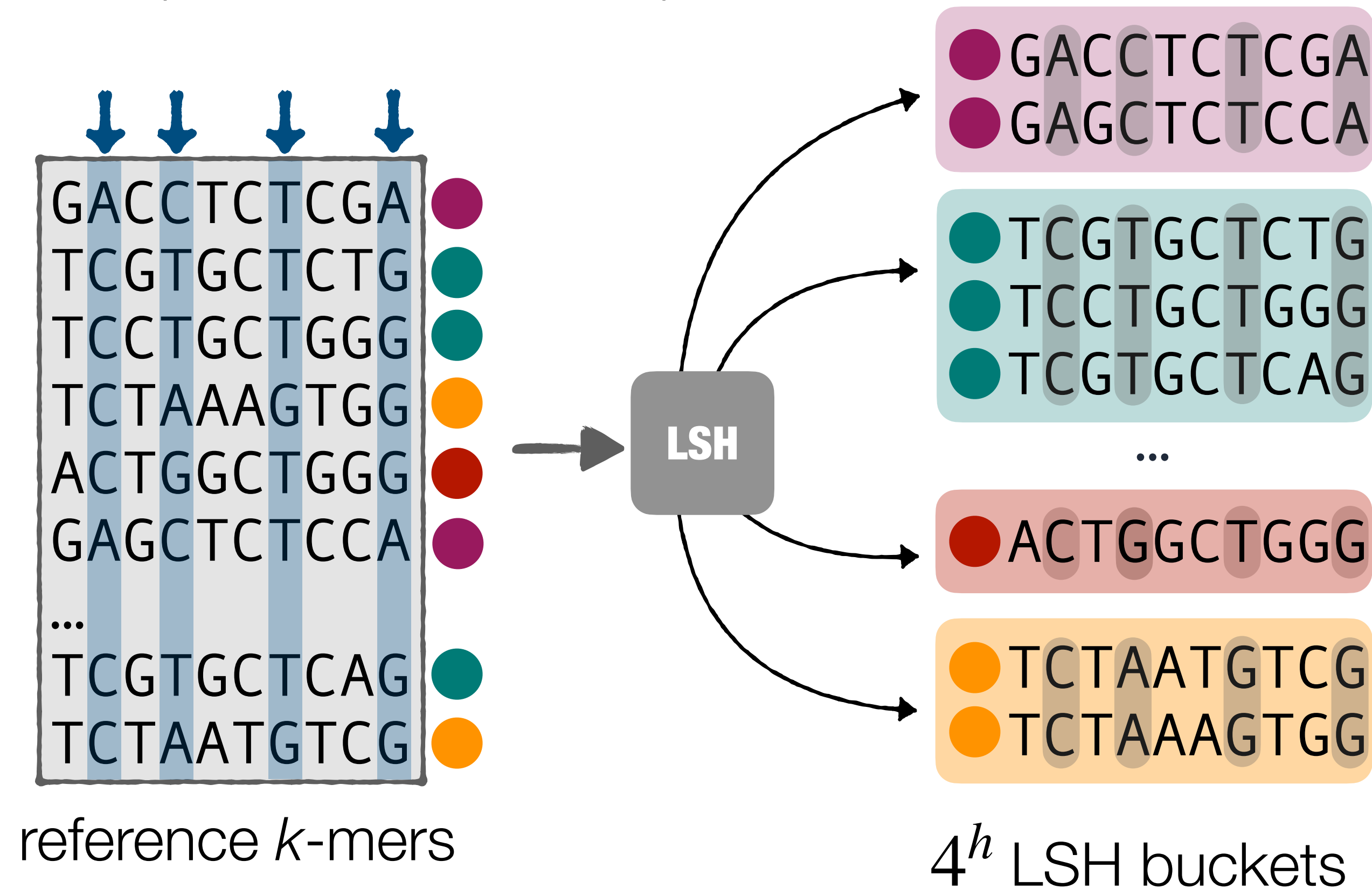
Select  $h$  random but fixed positions (default  $h$ : 14,  $k$ : 29)



reference  $k$ -mers

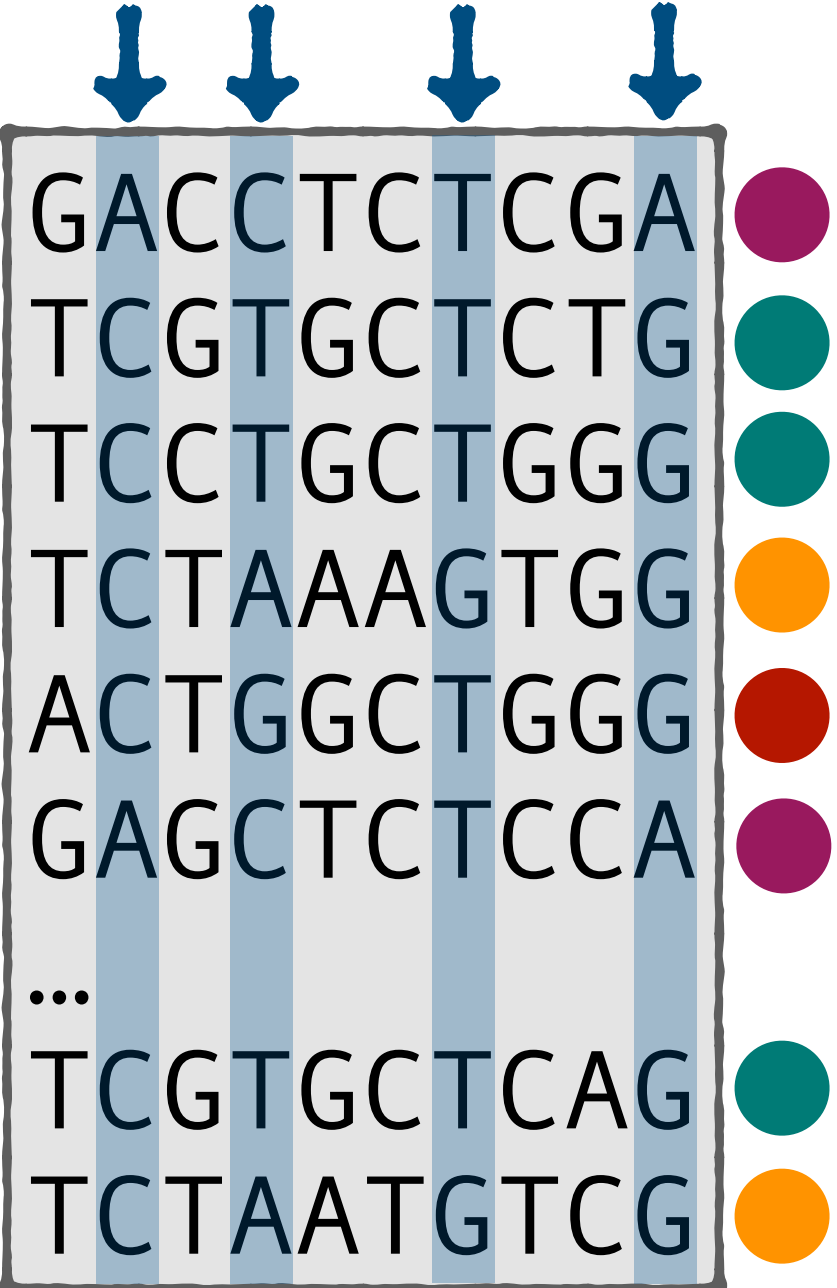
# Computing Hamming distances w/ locality sensitive hashing

Select  $h$  random but fixed positions (default  $h$ : 14,  $k$ : 29)

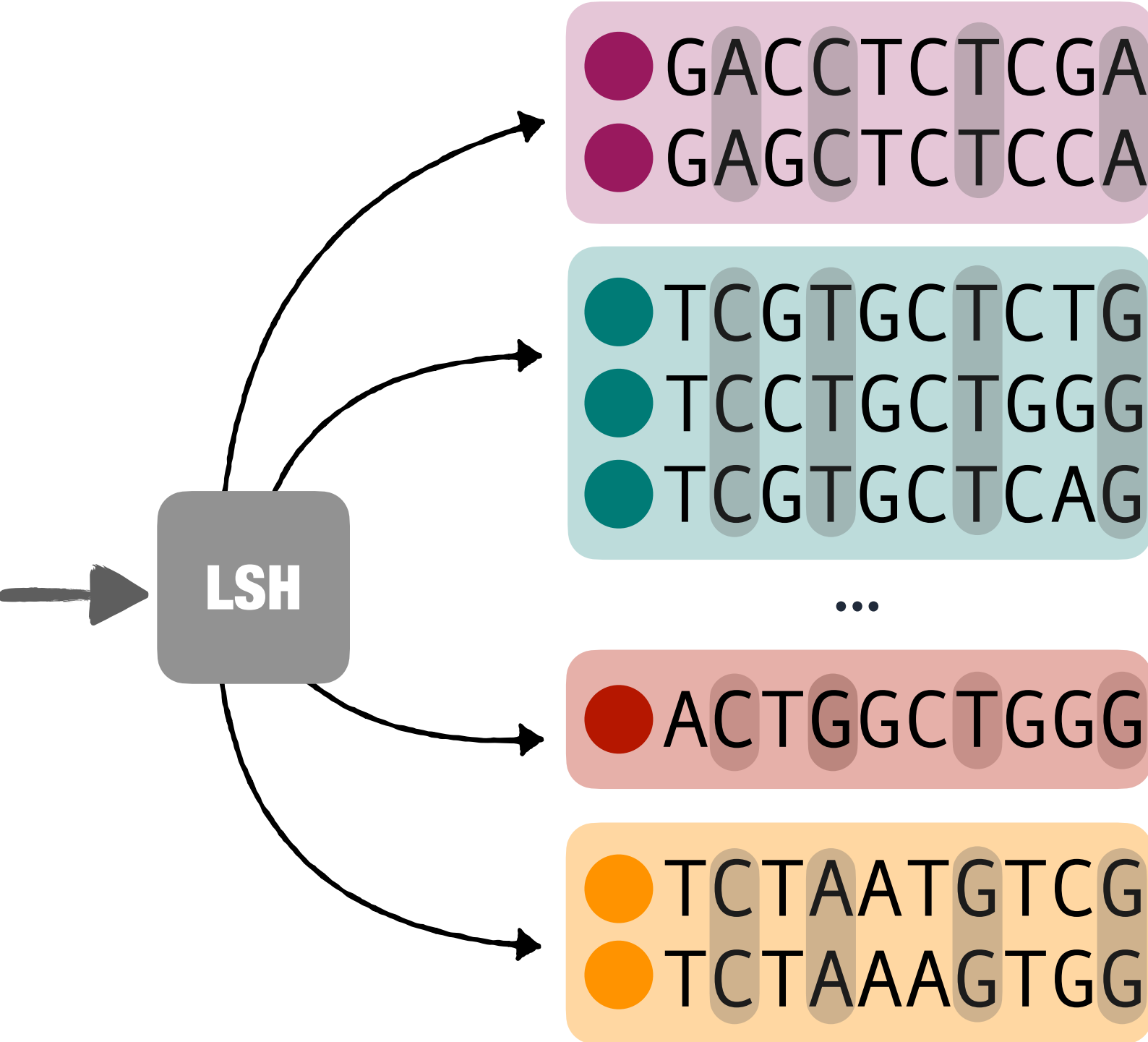


# Computing Hamming distances w/ locality sensitive hashing

Select  $h$  random but fixed positions (default  $h$ : 14,  $k$ : 29)



reference  $k$ -mers

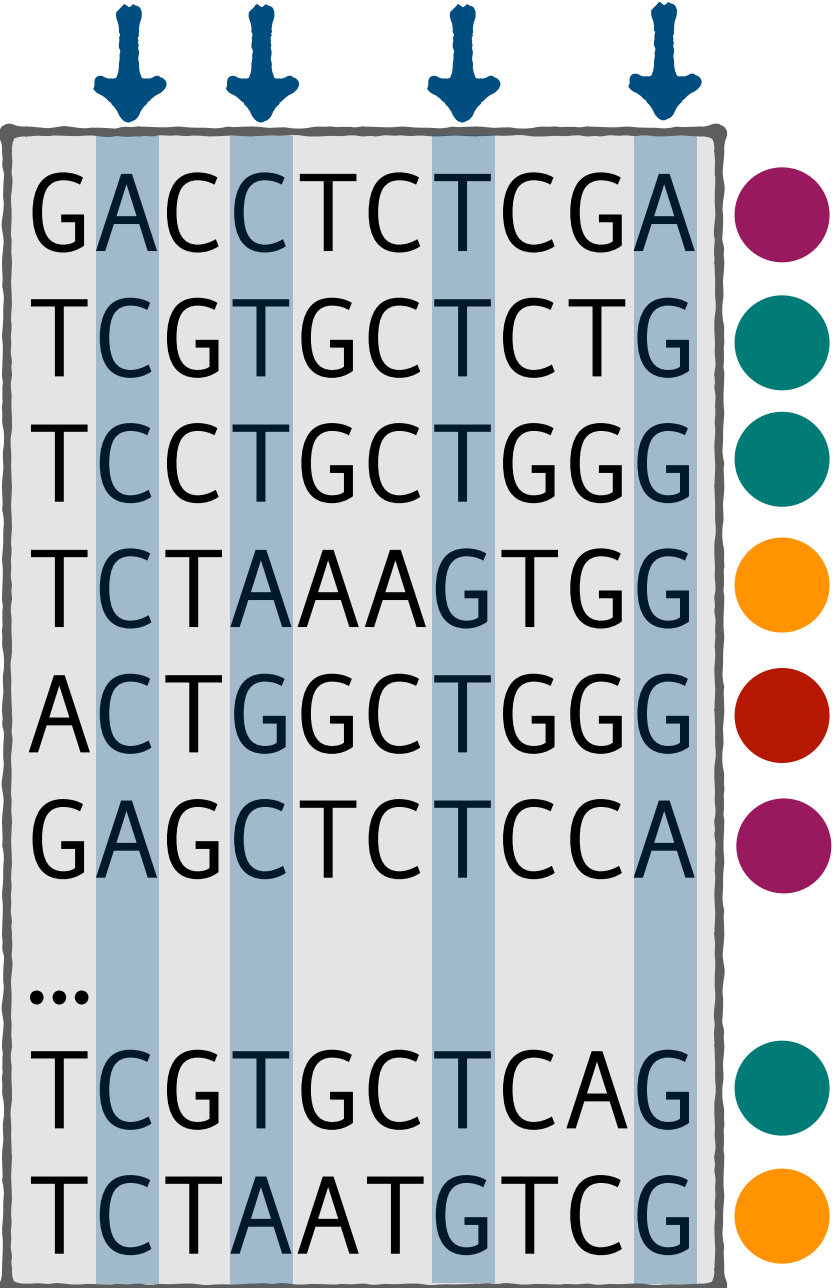


$4^h$  LSH buckets

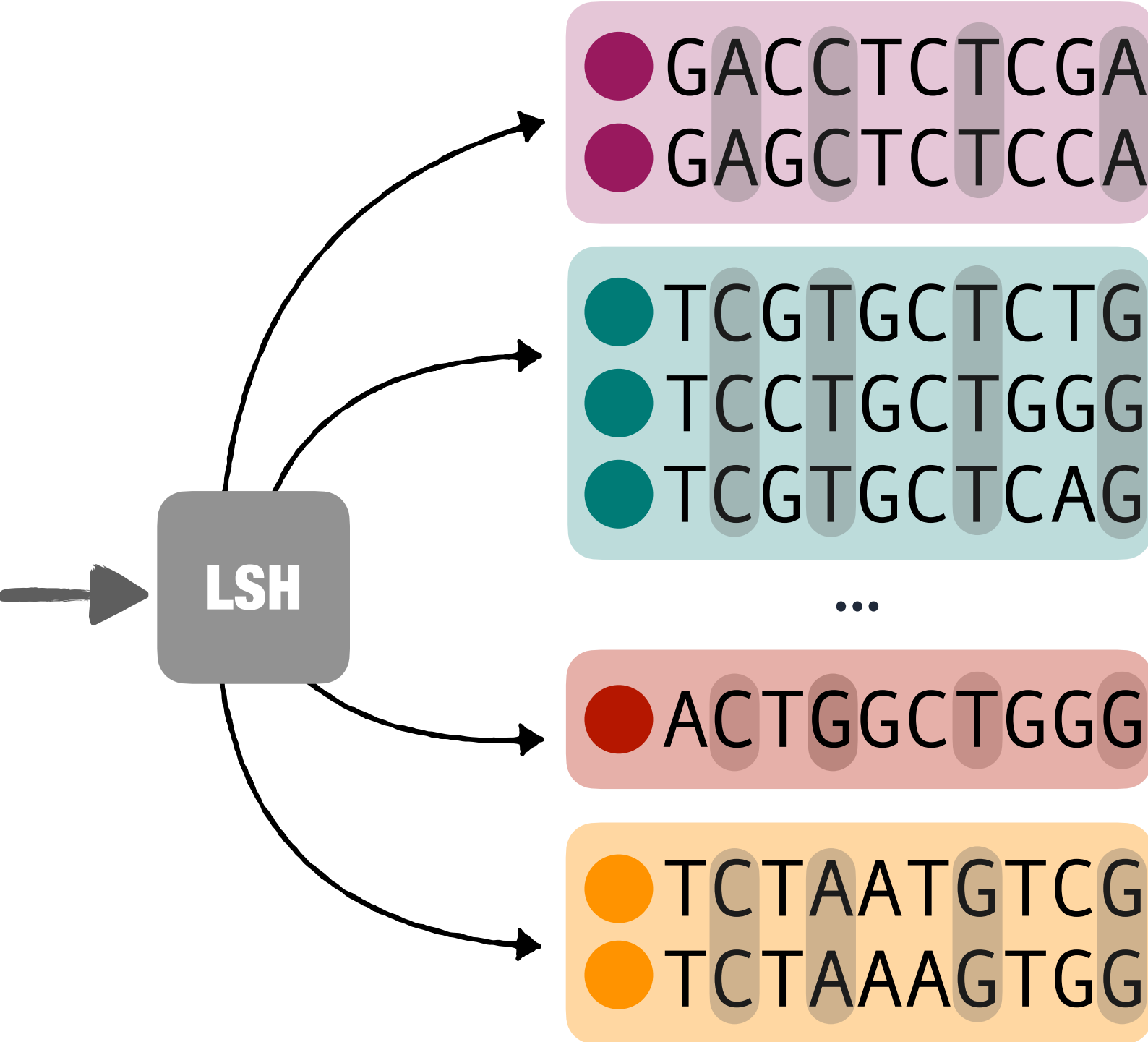
Given a query  $k$ -mer  
ACCTGCTGGG

# Computing Hamming distances w/ locality sensitive hashing

Select  $h$  random but fixed positions (default  $h$ : 14,  $k$ : 29)



reference  $k$ -mers



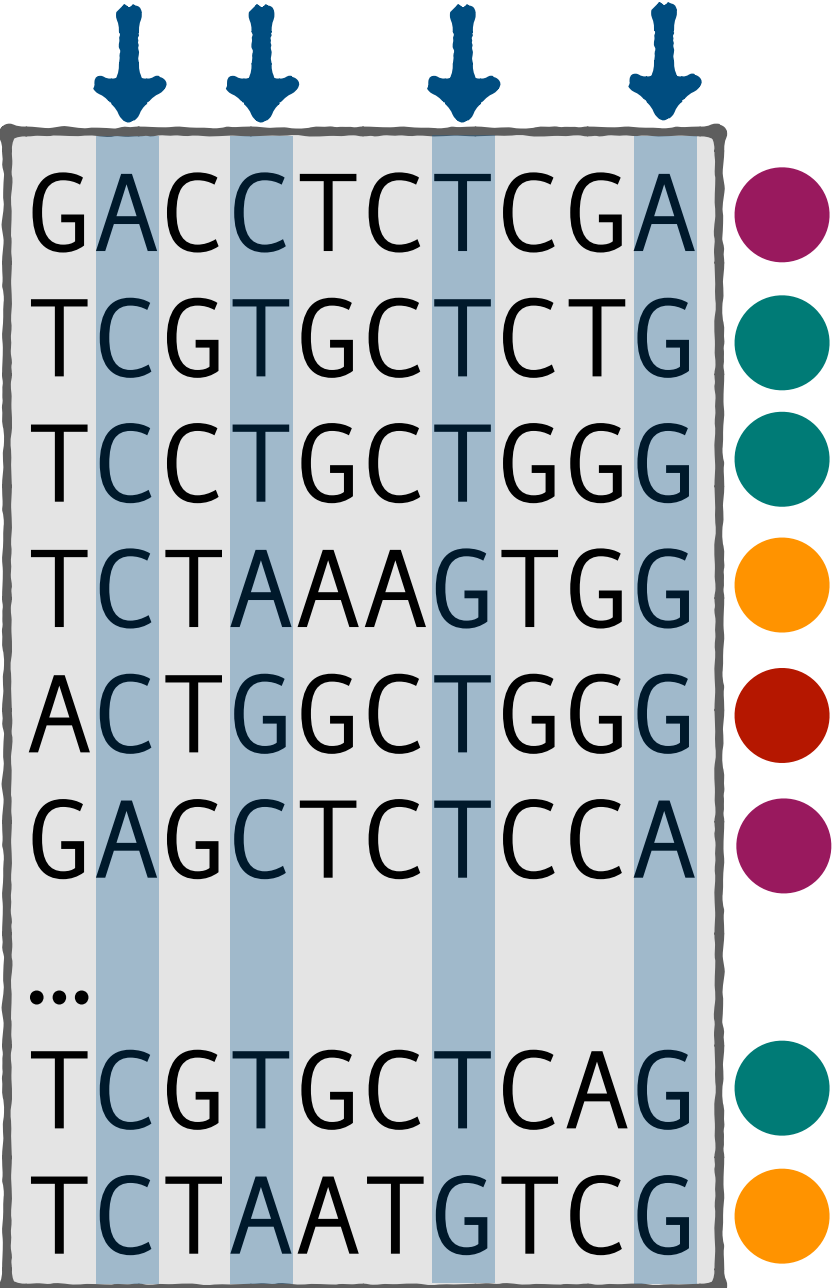
$4^h$  LSH buckets

Given a query  $k$ -mer

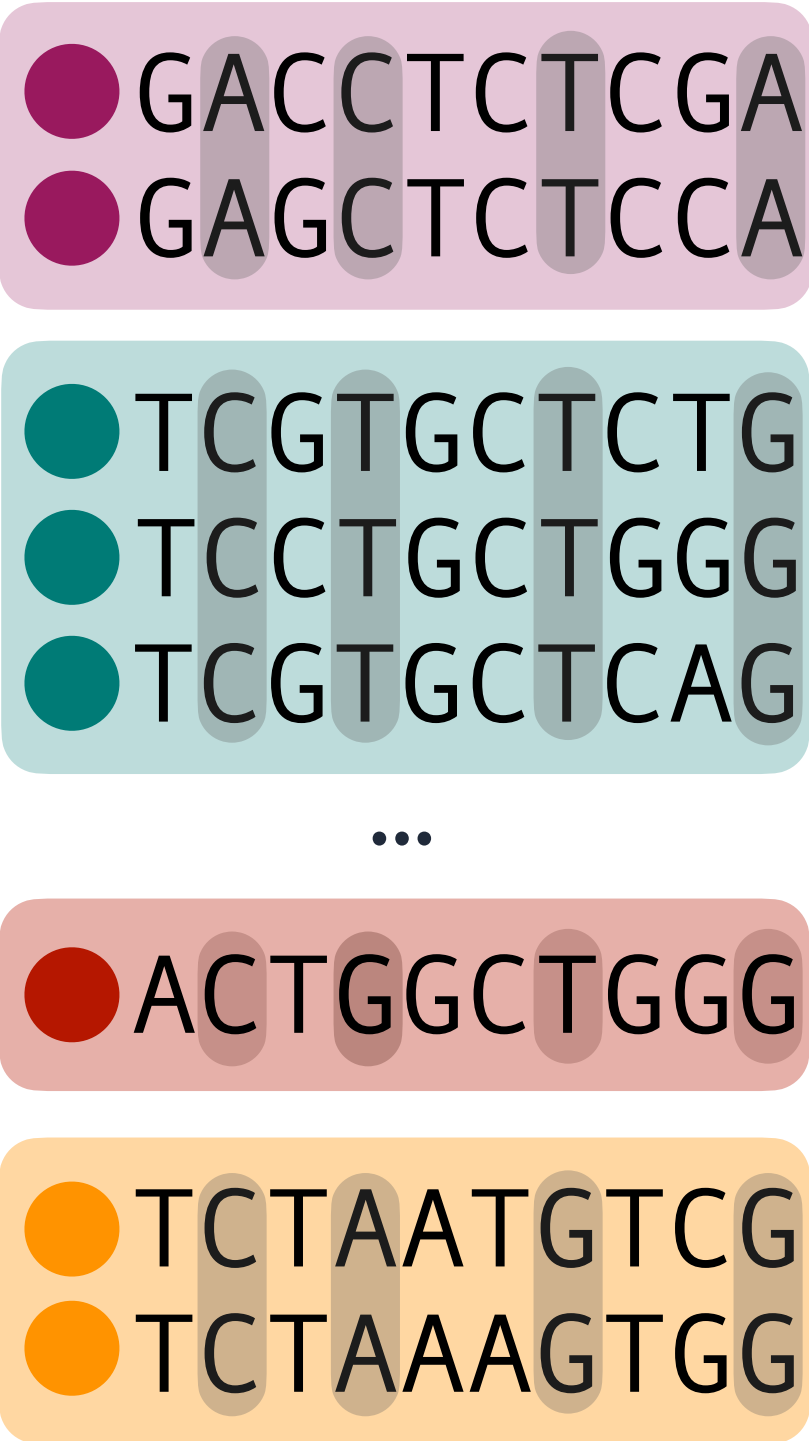
ACCTGCTGGG

# Computing Hamming distances w/ locality sensitive hashing

Select  $h$  random but fixed positions (default  $h$ : 14,  $k$ : 29)



LSH



**HD**

4  
1  
4

? FN at  
HD=2

Given a query  $k$ -mer

ACCTGCTGGG

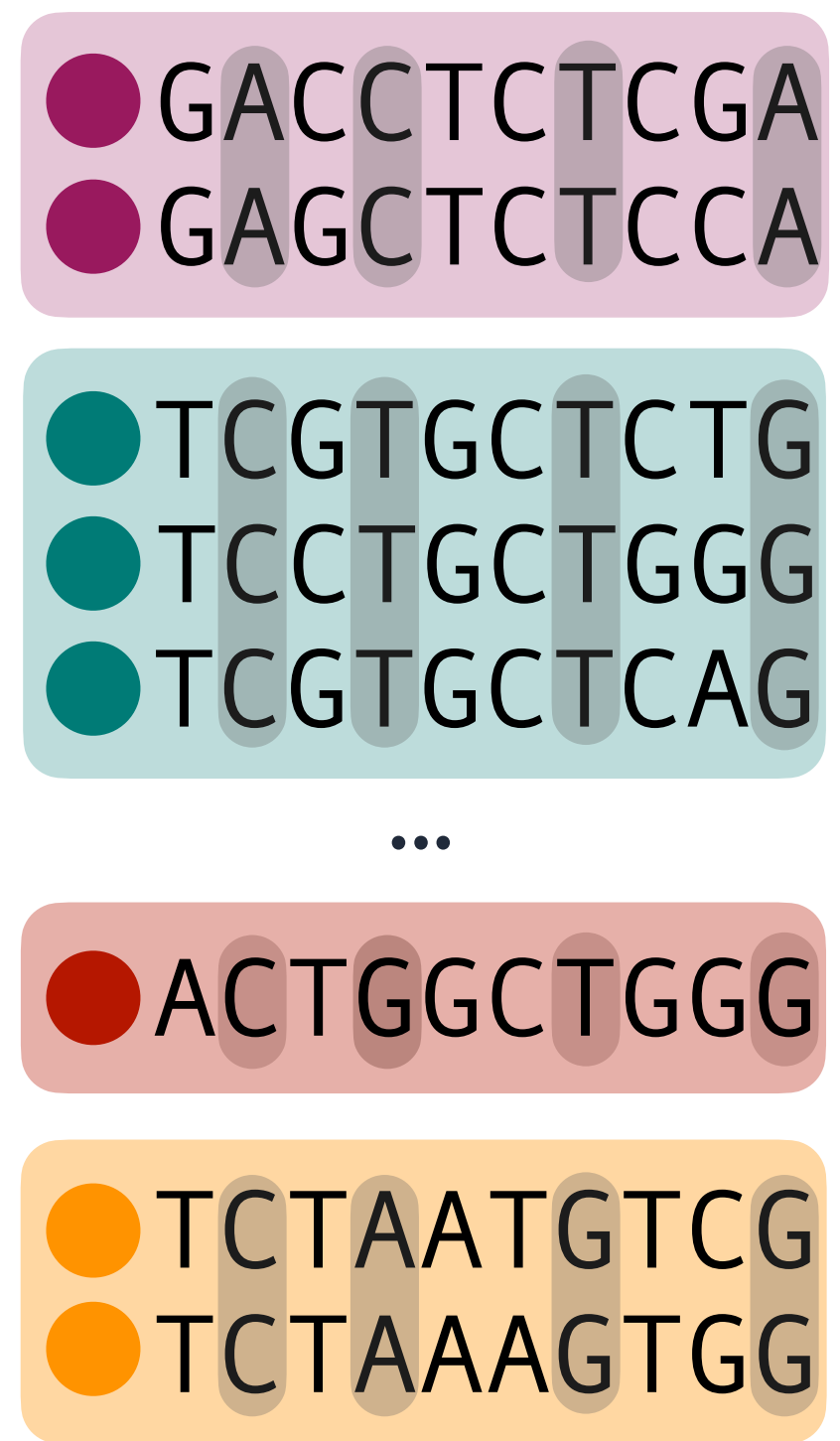


# Computing Hamming distances w/ locality sensitive hashing

Select  $h$  random but fixed positions (default  $h$ : 14,  $k$ : 29)



LSH

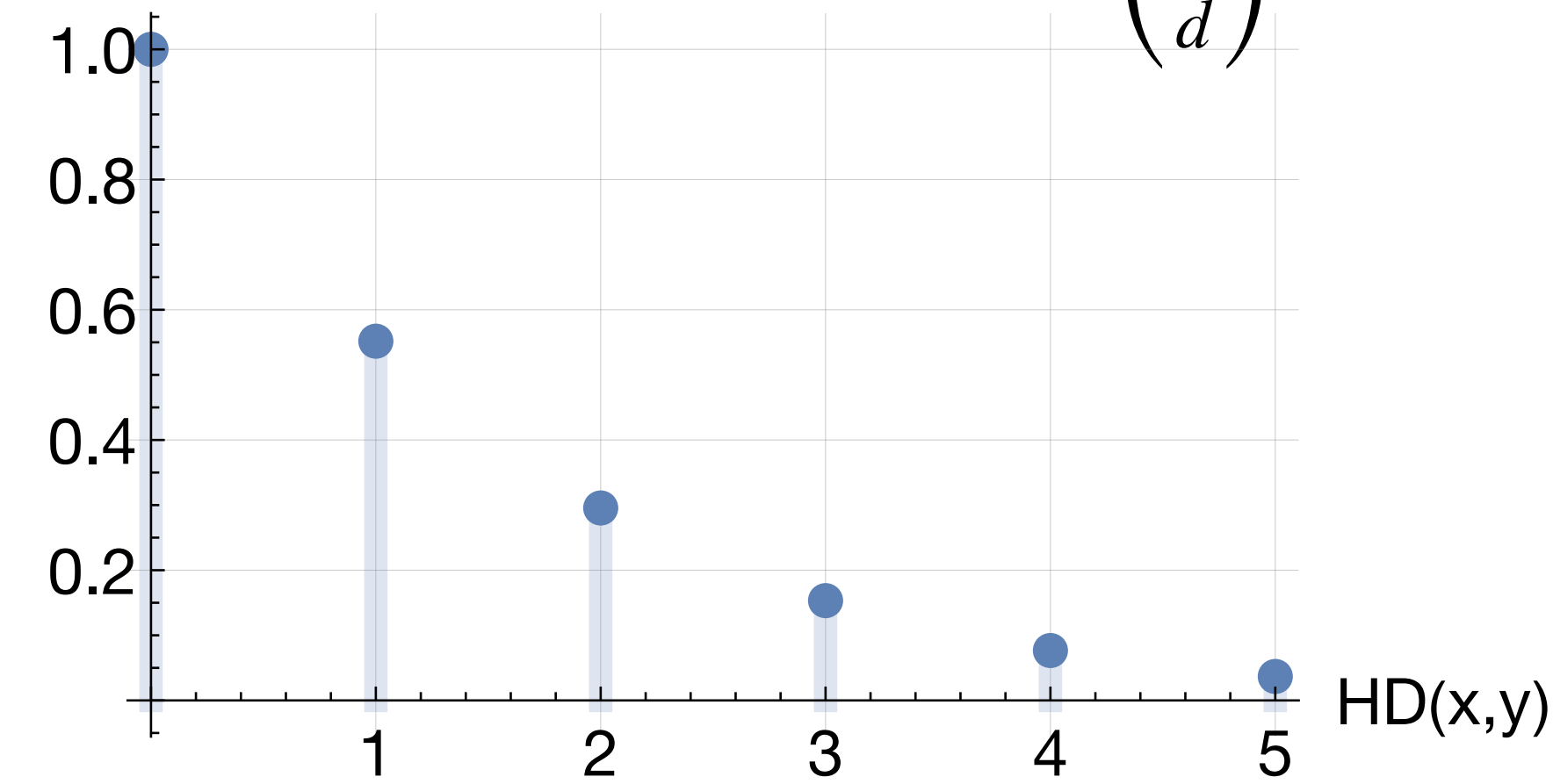


Given a query  $k$ -mer

ACCTGCTGGG

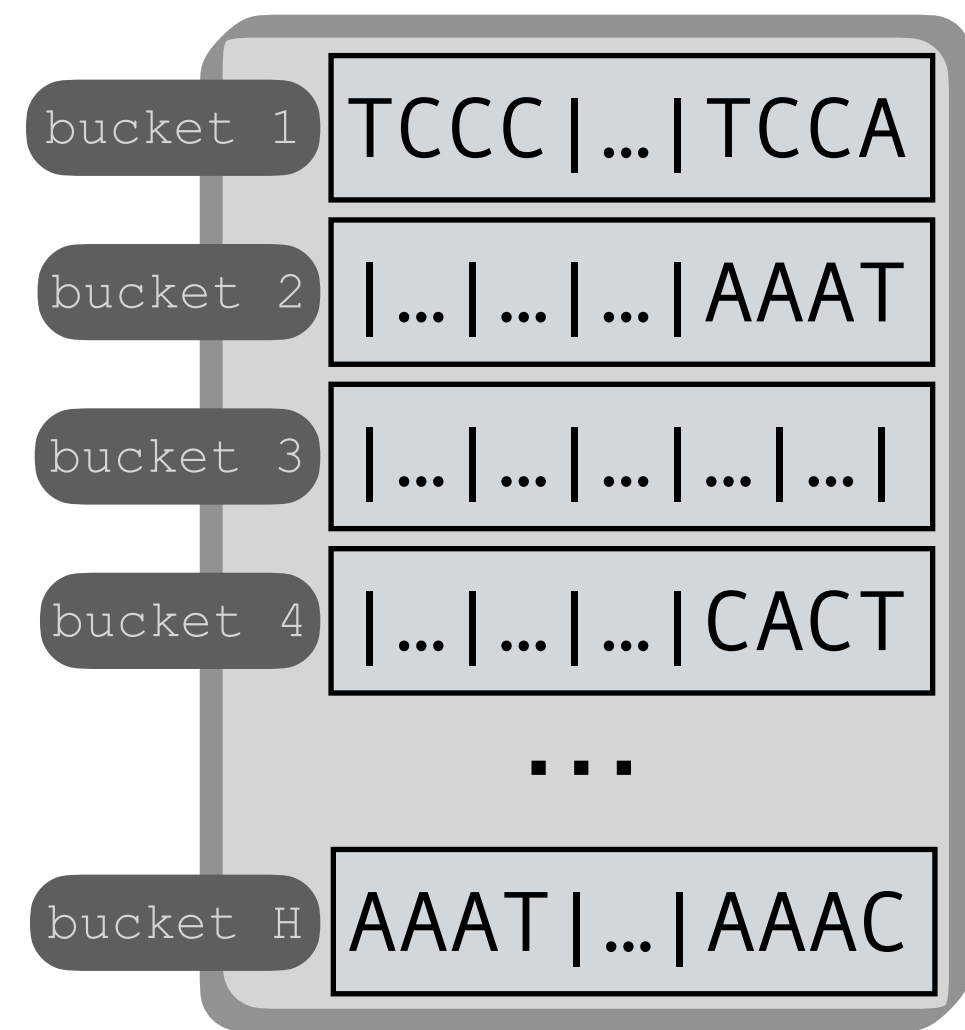
match at HD= $d$  with probability:  $\frac{\binom{k-h}{d}}{\binom{k}{d}}$

$P[\text{LSH}(x)=\text{LSH}(y)]$



# Mapping indexed k-mers to reference genomes

# Mapping indexed k-mers to reference genomes



LSH index

$R_1$

$R_2$

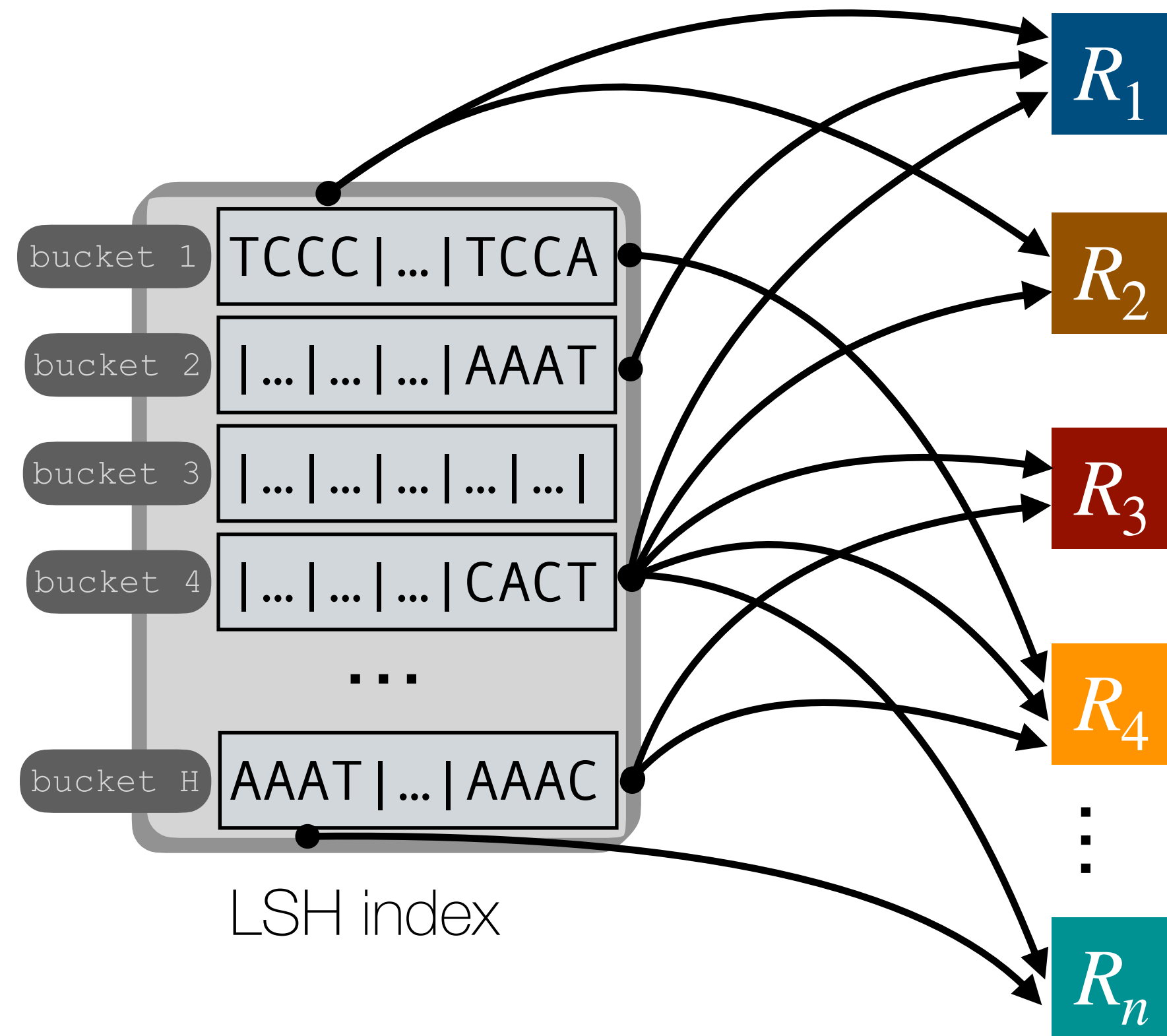
$R_3$

$R_4$

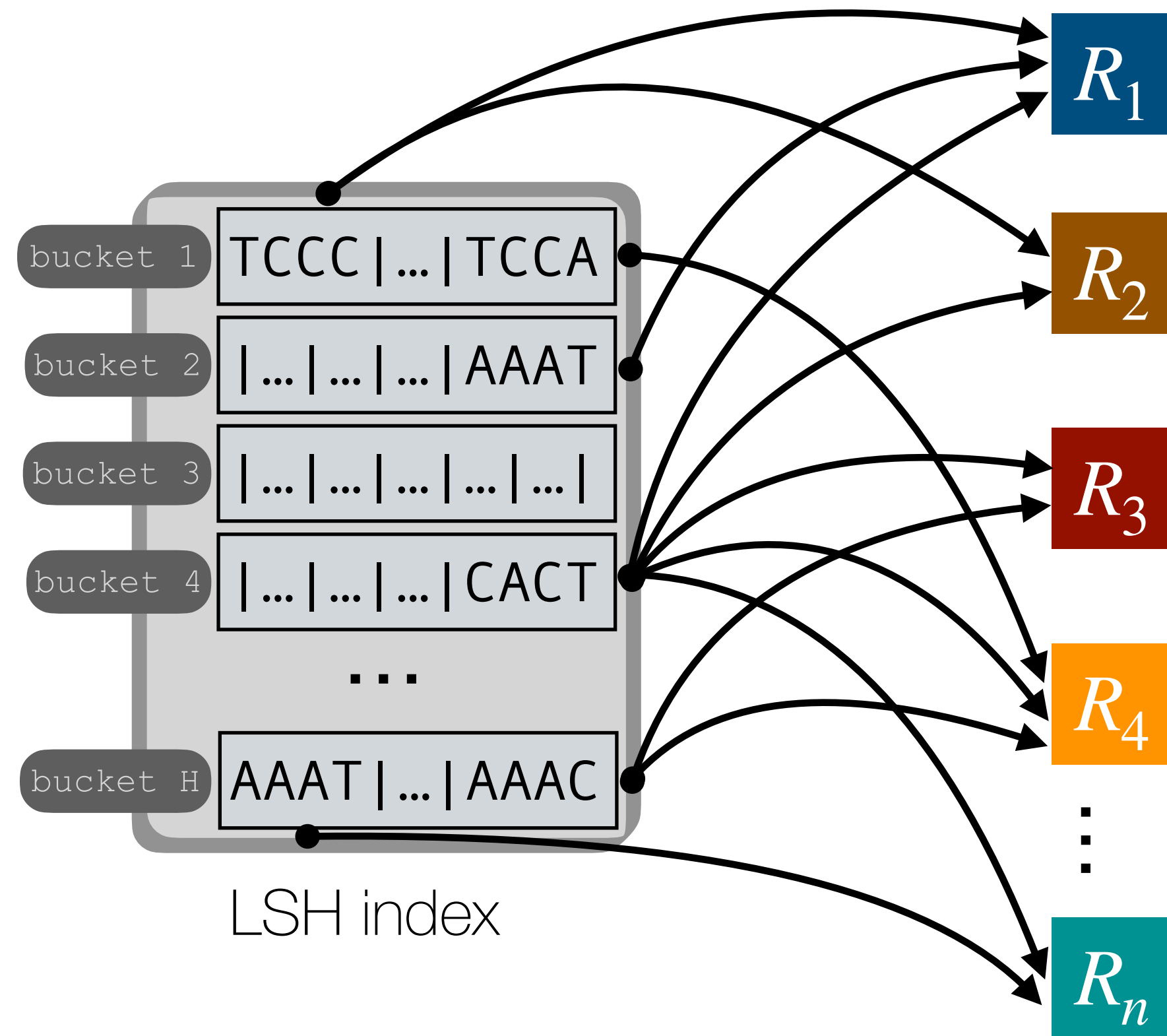
⋮

$R_n$

# Mapping indexed k-mers to reference genomes



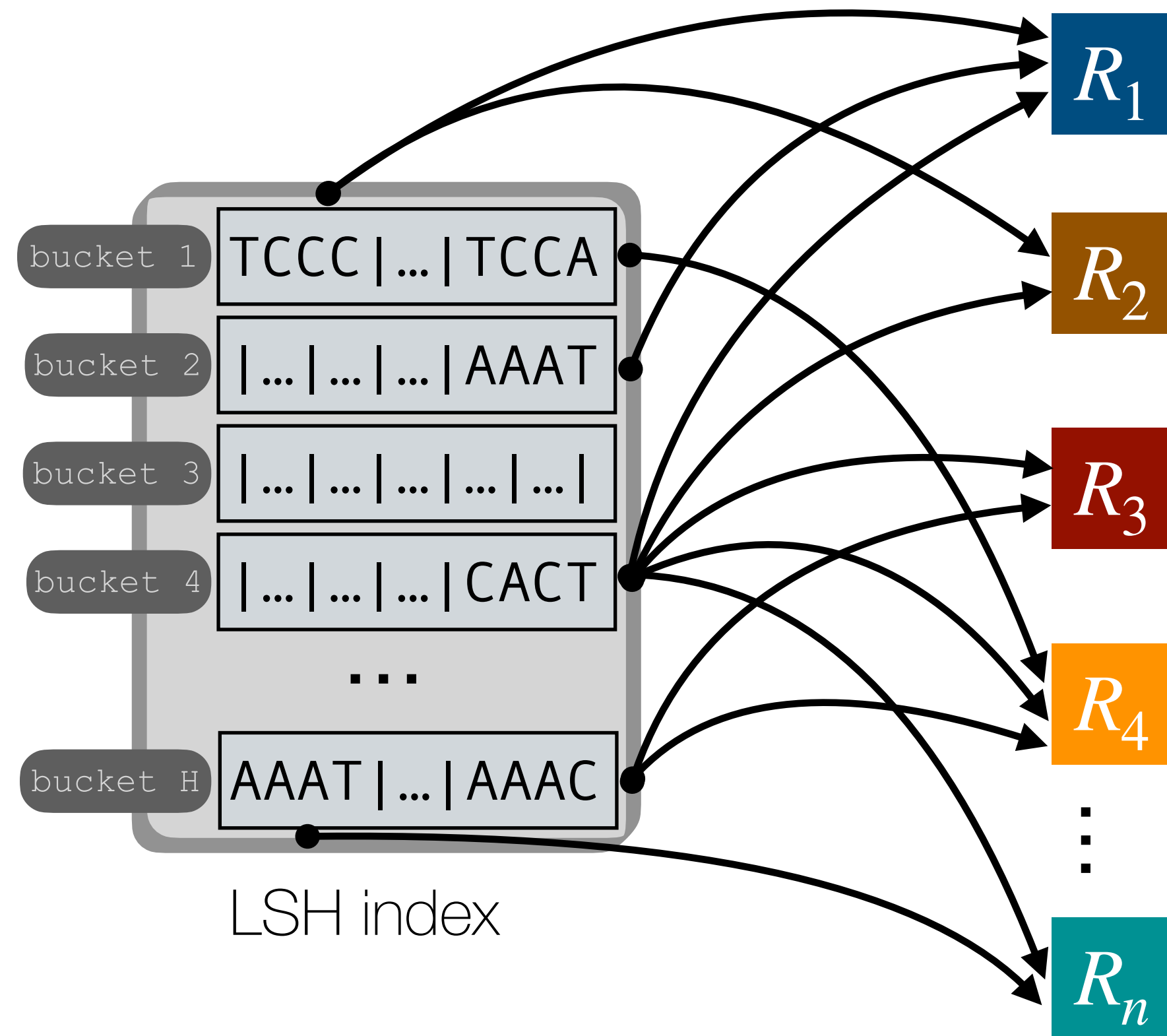
# Mapping indexed k-mers to reference genomes



well studied **colored k-mer** problem  
often represented as a sparse matrix

**color:** a subset of references  
(including singletons)

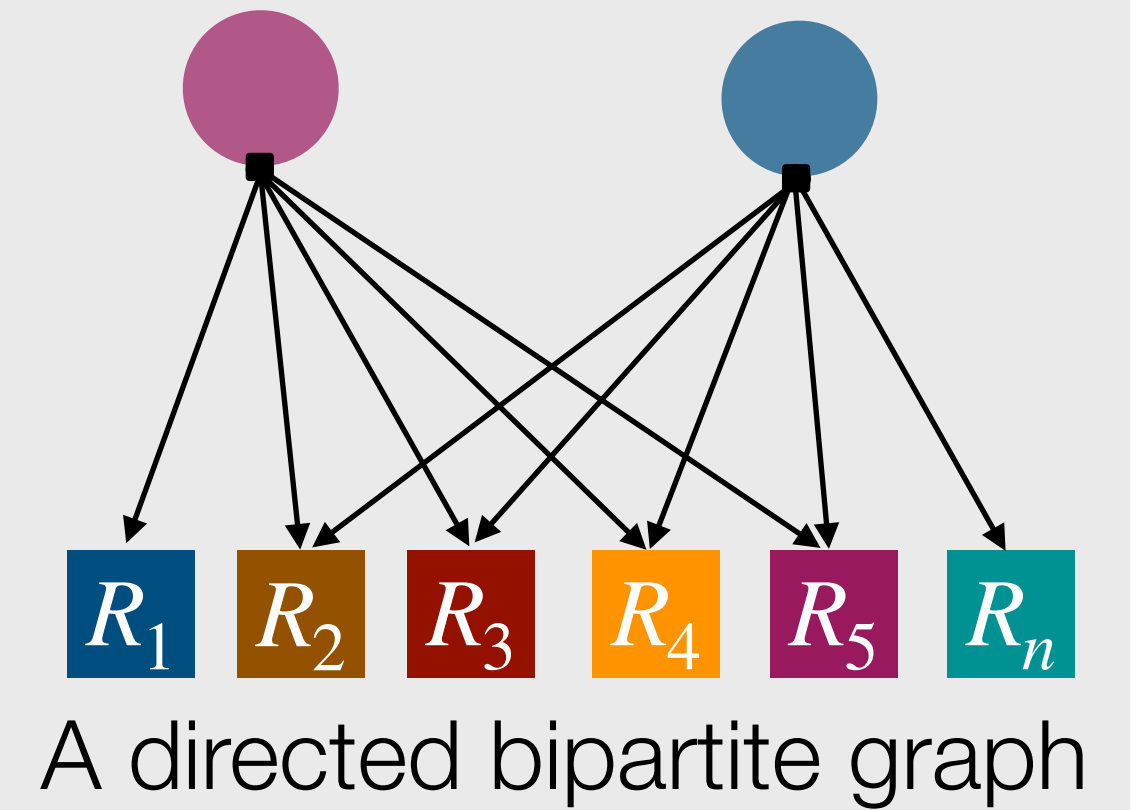
# Mapping indexed k-mers to reference genomes



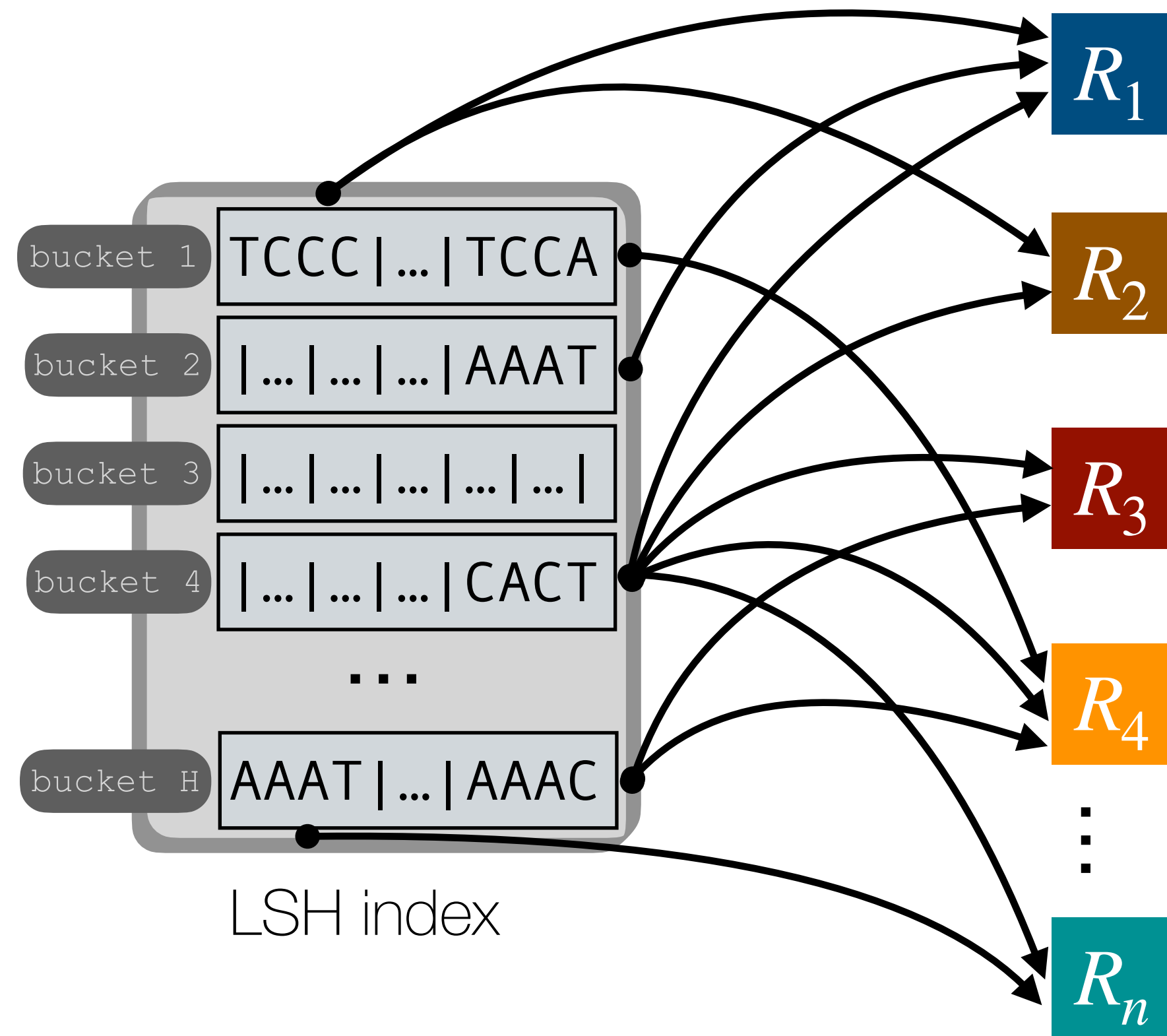
well studied **colored k-mer** problem  
often represented as a sparse matrix

**color:** a subset of references  
(including singletons)

$$|V| + |E| = 18$$



# Mapping indexed k-mers to reference genomes

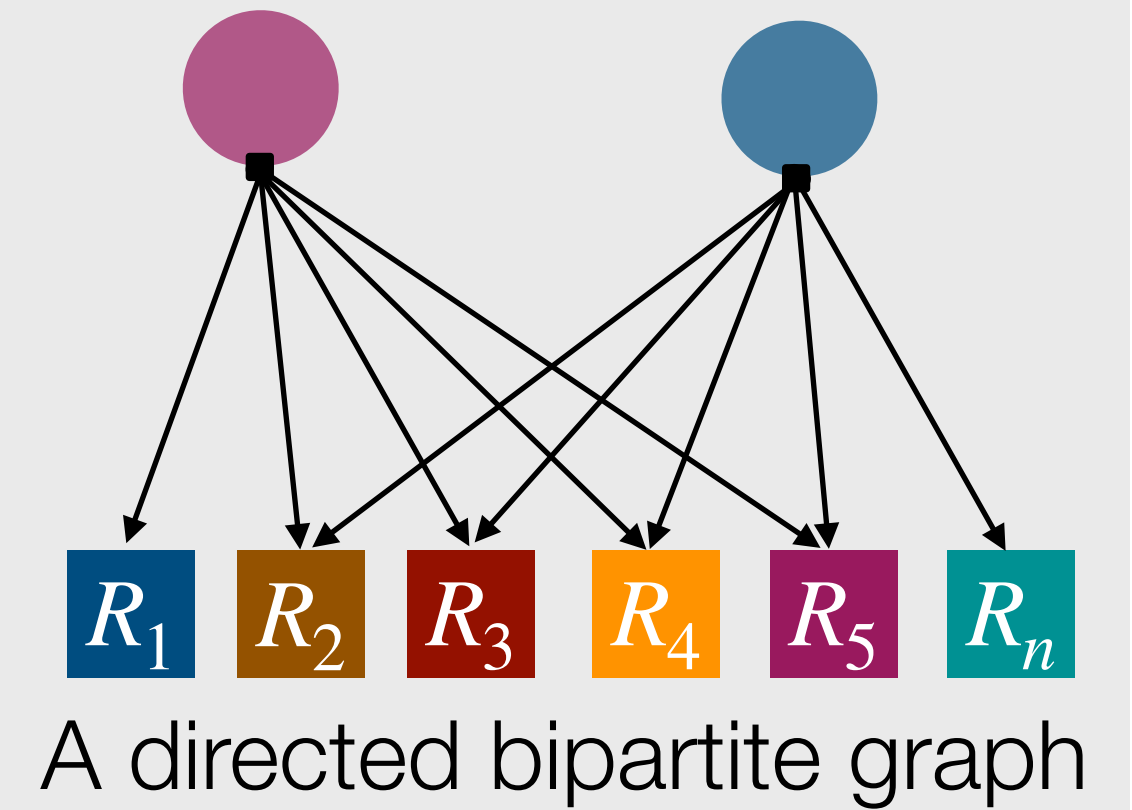


Minimize  $|V| + |E|$ ?

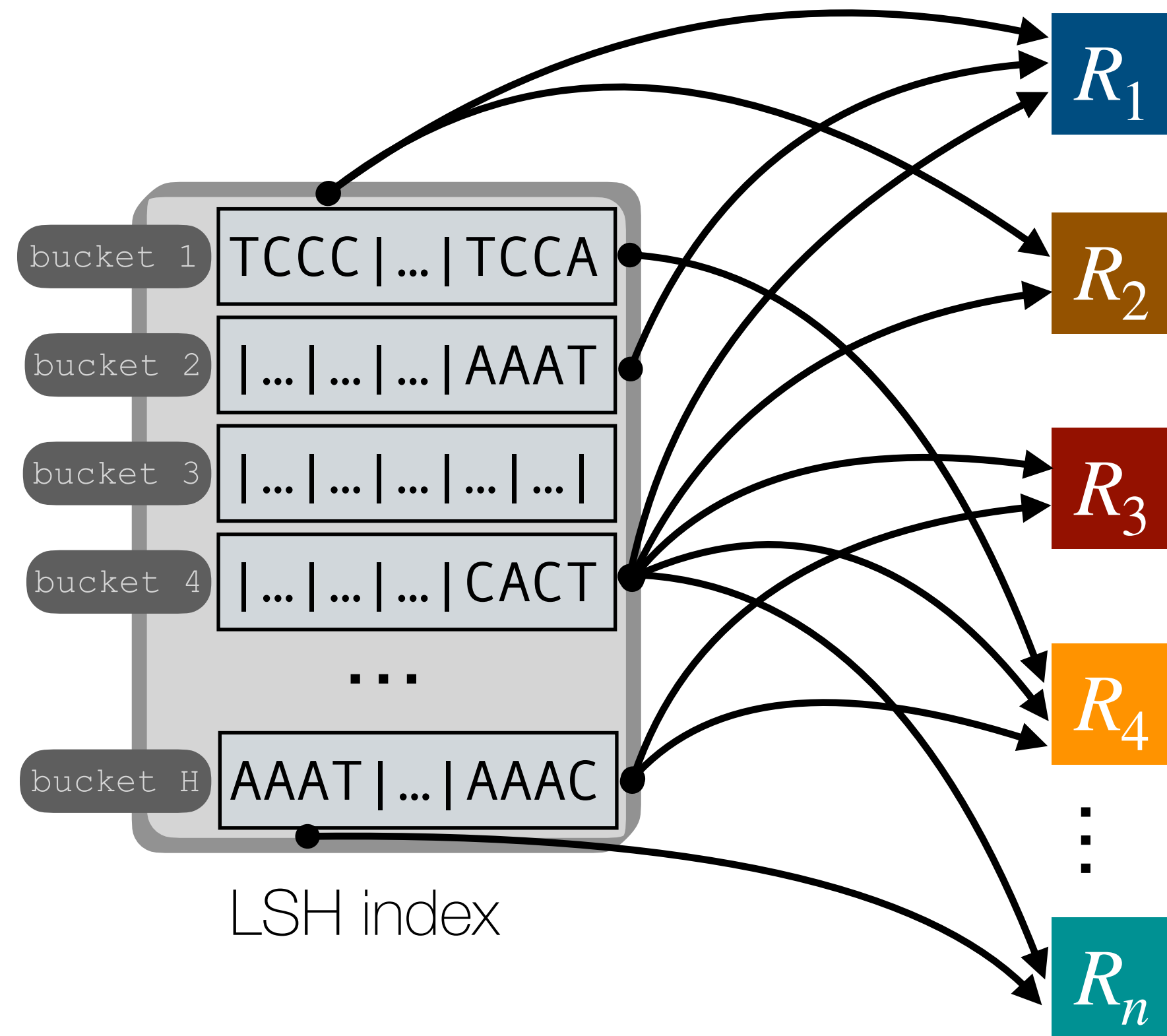
well studied **colored k-mer** problem  
often represented as a sparse matrix

**color:** a subset of references  
(including singletons)

$$|V| + |E| = 18$$



# Mapping indexed k-mers to reference genomes



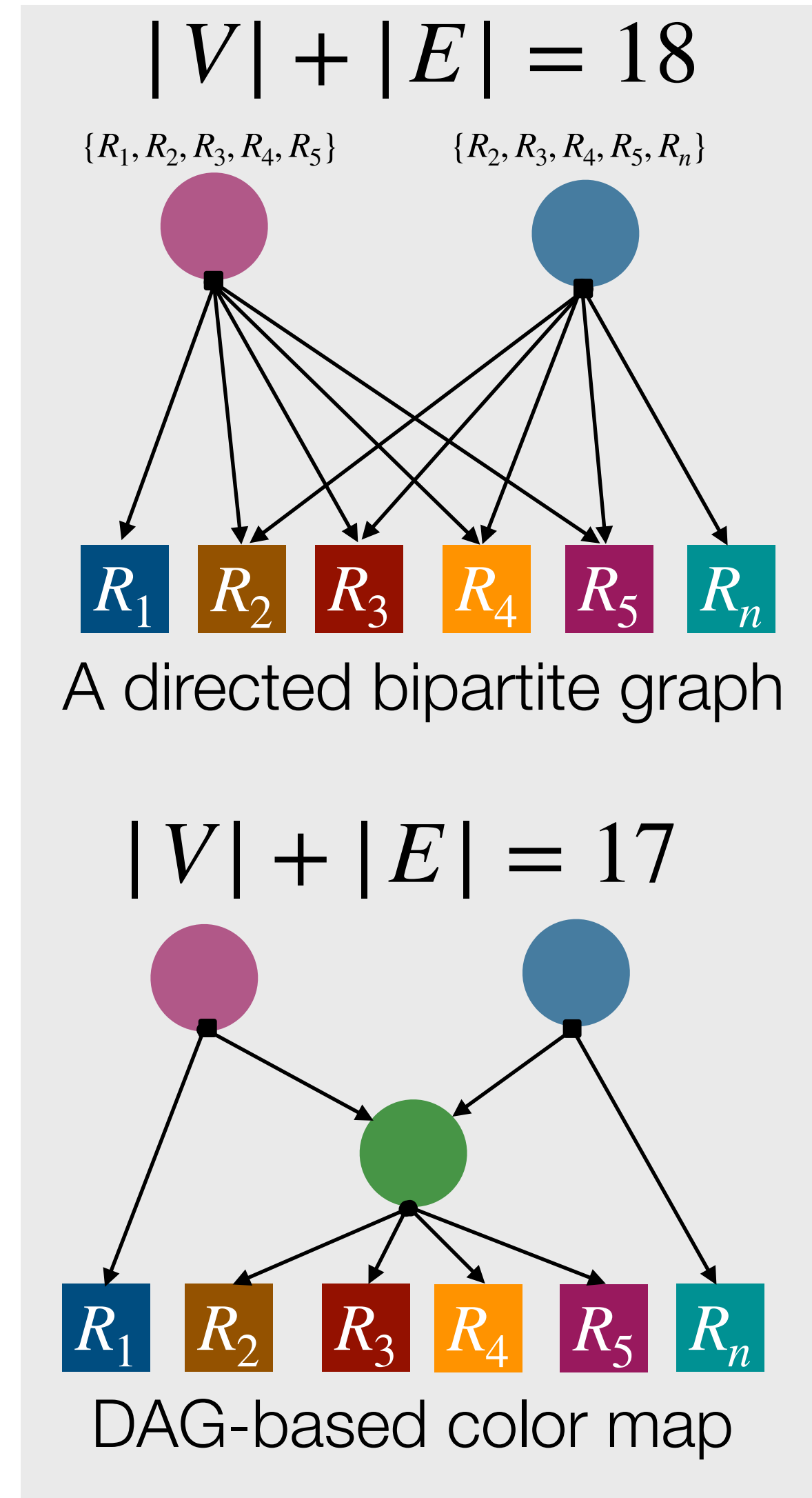
LSH index

well studied **colored k-mer** problem  
often represented as a sparse matrix

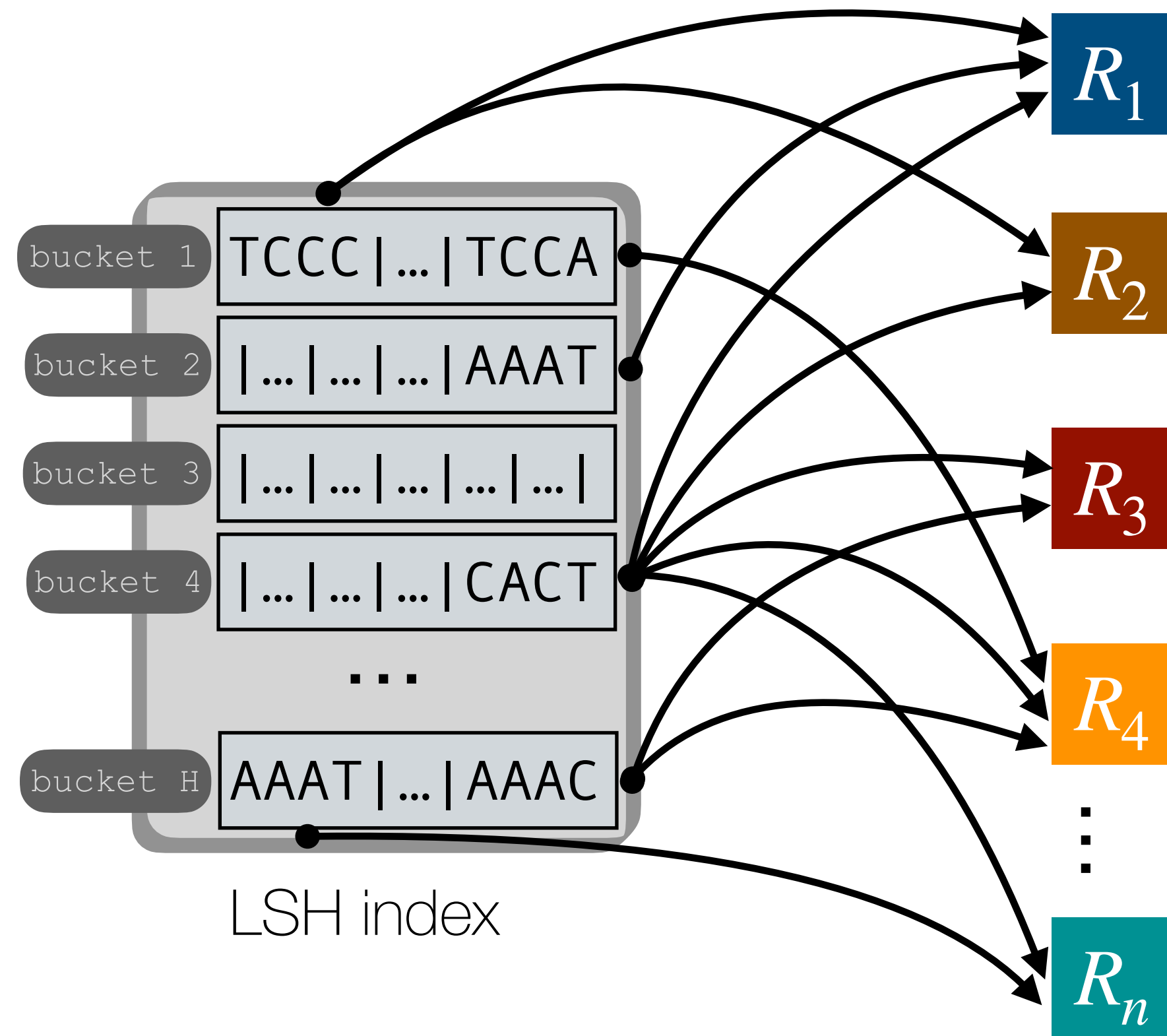
**color:** a subset of references  
(including singletons)

Minimize  $|V| + |E|$ ?

- i. reduce the size by adding nodes for frequently **shared subsets** (*meta-colors*)



# Mapping indexed k-mers to reference genomes

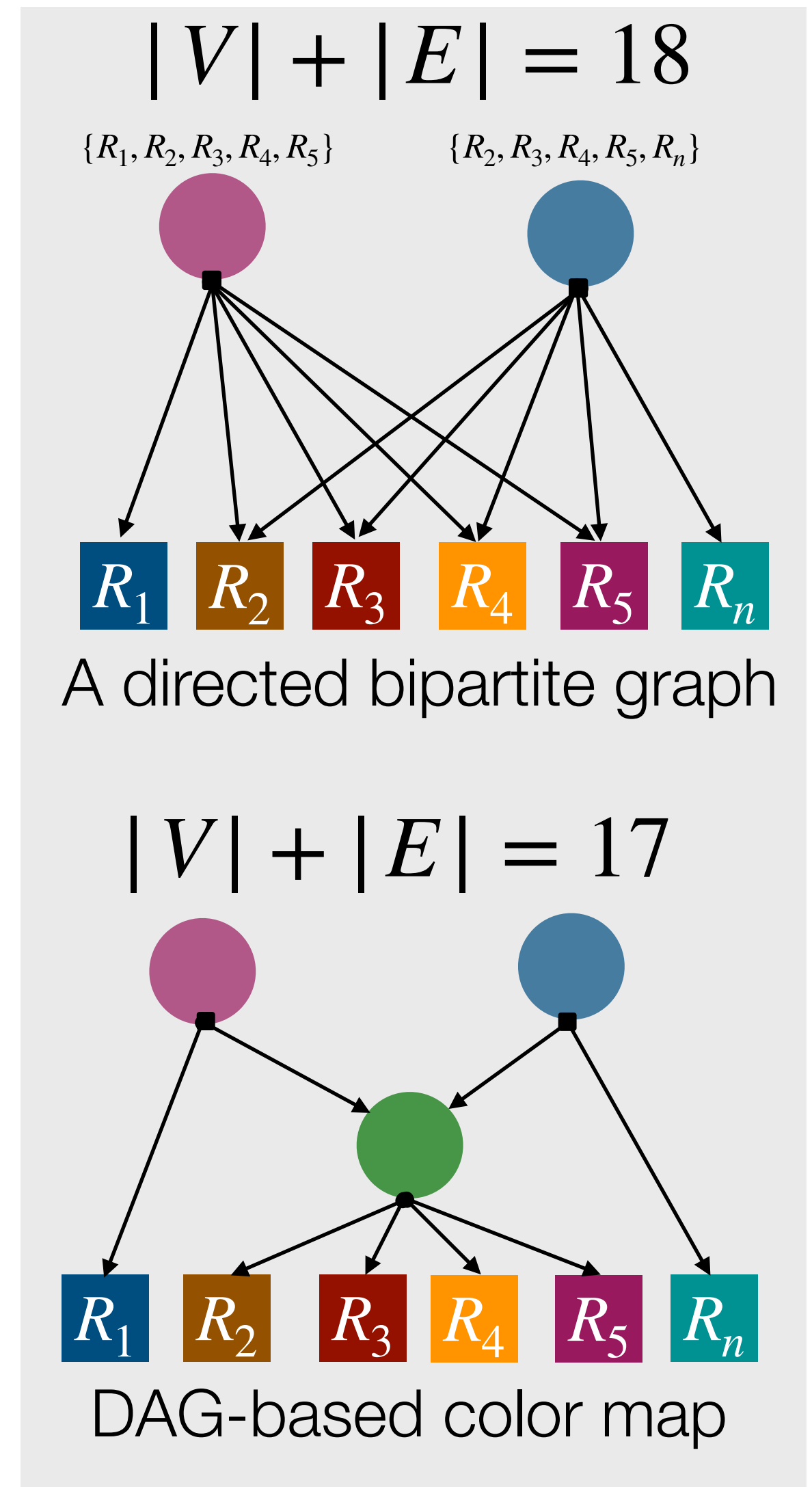


well studied **colored k-mer** problem  
often represented as a sparse matrix

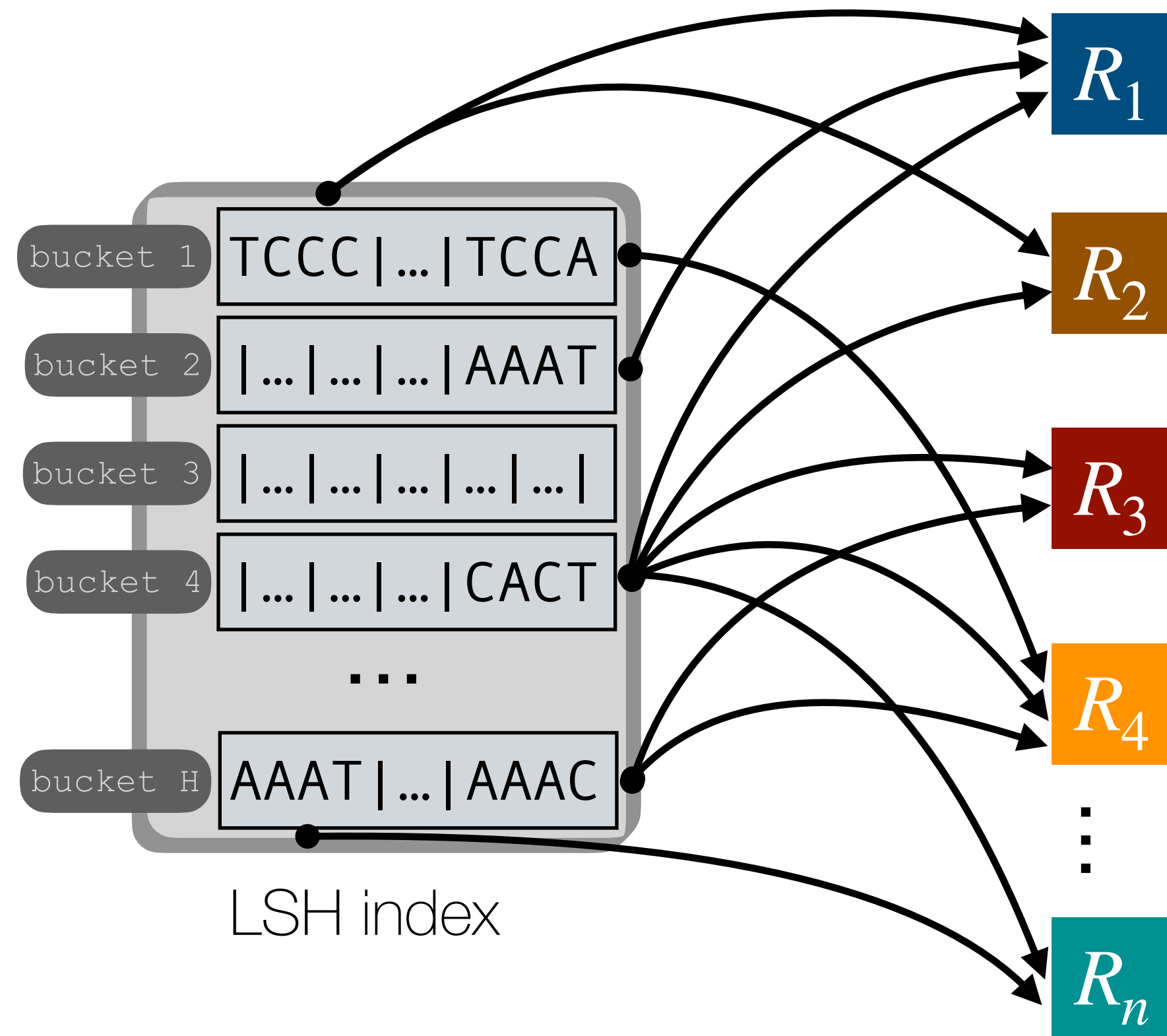
**color:** a subset of references  
(including singletons)

Minimize  $|V| + |E|$ ?

- reduce the size by adding nodes for frequently **shared subsets** (*meta-colors*)
- explain larger color w/ smaller colors



# Mapping indexed k-mers to reference genomes



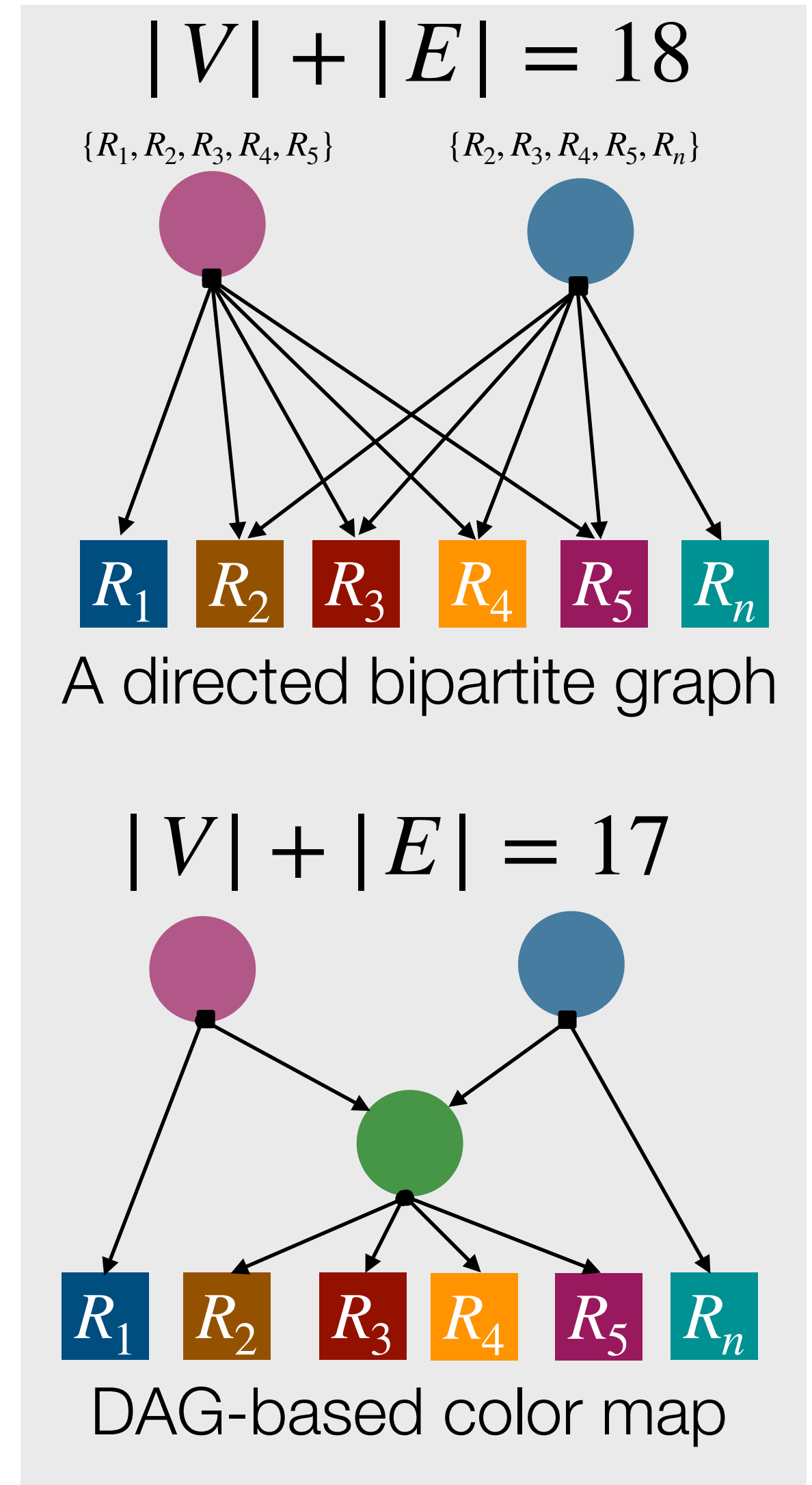
LSH index

well studied **colored k-mer** problem often represented as a sparse matrix

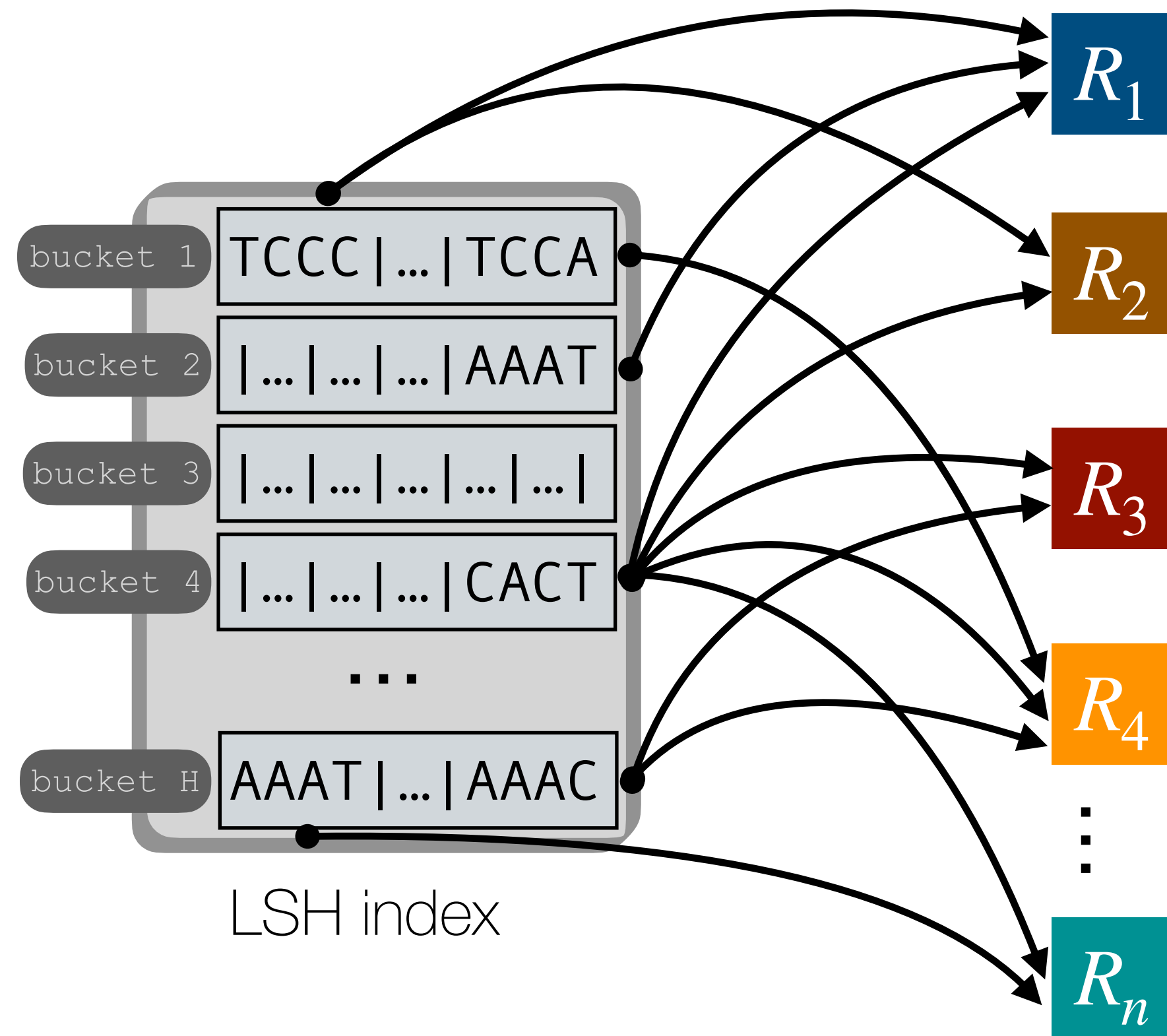
**color:** a subset of references (including singletons)

Minimize  $|V| + |E|$ ?

- i. reduce the size by adding nodes for frequently **shared subsets** (*meta-colors*)
- ii. explain larger color w/ smaller colors
- iii. use this **DAG** to **reconstruct colors** of matching *k*-mers during query time



# Mapping indexed k-mers to reference genomes



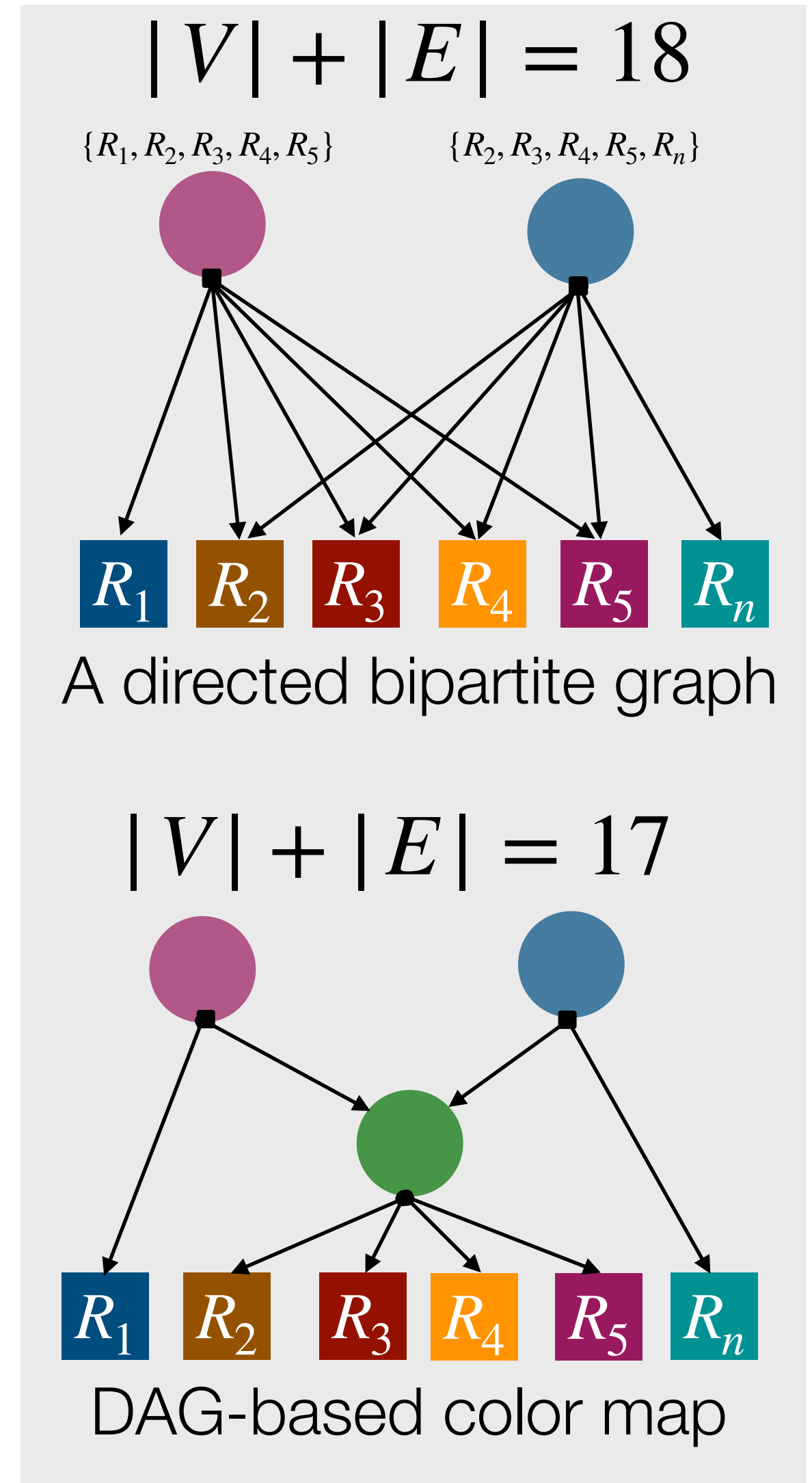
well studied **colored k-mer** problem often represented as a sparse matrix

**color:** a subset of references (including singletons)

Minimize  $|V| + |E|$ ?

- i. reduce the size by adding nodes for frequently **shared subsets** (*meta-colors*)
- ii. explain larger color w/ smaller colors
- iii. use this **DAG** to **reconstruct colors** of matching *k*-mers during query time

We use a heuristic guided by the phylogeny to build an efficient DAG.



# **Finding homologous k-mers of indexed references**

# Finding homologous k-mers of indexed references

Given a query sequence;

```
>ID_XYZ  
ATACCTAGGAGTACGGGAC
```

# Finding homologous k-mers of indexed references

Given a query sequence;

```
>ID_XYZ  
ATACCTAGGAGTACGGGAC
```

```
1: ATAC  
2:  TACC  
3:  ACCT  
4:  CCTA  
5:  CTAG  
   ...  
L-k+1:          GGAC
```

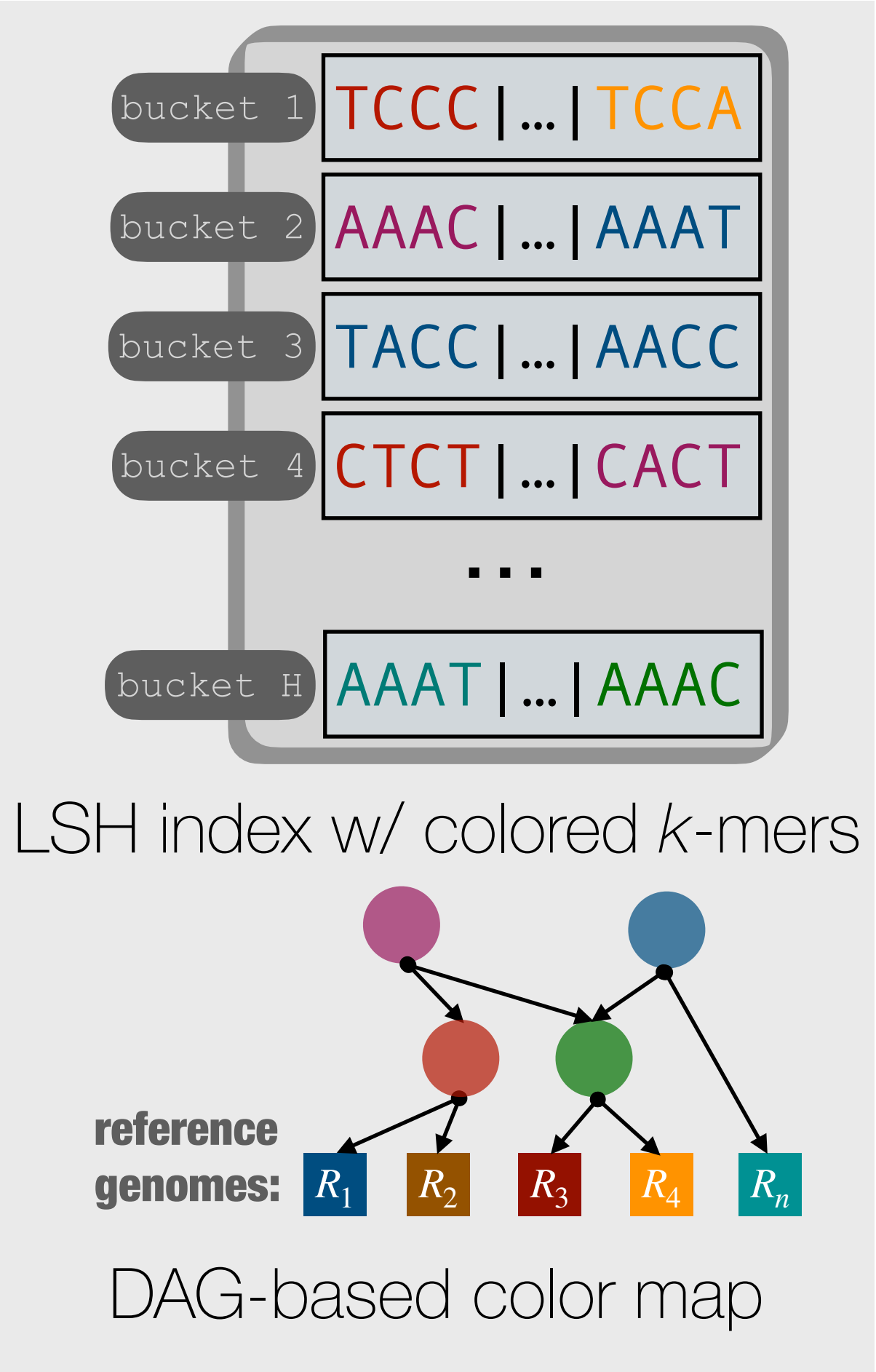
# Finding homologous k-mers of indexed references

Given a query sequence;

```
>ID_XYZ  
ATACCTAGGAGTACGGGAC
```

```
1: ATAC  
2:  TACC  
3:   ACCT  
4:    CCTA  
5:     CTAG  
...  
L-k+1: GGAC
```

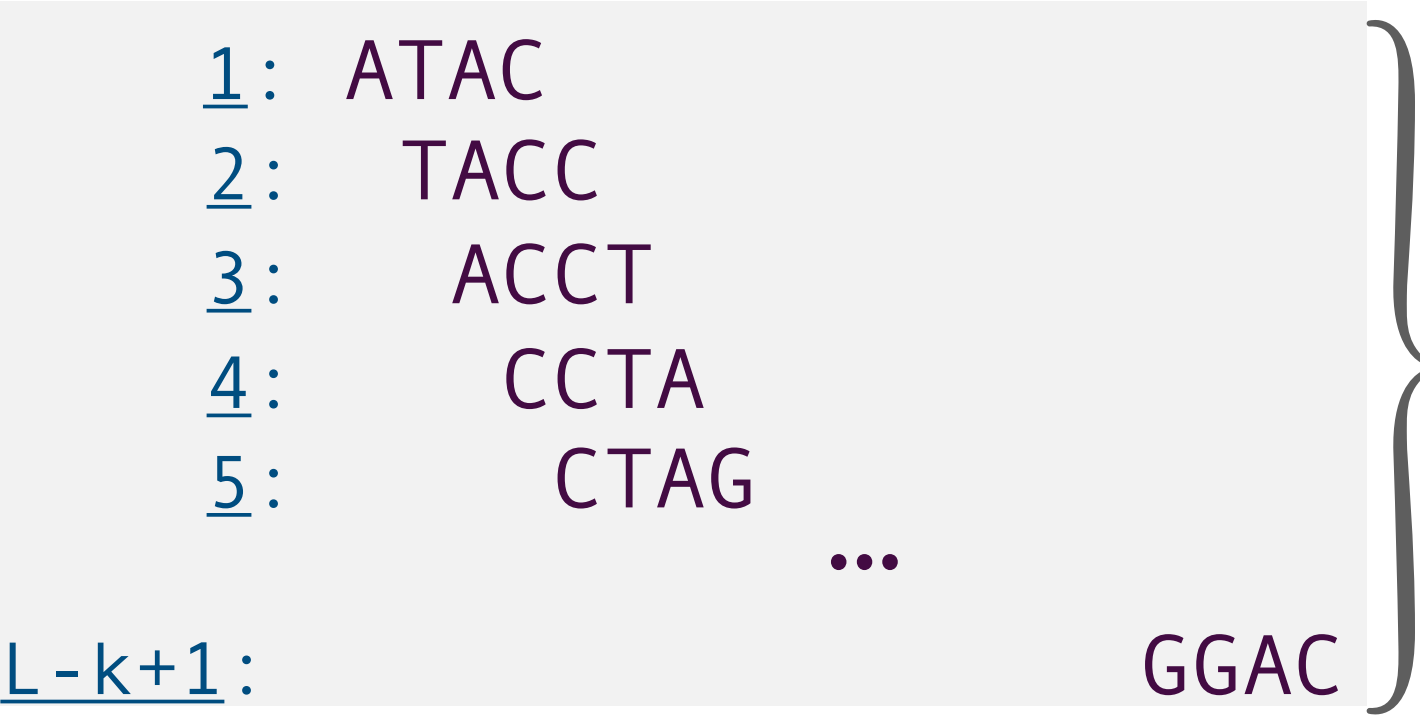
search *k*-mer matches up to a HD threshold  $\delta$



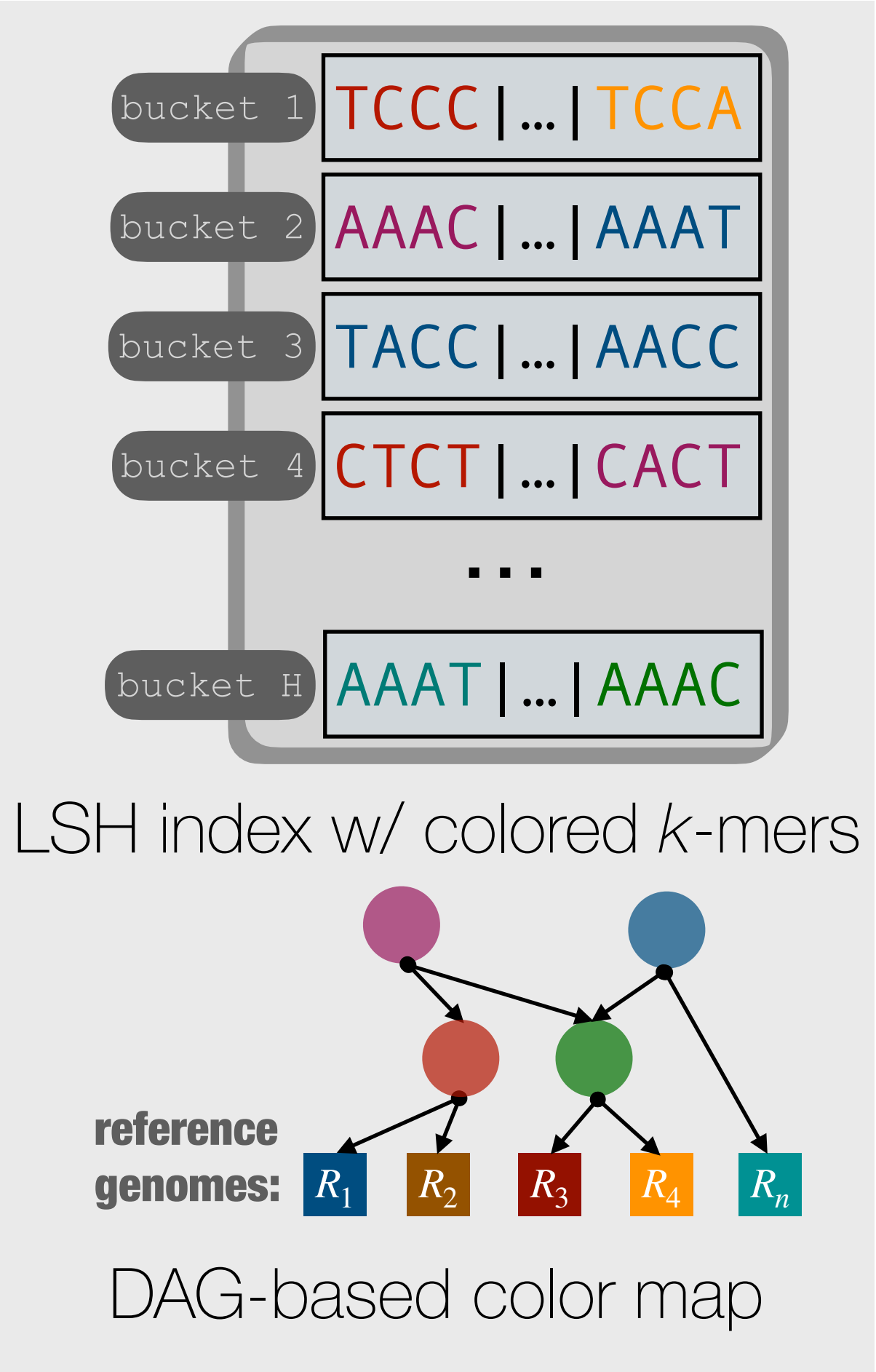
# Finding homologous k-mers of indexed references

Given a query sequence;

```
>ID_XYZ
ATACCTAGGAGTACGGGAC
```



search  $k$ -mer matches up to a HD threshold  $\delta$



keep closest HD match as homologous

sparse table

	$R_1$	$R_2$	...	$R_n$
<u>1</u>	0	1	...	-
<u>2</u>	1	4	...	4
<u>3</u>	-	-	...	-
<u>4</u>	-	2	...	-
<u>5</u>	0	0	...	3
...	...	...	...	...
<u>L-k+1</u>	3	0	...	4

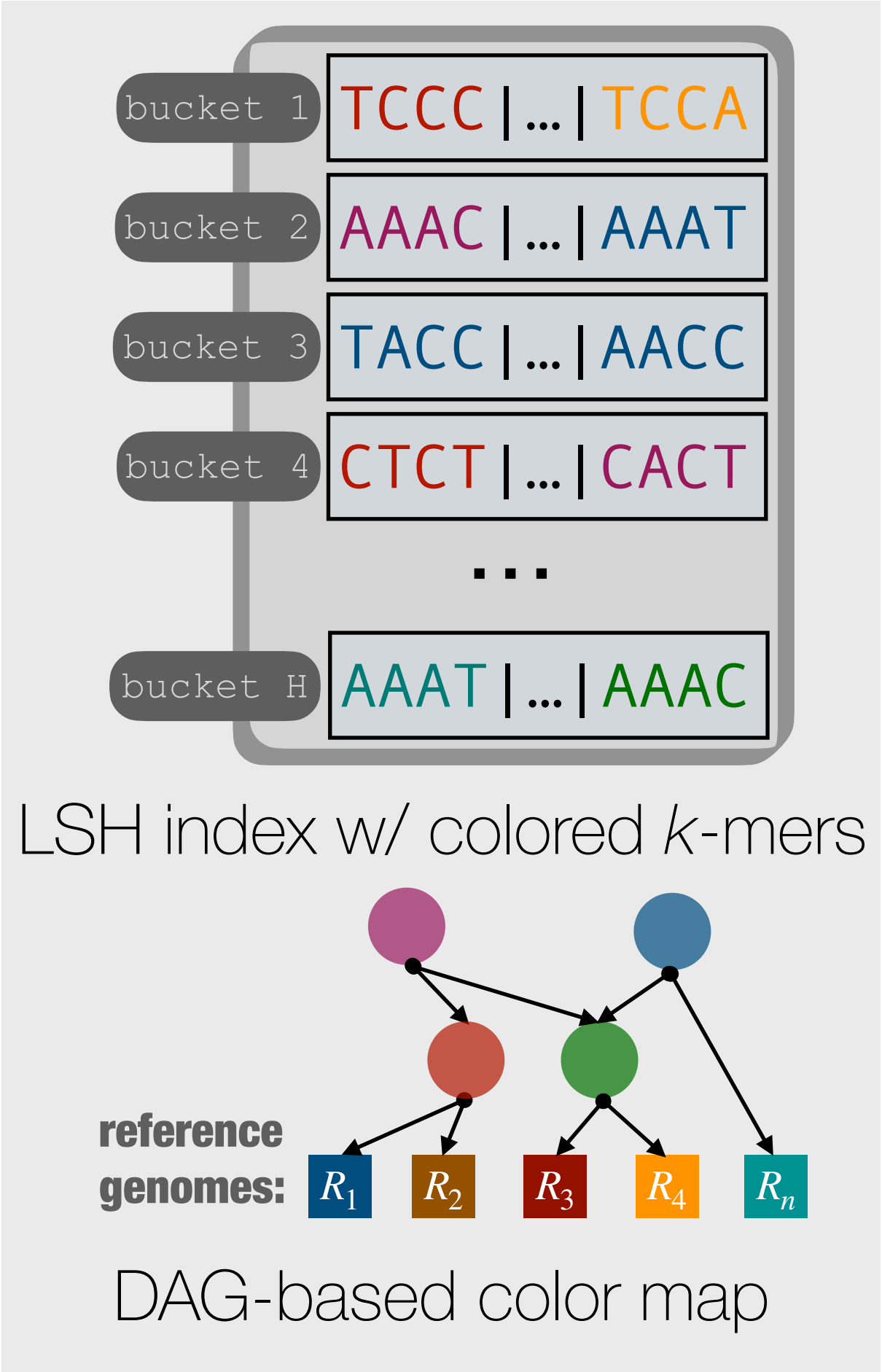
# Finding homologous k-mers of indexed references

Given a query sequence;

```
>ID_XYZ
ATACCTAGGAGTACGGGAC
```



search  $k$ -mer matches up to a HD threshold  $\delta$



keep closest HD match as homologous

sparse table

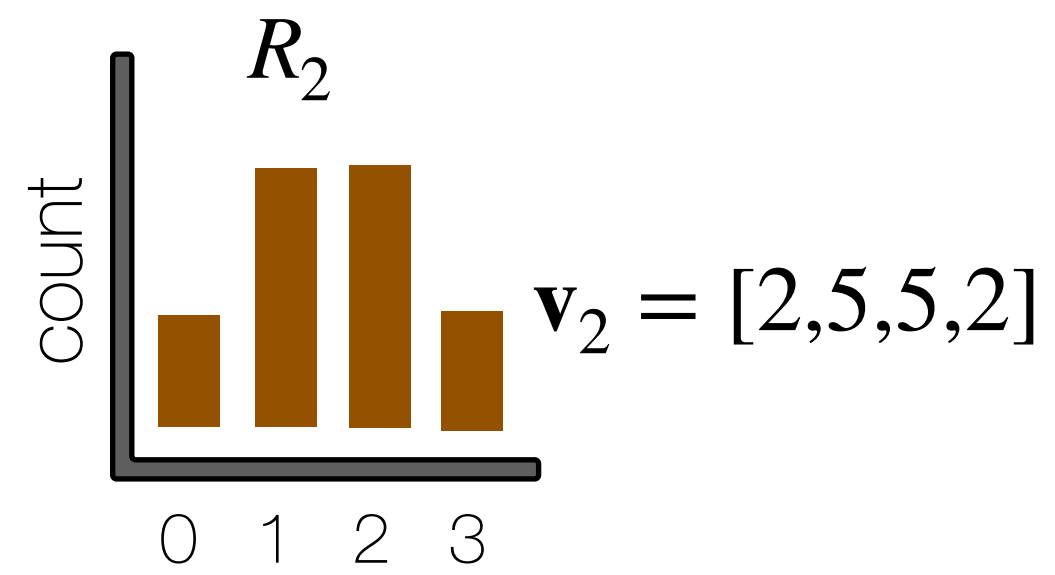
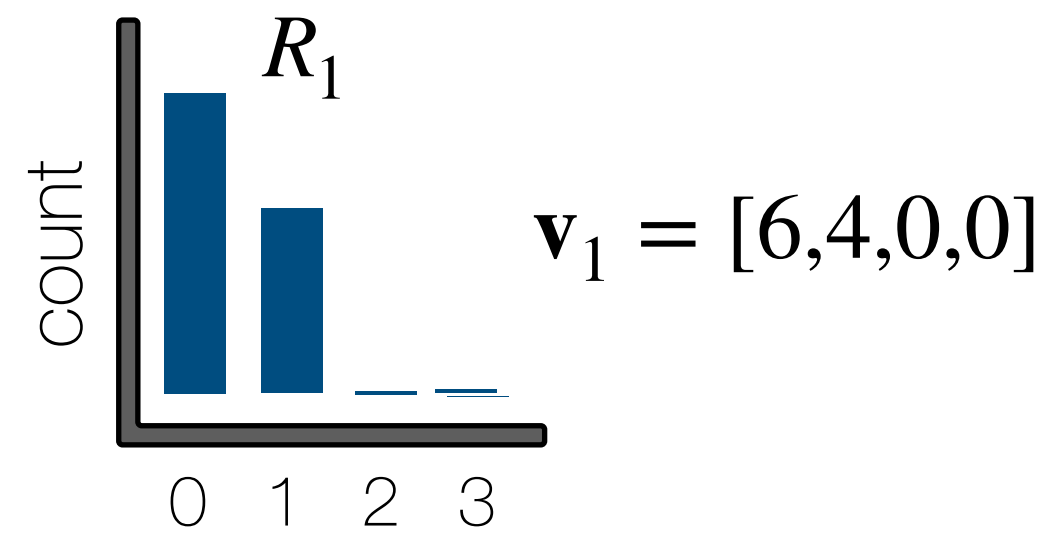
	$R_1$	$R_2$	...	$R_n$
<u>1</u>	0	1	...	-
<u>2</u>	1	4	...	4
<u>3</u>	-	-	...	-
<u>4</u>	-	2	...	-
<u>5</u>	0	0	...	3
...	...	...	...	...
<u>L-k+1</u>	3	0	...	4

summarize as a histogram for each reference

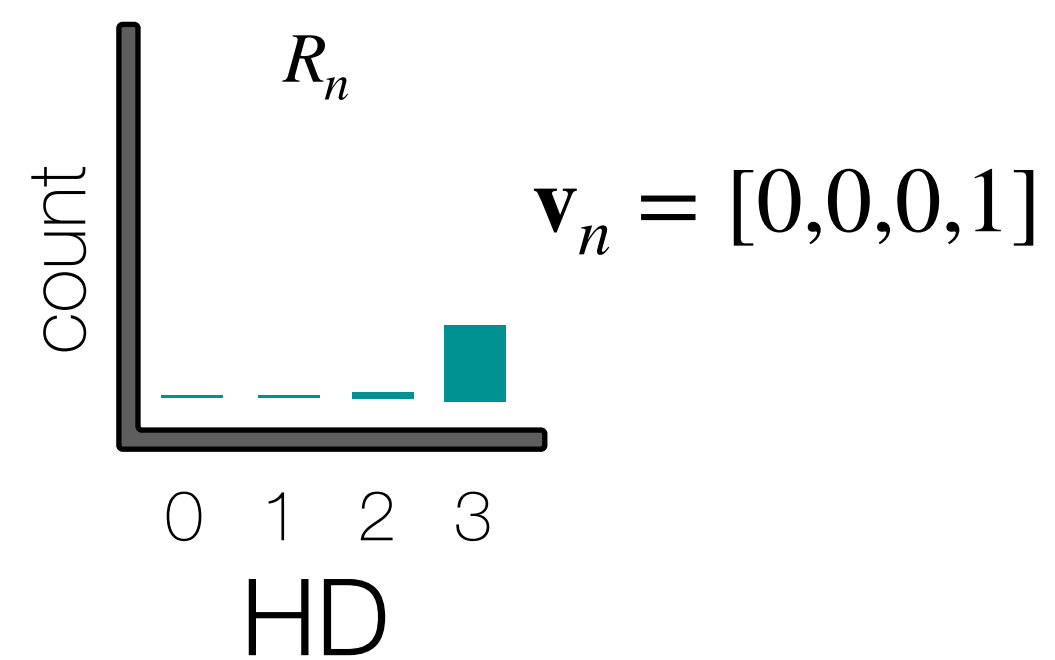
# **Likelihood of k-mer matches & Hamming distances**

# Likelihood of k-mer matches & Hamming distances

## Hamming distance histograms

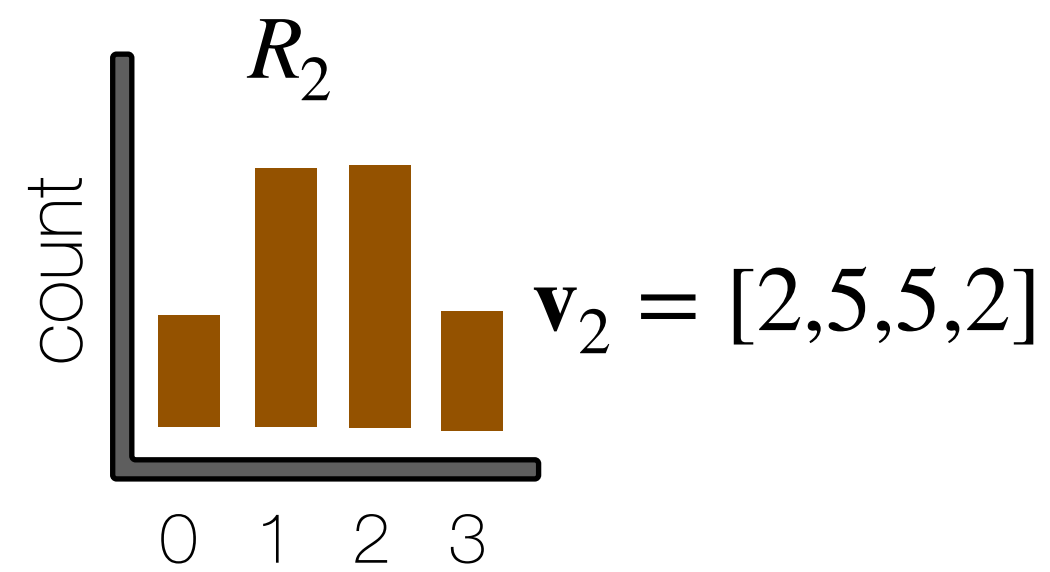
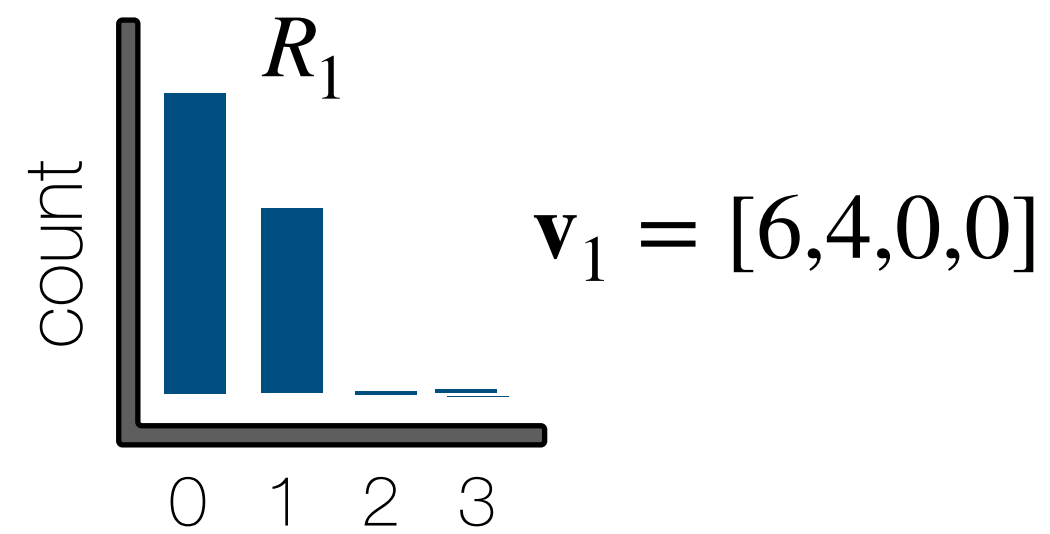


...

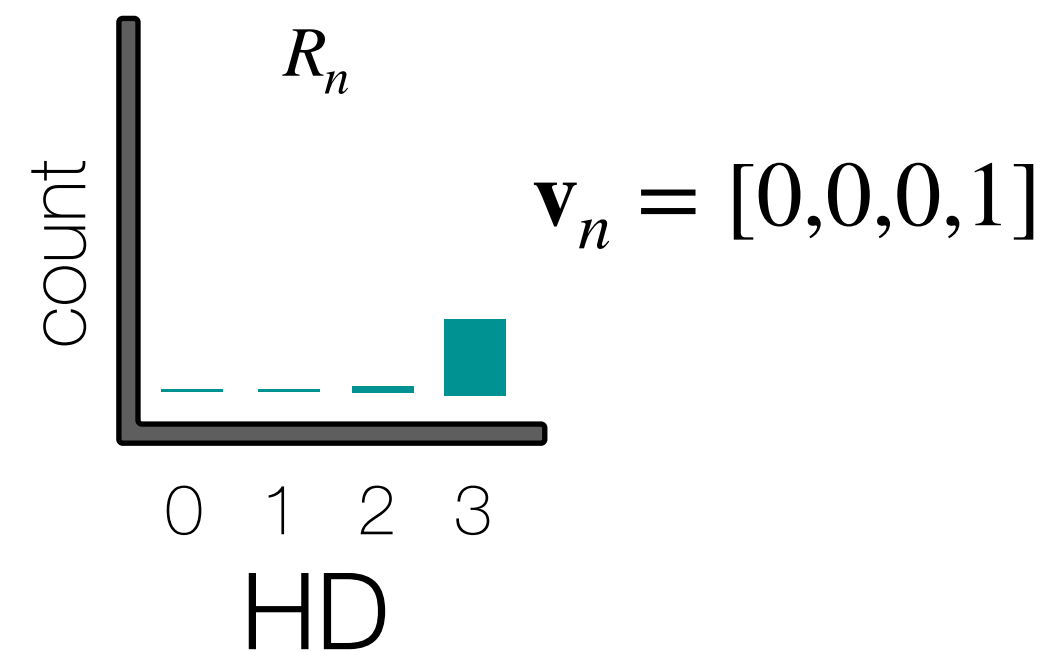


# Likelihood of k-mer matches & Hamming distances

## Hamming distance histograms



...



**Goal:** compute the probability of  $q$  having distance  $D$  to  $R_i$

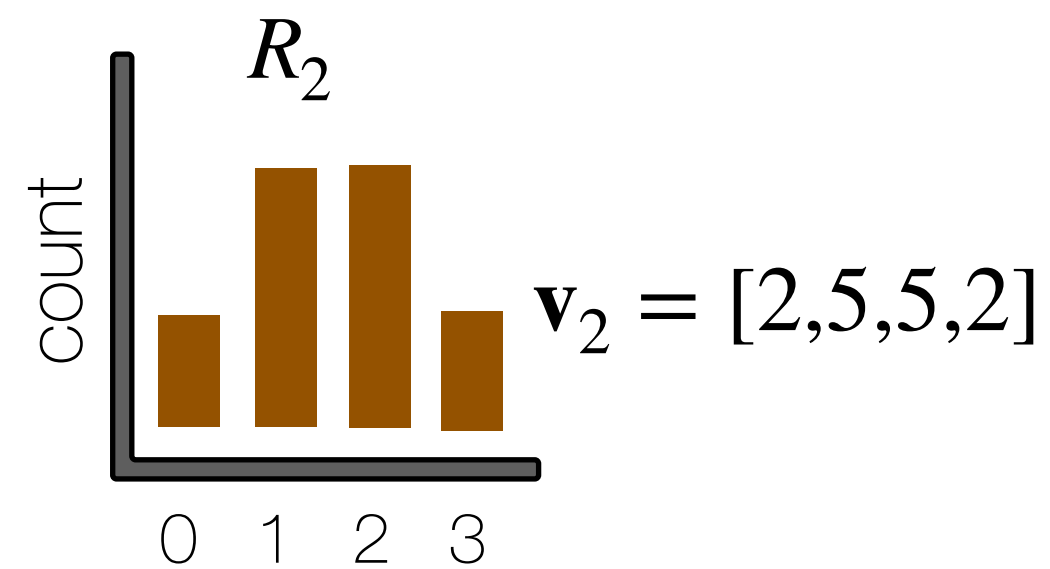
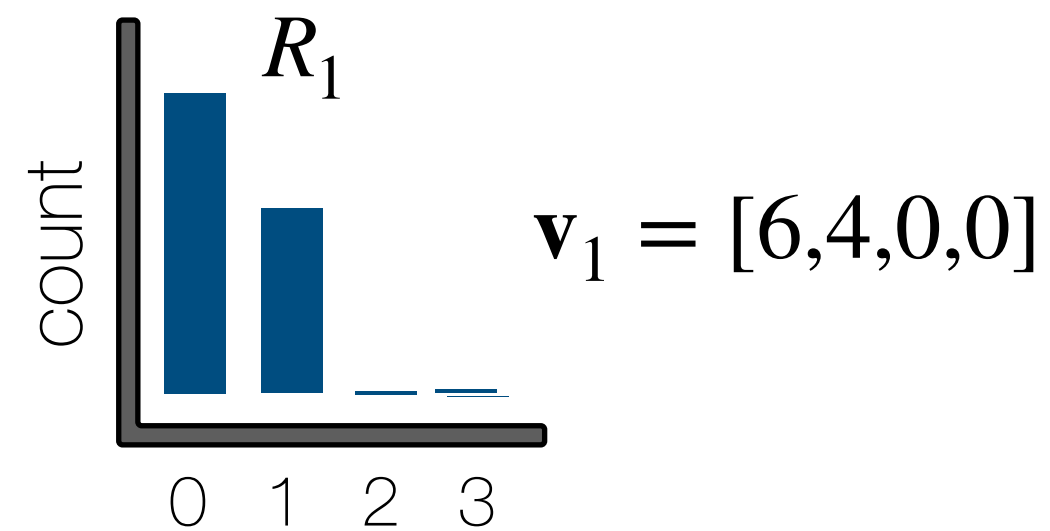
- $\mathbf{v}_i$ : match count for each HD up to  $\delta$
- $u_i$ : number of mismatches  $(L - k + 1) - \sum v_i$

Likelihood of distance  $D$  to reference  $R_i$

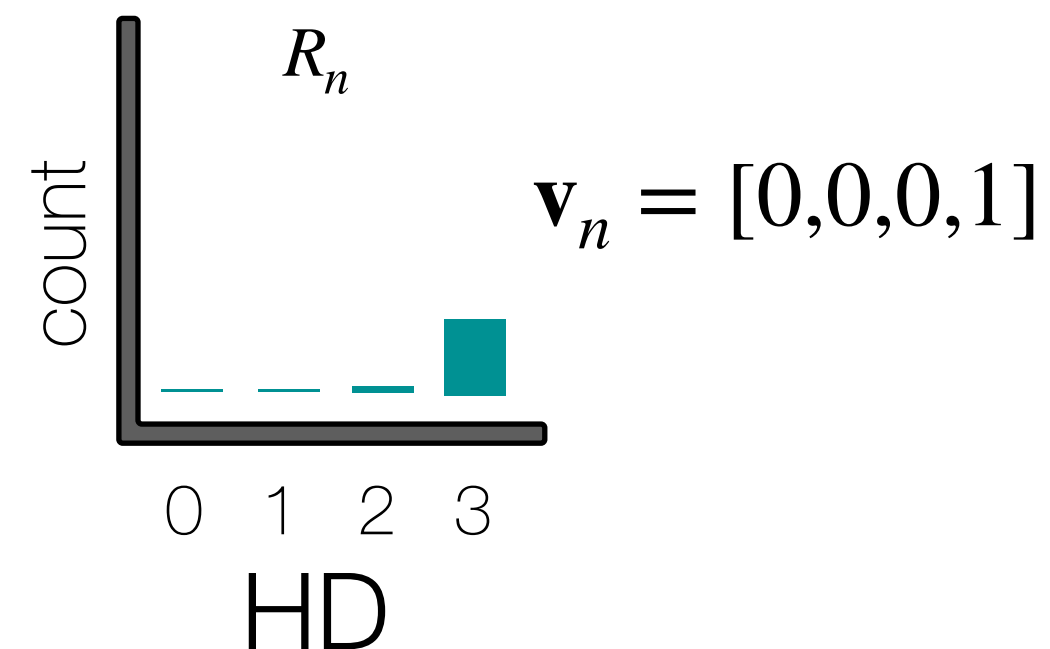
$$\mathcal{L}_i(D; k, h, \delta, u_i, \mathbf{v}_i) =$$

# Likelihood of $k$ -mer matches & Hamming distances

## Hamming distance histograms



...



**Goal:** compute the probability of  $q$  having distance  $D$  to  $R_i$

- $\mathbf{v}_i$ : match count for each HD up to  $\delta$
- $u_i$ : number of mismatches  $(L - k + 1) - \sum v_i$

## The Statistics of $k$ -mers from a Sequence Undergoing a Simple Mutation Process Without Spurious Matches

Authors: [Antonio Blanca](#), [Robert S. Harris](#), [David Koslicki](#), and [Paul Medvedev](#) | [AUTHORS INFO & AFFILIATIONS](#)

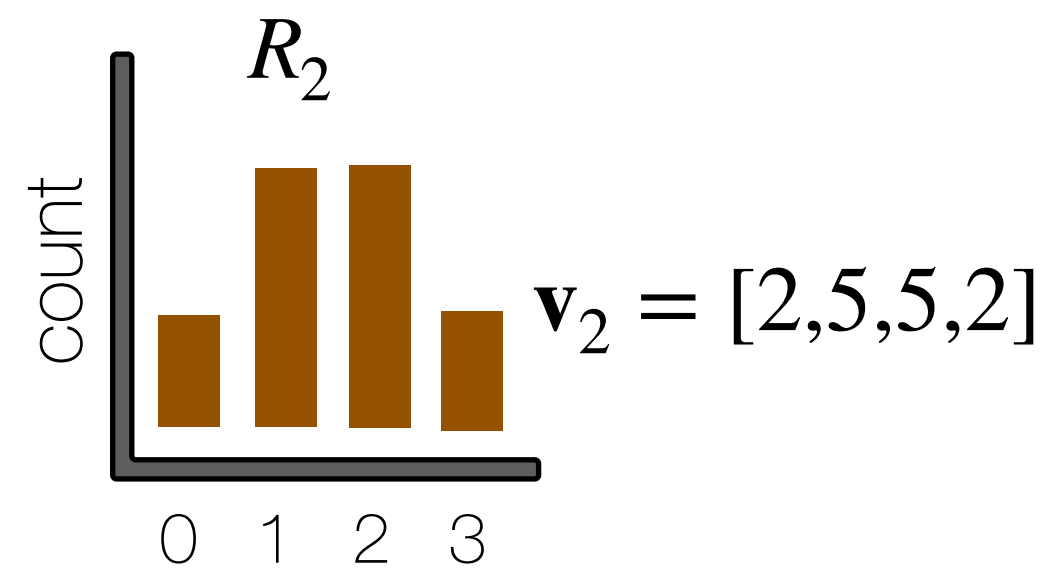
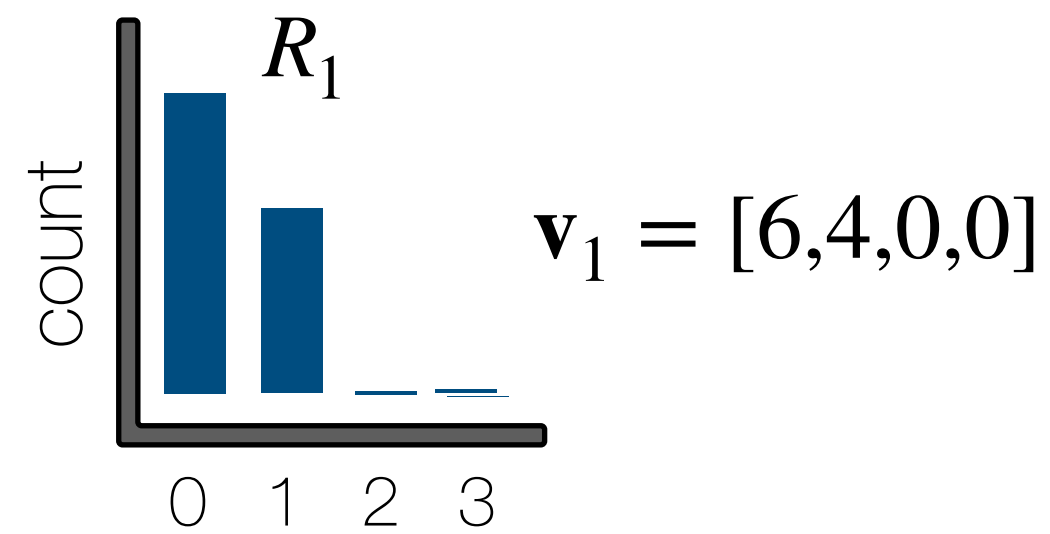
Publication: Journal of Computational Biology • <https://doi.org/10.1089/cmb.2021.0431>

Likelihood of distance  
 $D$  to reference  $R_i$

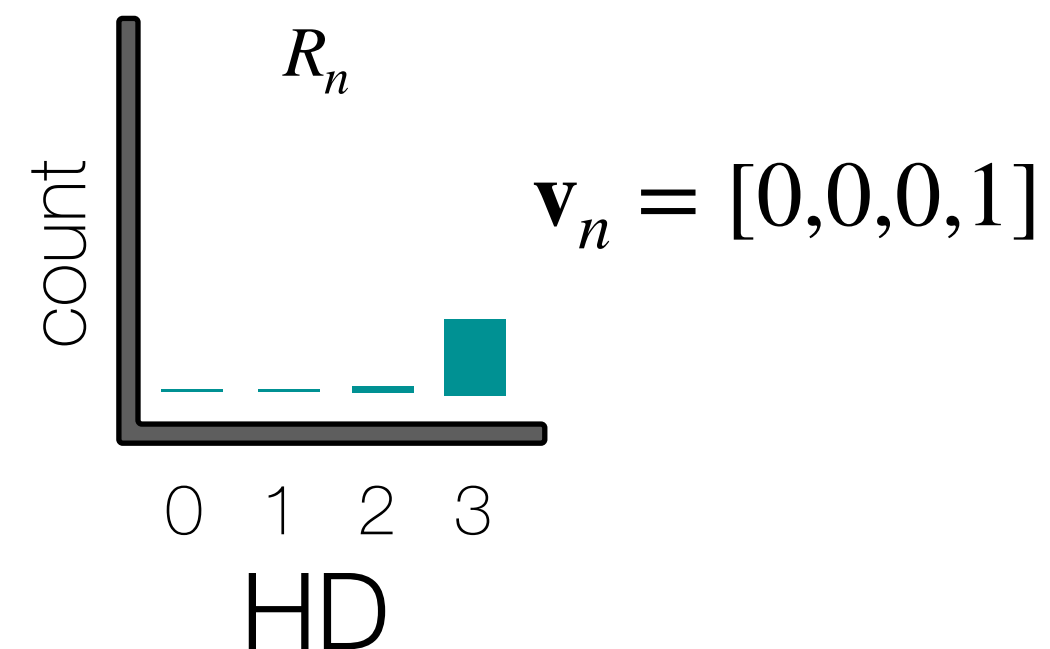
$$\mathcal{L}_i(D; k, h, \delta, u_i, \mathbf{v}_i) =$$

# Likelihood of k-mer matches & Hamming distances

## Hamming distance histograms



...



**Goal:** compute the probability of  $q$  having distance  $D$  to  $R_i$

- $\mathbf{v}_i$ : match count for each HD up to  $\delta$
- $u_i$ : number of mismatches  $(L - k + 1) - \sum v_i$

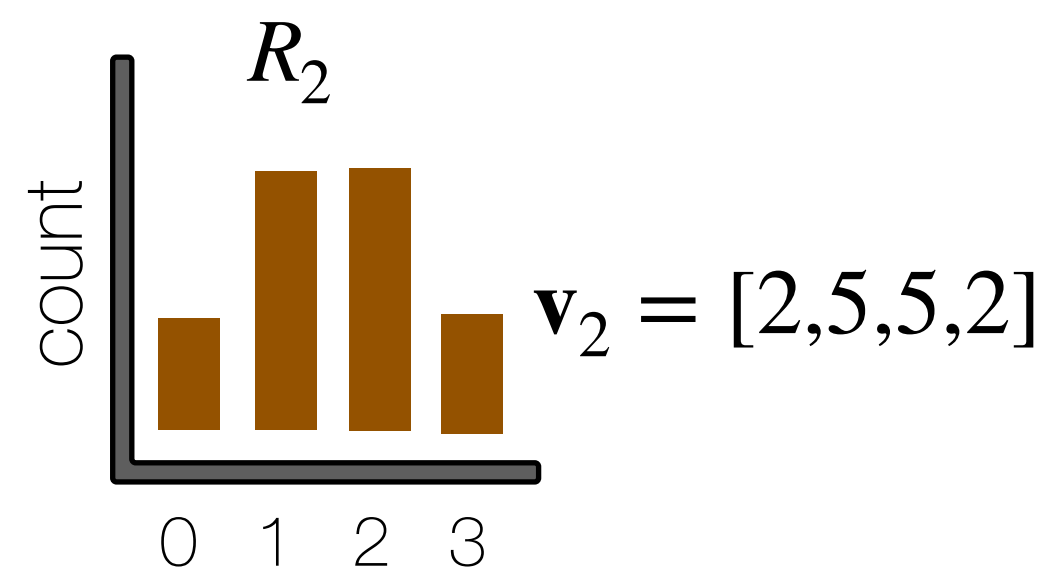
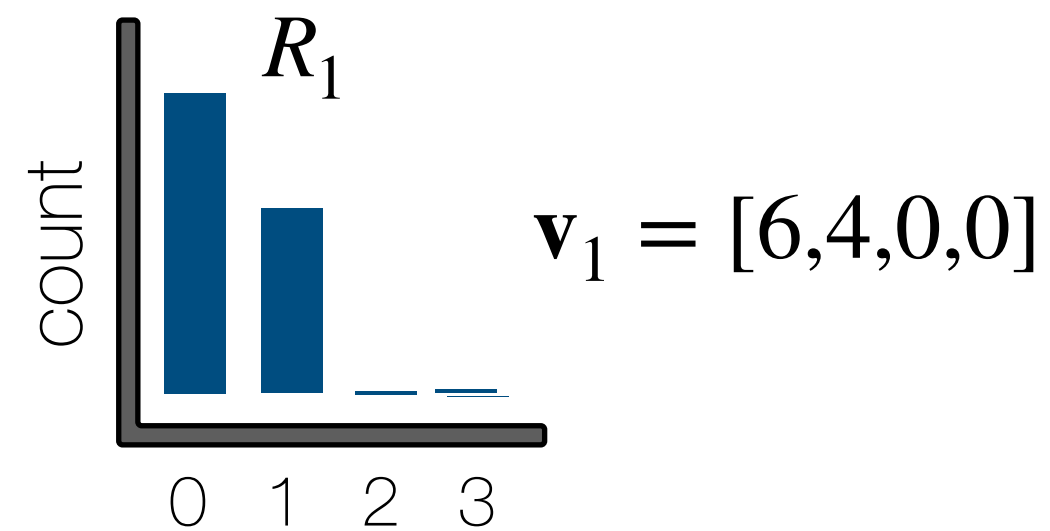
**Independence assumption:** treat  $q$  as a bag of  $L - k + 1$  k-mers

Likelihood of distance  
 $D$  to reference  $R_i$

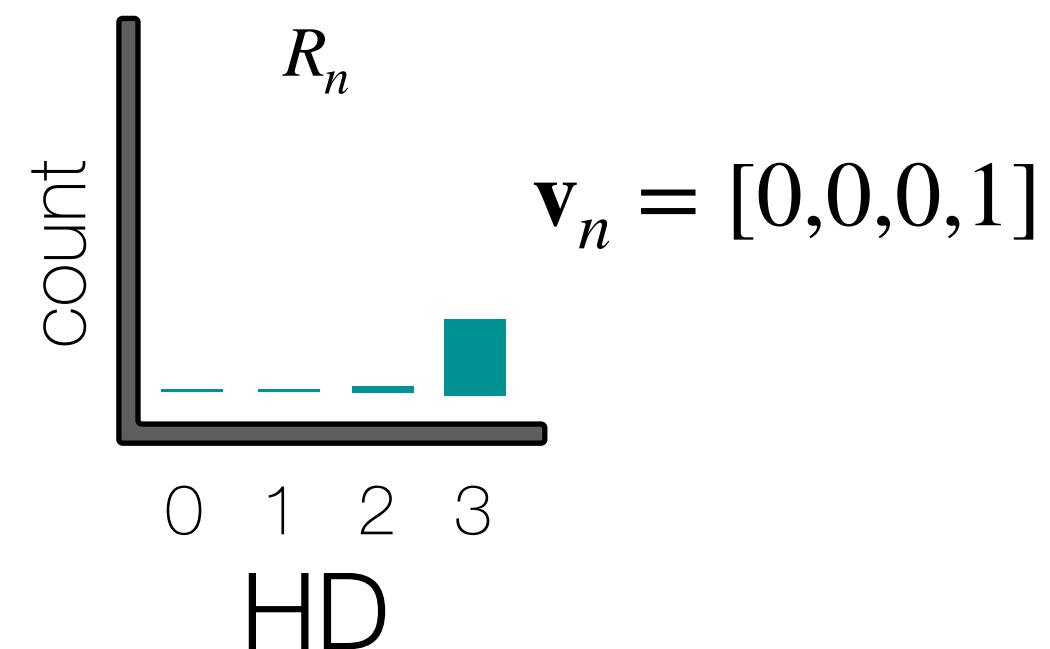
$$\mathcal{L}_i(D; k, h, \delta, u_i, \mathbf{v}_i) =$$

# Likelihood of k-mer matches & Hamming distances

## Hamming distance histograms



...



**Goal:** compute the probability of  $q$  having distance  $D$  to  $R_i$

- $\mathbf{v}_i$ : match count for each HD up to  $\delta$
- $u_i$ : number of mismatches  $(L - k + 1) - \sum v_i$

**Independence assumption:** treat  $q$  as a bag of  $L - k + 1$  k-mers

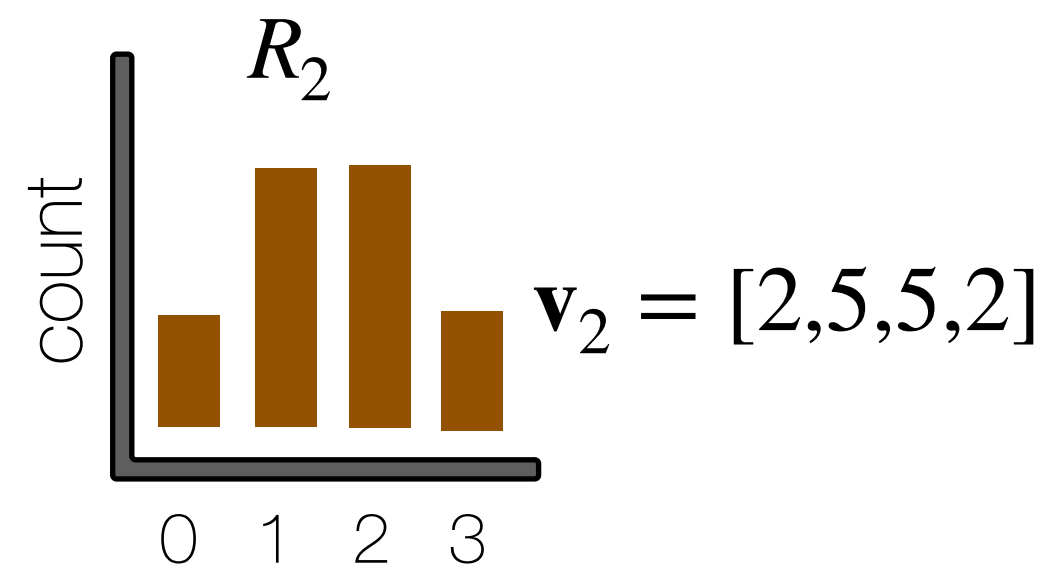
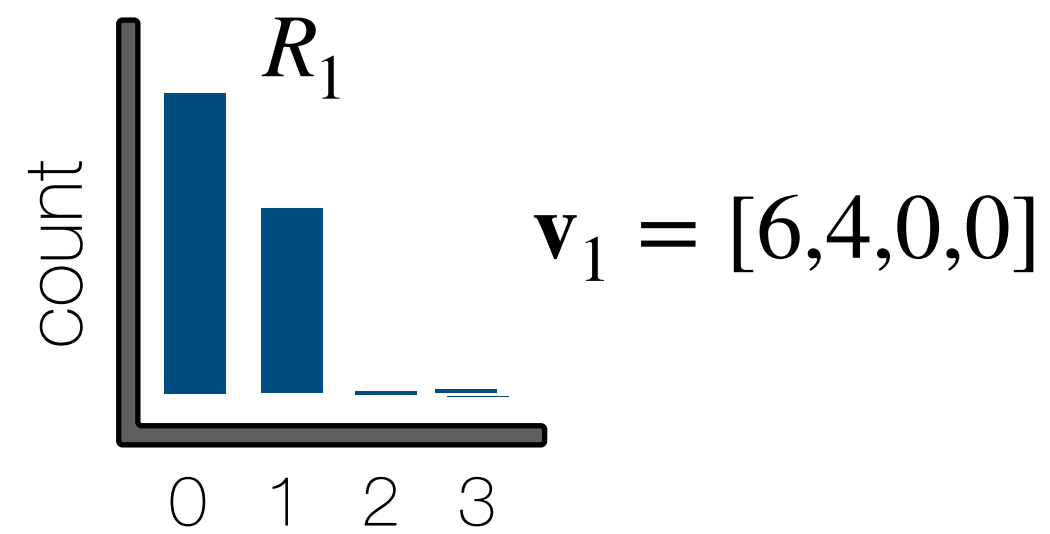
Likelihood of distance  
 $D$  to reference  $R_i$

$$\mathcal{L}_i(D; k, h, \delta, u_i, \mathbf{v}_i) = P_{\text{miss}}(D; k, h, \delta)^{u_i}$$

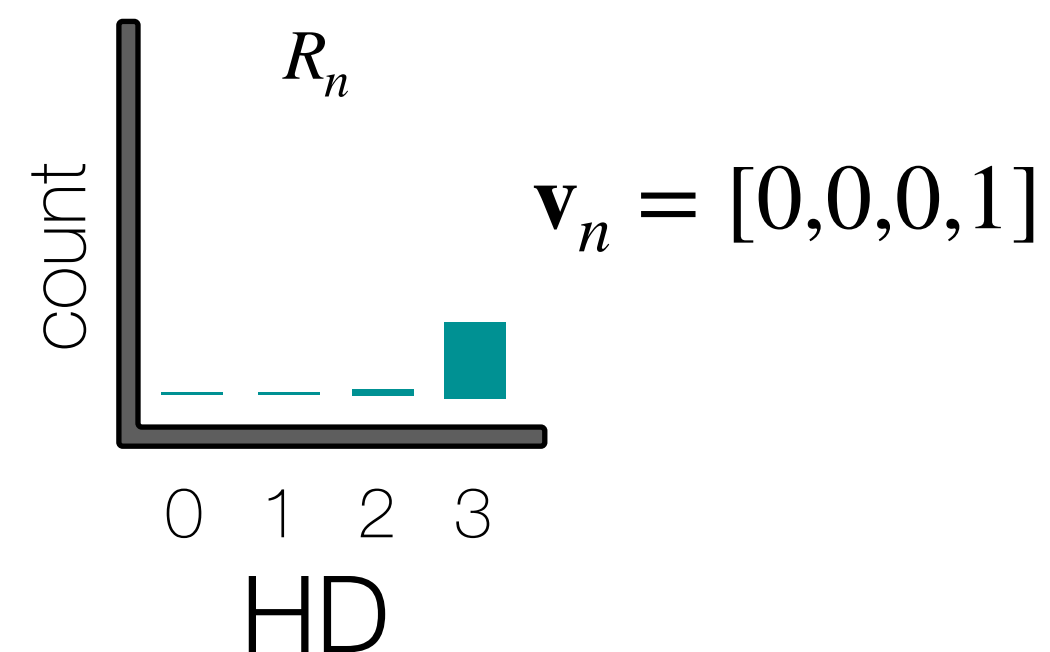
Probability of having  $u_i$   
mismatches in total

# Likelihood of k-mer matches & Hamming distances

## Hamming distance histograms



...



**Goal:** compute the probability of  $q$  having distance  $D$  to  $R_i$

- $\mathbf{v}_i$ : match count for each HD up to  $\delta$
- $u_i$ : number of mismatches  $(L - k + 1) - \sum v_i$

**Independence assumption:** treat  $q$  as a bag of  $L - k + 1$  k-mers

Likelihood of distance  
 $D$  to reference  $R_i$

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

Probability of having  $u_i$   
mismatches in total

Probability of having  $v_{i,d}$   
matches w/ HD =  $d$

# **Observing homologous k-mers with varying HDs**

# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) =$$

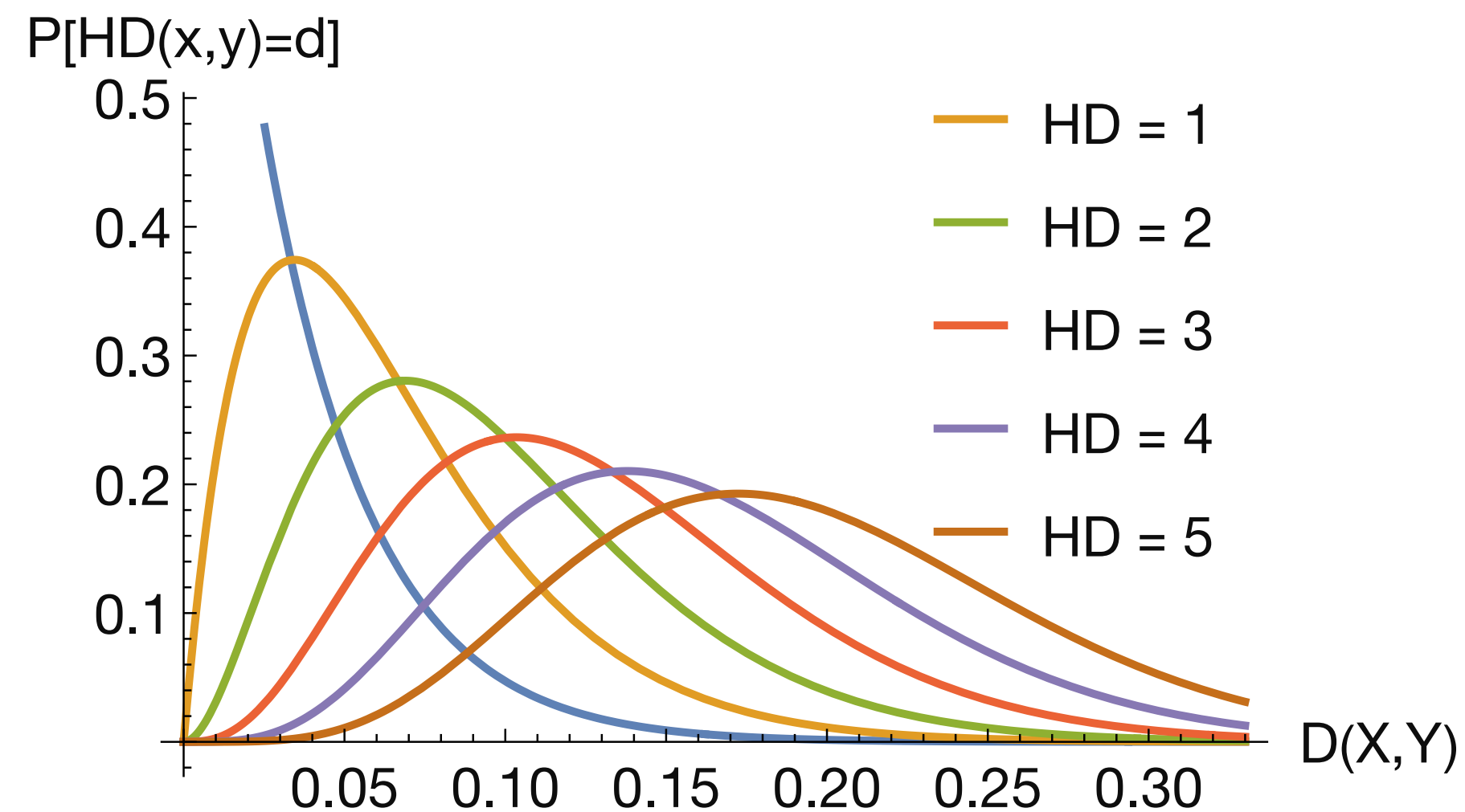
# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) = P_{mutate}(D; d, k)$$

# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) = P_{mutate}(D; d, k)$$

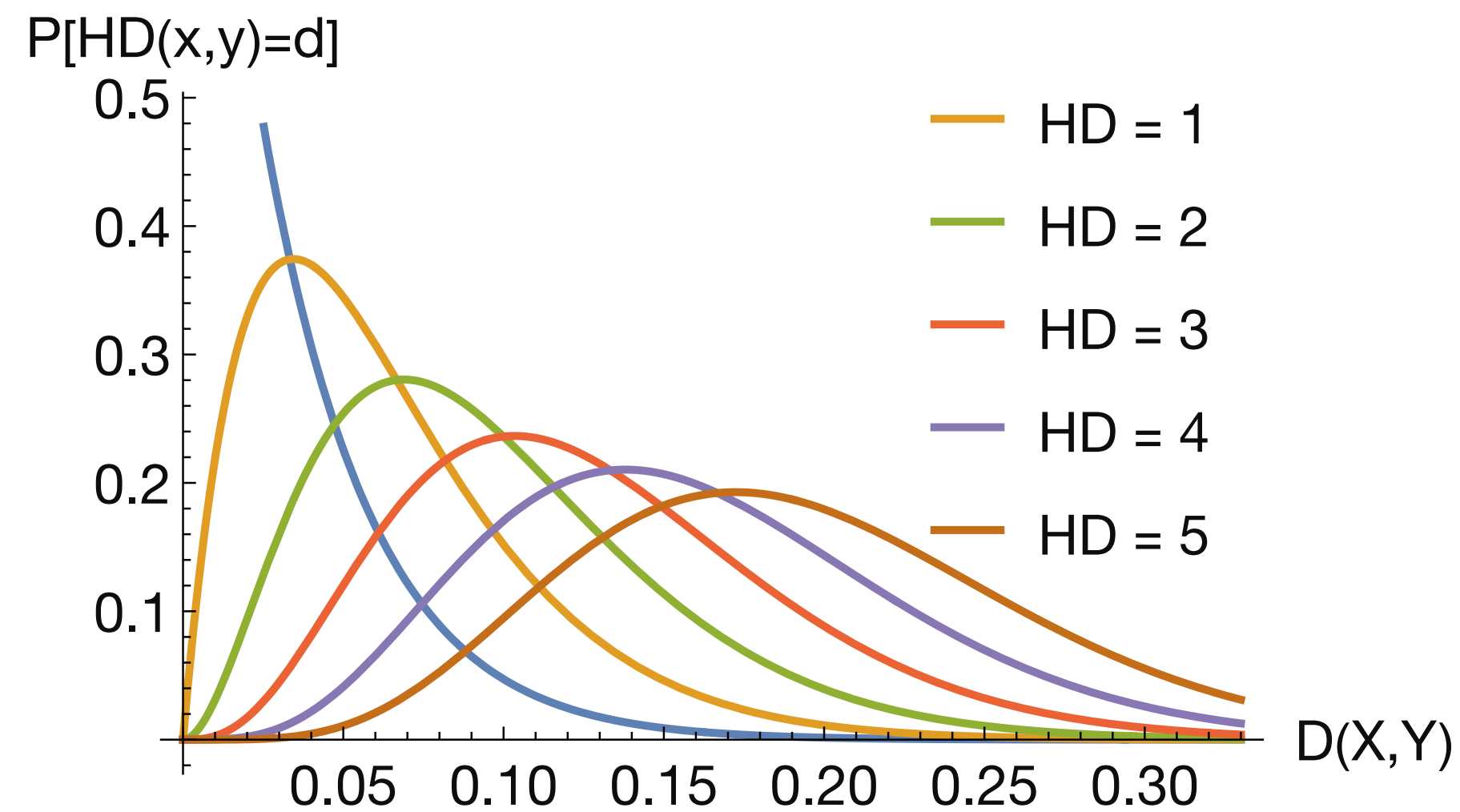
$$D^d(1 - D)^{(k-d)} \binom{k}{d}$$



# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) = P_{mutate}(D; d, k) P_{collide}(d, k, h)$$

$$D^d(1 - D)^{(k-d)} \binom{k}{d}$$

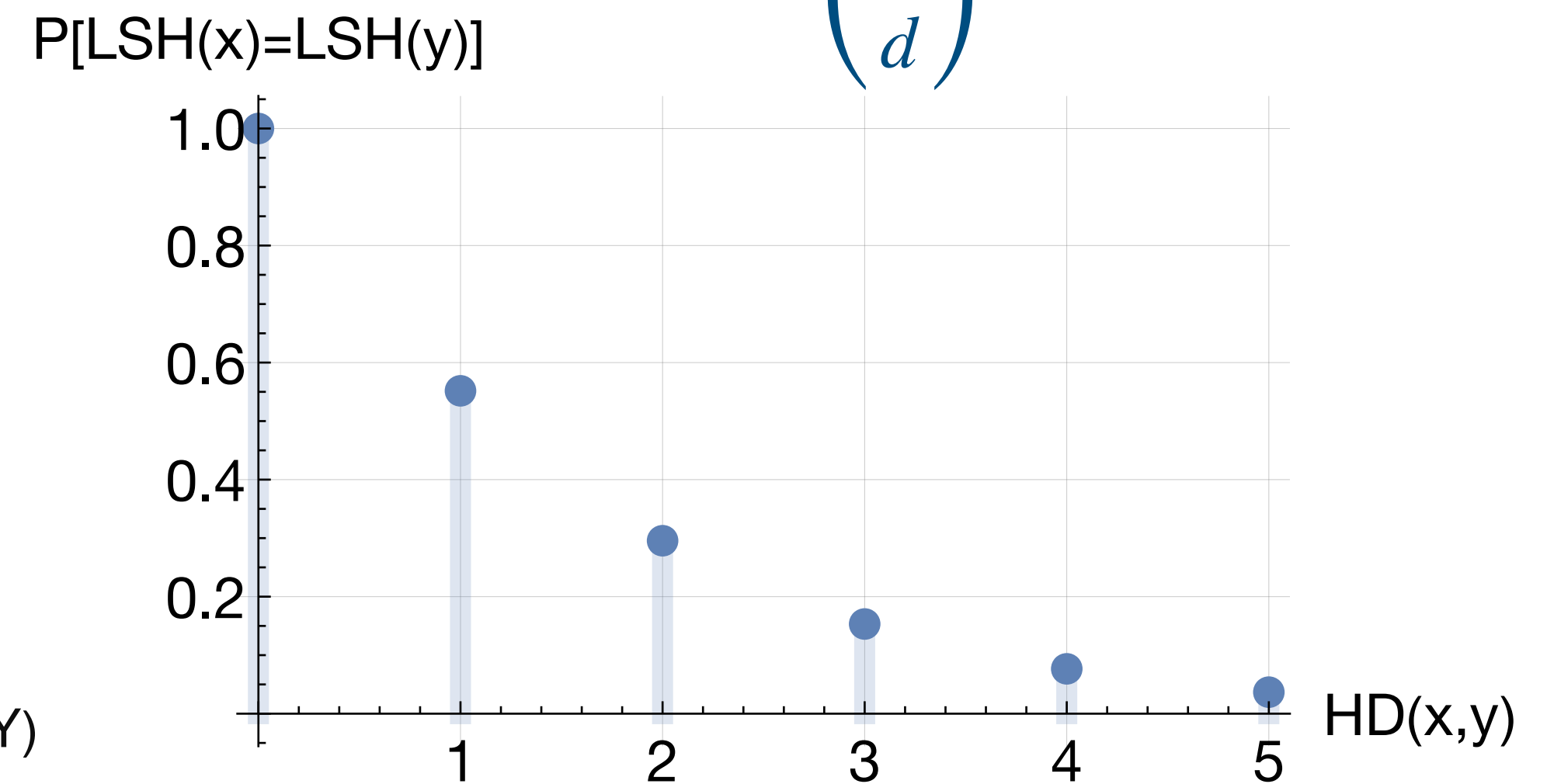
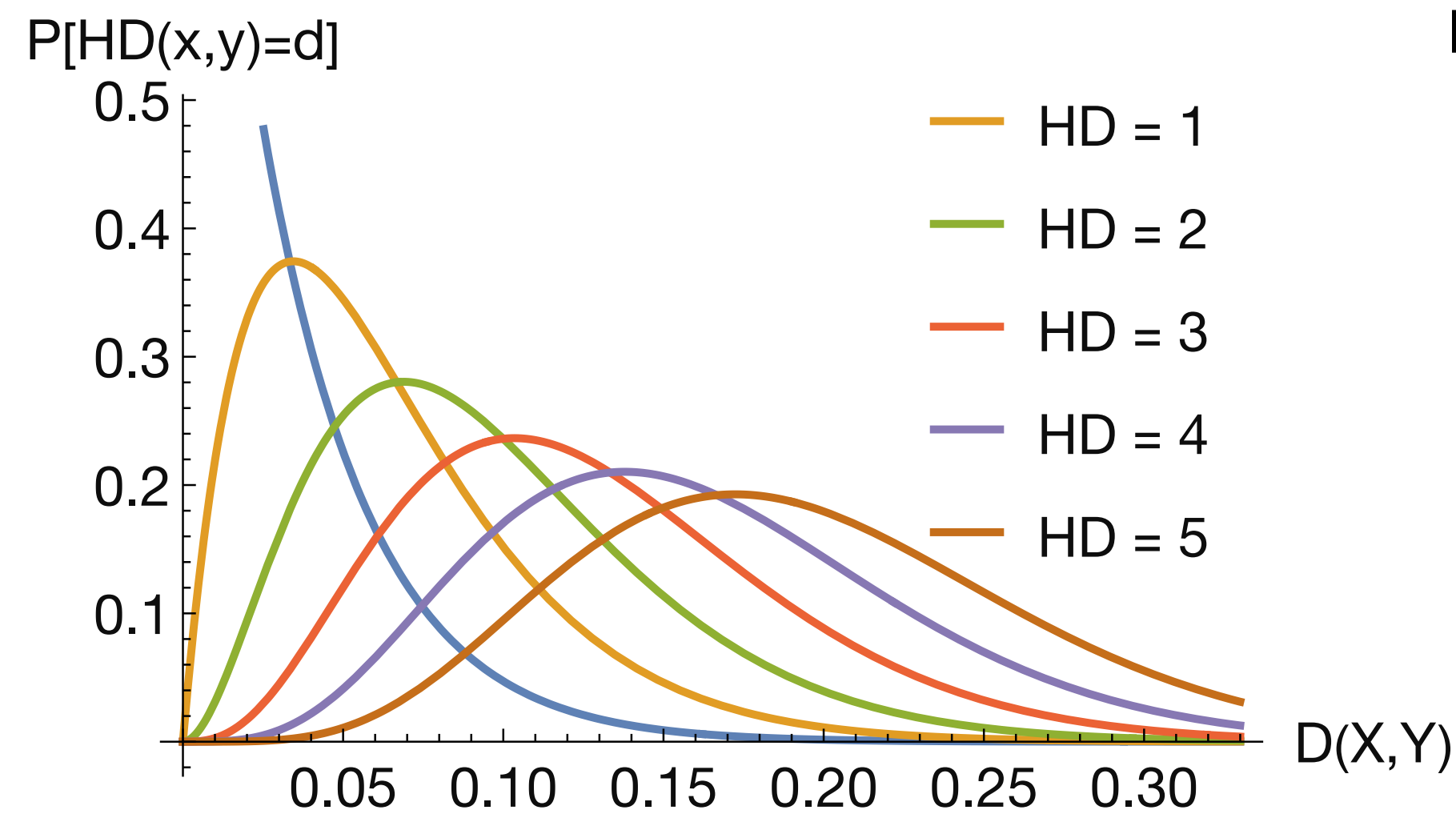


# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) = P_{mutate}(D; d, k) P_{collide}(d, k, h)$$

$$D^d (1 - D)^{(k-d)} \binom{k}{d}$$

$$\frac{\binom{k-h}{d}}{\binom{k}{d}}$$

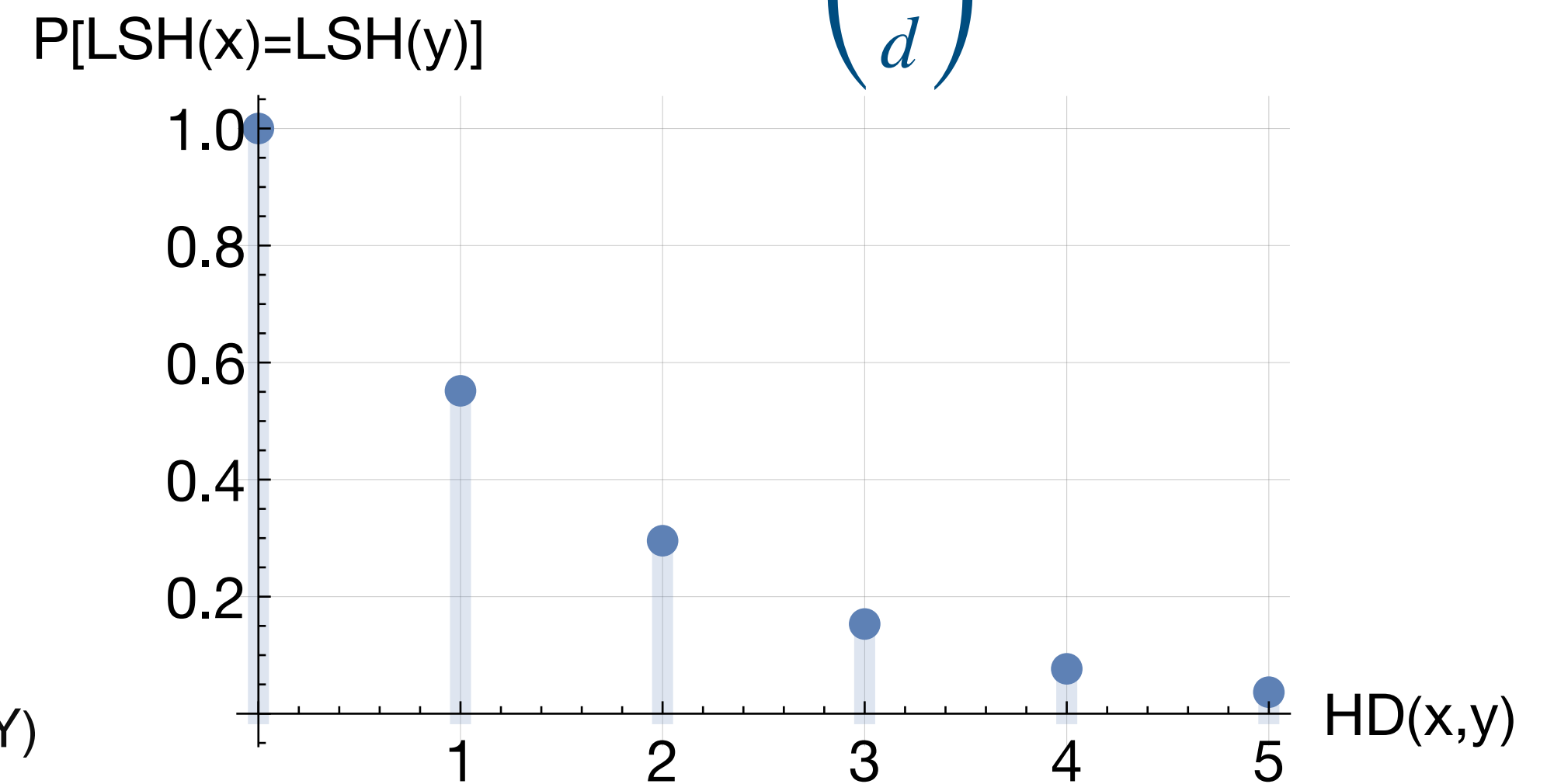
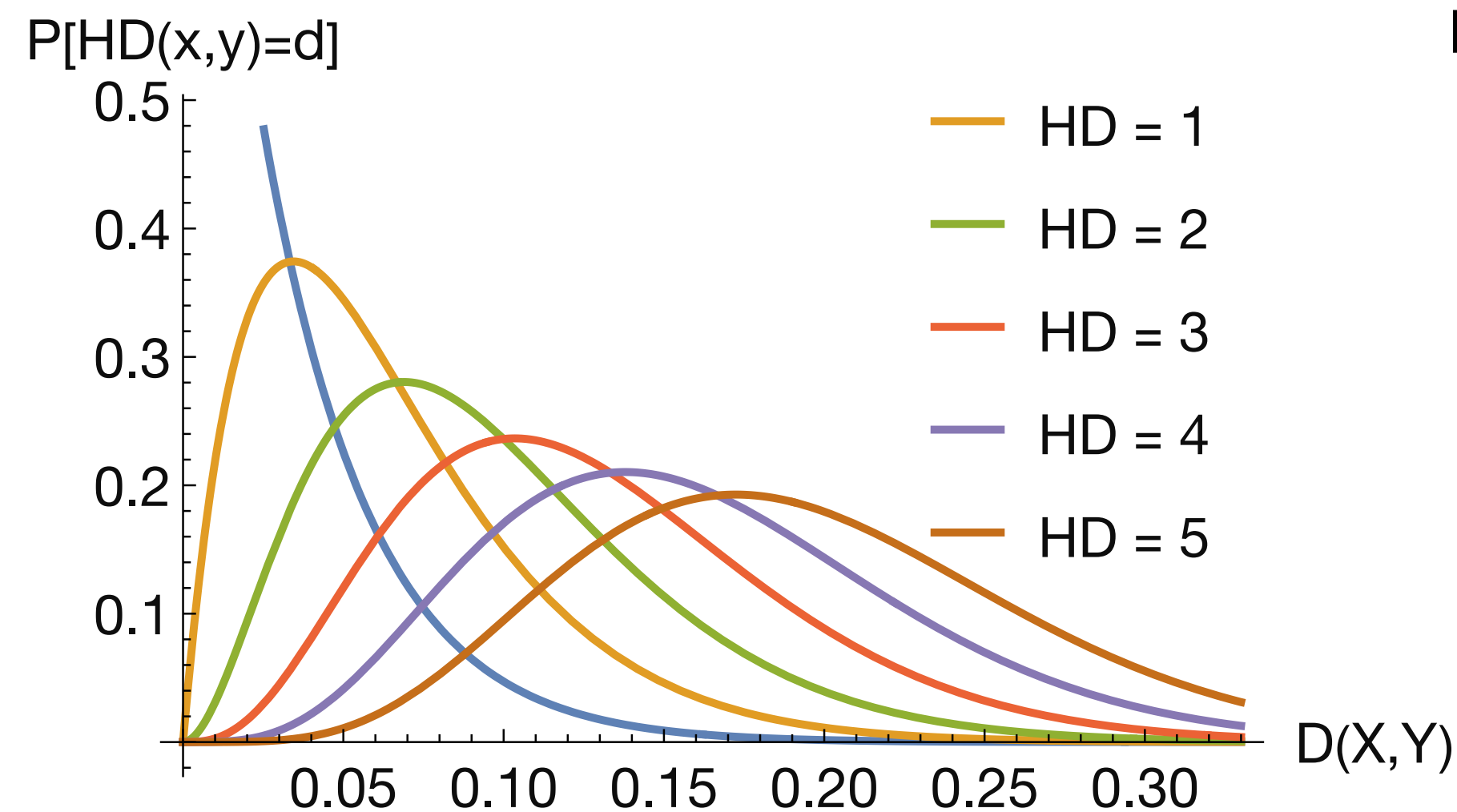


# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) = \rho_i P_{mutate}(D; d, k) P_{collide}(d, k, h)$$

$$D^d (1 - D)^{(k-d)} \binom{k}{d}$$

$$\frac{\binom{k-h}{d}}{\binom{k}{d}}$$



# Observing homologous k-mers with varying HDs

$$P_{match}(D; d, k, h) = \rho_i P_{mutate}(D; d, k) P_{collide}(d, k, h)$$

$$\rho_i = \frac{\text{\# of subsampled}}{\text{\# of distinct}}$$

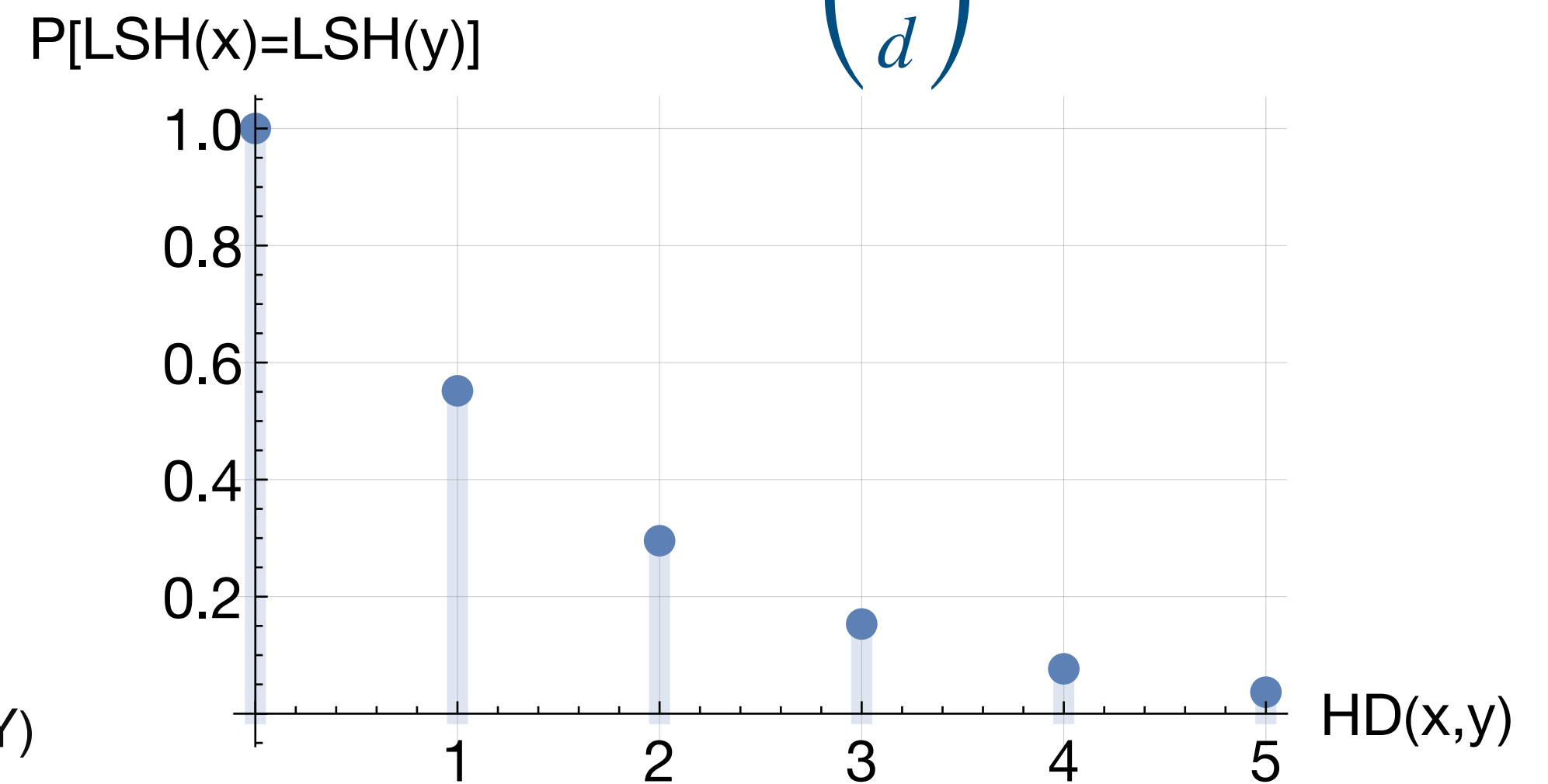
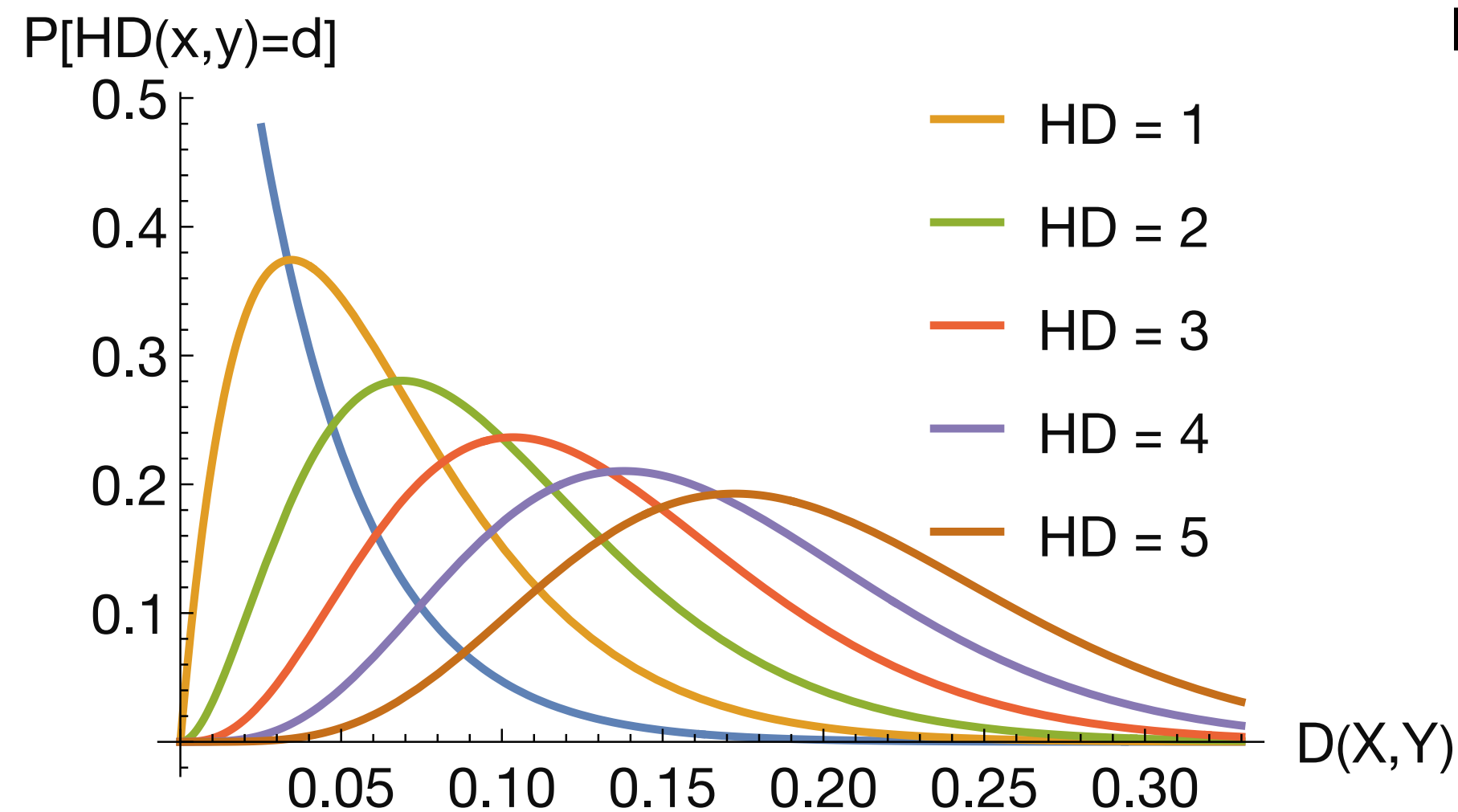
$$D^d (1 - D)^{(k-d)} \binom{k}{d}$$

$$\frac{\binom{k-h}{d}}{\binom{k}{d}}$$

precomputed for  $R_i$

→ not all  $k$ -mers have to be indexed:

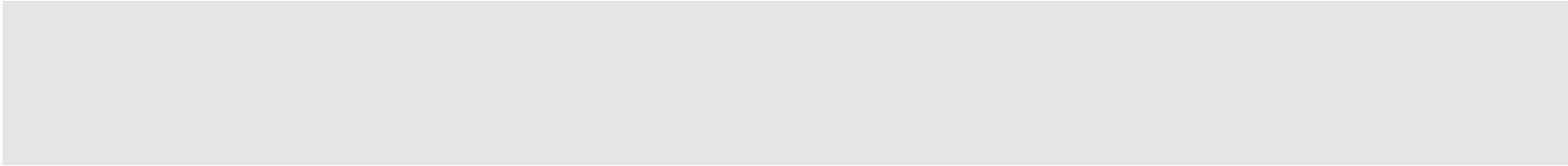
- minimizers
- FracMinHash
- ...



**Multiple events could lead to a mismatch might**

# Multiple events could lead to a mismatch might

A mismatch occurs for two homologous  $k$ -mers (query  $x$  and reference  $y$ )

$$P_{miss}(D; d, k, h, \delta) =$$


# Multiple events could lead to a mismatch might

A mismatch occurs for two homologous  $k$ -mers (query  $x$  and reference  $y$ )

- if  $y$  is not indexed:  $1 - \rho$

$$P_{miss}(D; d, k, h, \delta) =$$

$$(1 - \rho)$$

# Multiple events could lead to a mismatch might

A mismatch occurs for two homologous  $k$ -mers (query  $x$  and reference  $y$ )

- if  $y$  is not indexed:  $1 - \rho$
- if  $y$  is indexed (with probability  $\rho$ ), but either:

$$P_{miss}(D; d, k, h, \delta) =$$

$$(1 - \rho) + \rho \left( \right)$$

# Multiple events could lead to a mismatch might

A mismatch occurs for two homologous  $k$ -mers (query  $x$  and reference  $y$ )

- if  $y$  is not indexed:  $1 - \rho$
- if  $y$  is indexed (with probability  $\rho$ ), but either:
  - $\text{HD}(x, y) > \delta$  or  $\text{LSH}(x) \neq \text{LSH}(y)$ !

$$P_{\text{miss}}(D; d, k, h, \delta) =$$

$$(1 - \rho) + \rho \left( \sum_{d=\delta+1}^k P_{\text{mutate}}(D; d, k) \right)$$

# Multiple events could lead to a mismatch might

A mismatch occurs for two homologous  $k$ -mers (query  $x$  and reference  $y$ )

- if  $y$  is not indexed:  $1 - \rho$
- if  $y$  is indexed (with probability  $\rho$ ), but either:
  - $\text{HD}(x, y) > \delta$  or  $\text{LSH}(x) \neq \text{LSH}(y)$ !

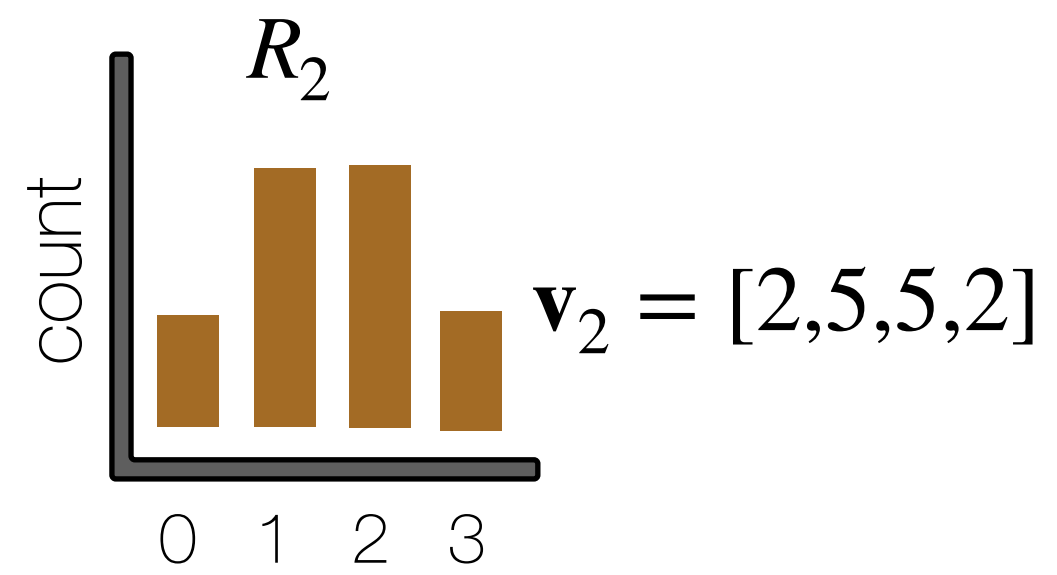
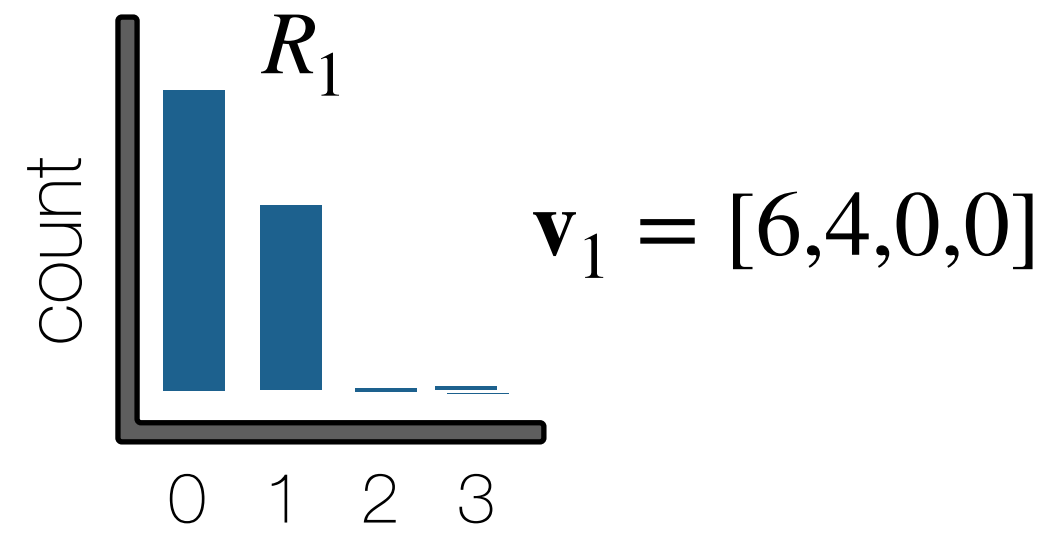
$$P_{\text{miss}}(D; d, k, h, \delta) =$$

$$(1 - \rho) + \rho \left( \sum_{d=\delta+1}^k P_{\text{mutate}}(D; d, k) + \sum_{d=0}^{\delta} P_{\text{mutate}}(D; d, k) (1 - P_{\text{collide}}(d, k, h)) \right)$$

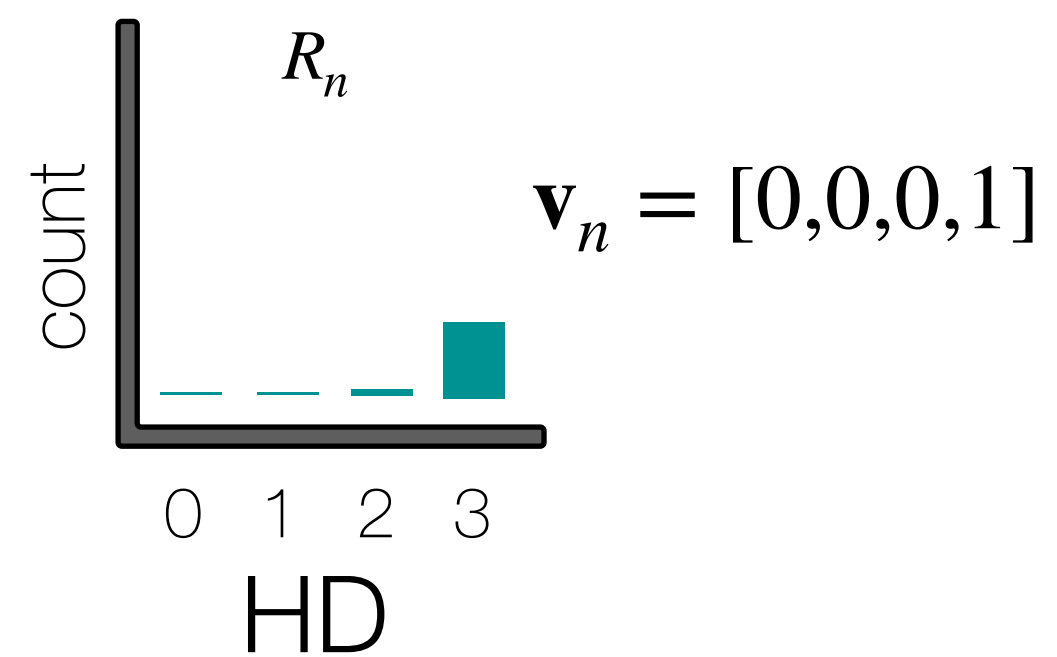
# **Maximum likelihood estimation of distances**

# Maximum likelihood estimation of distances

## Hamming distance histograms

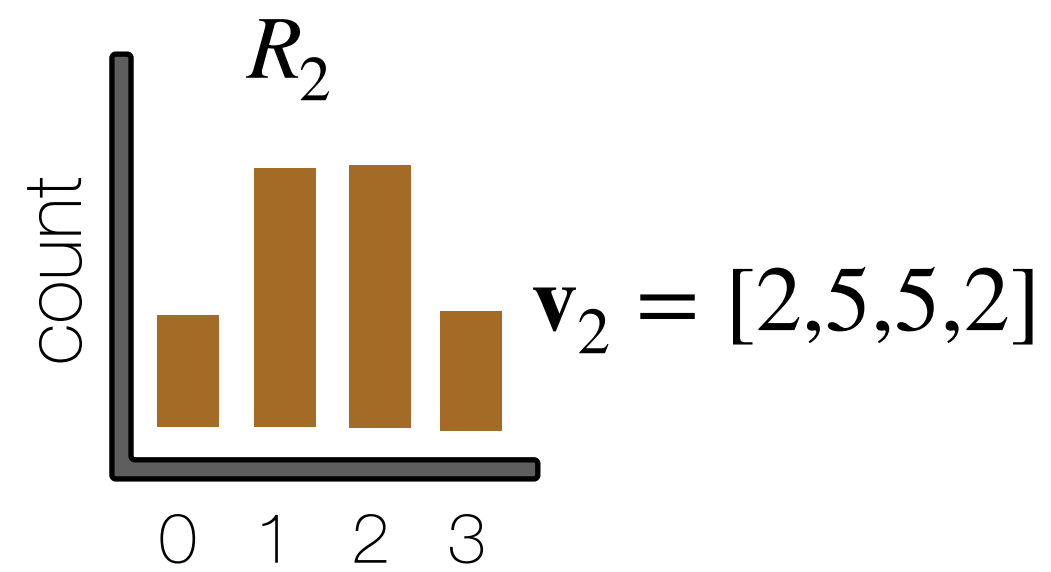
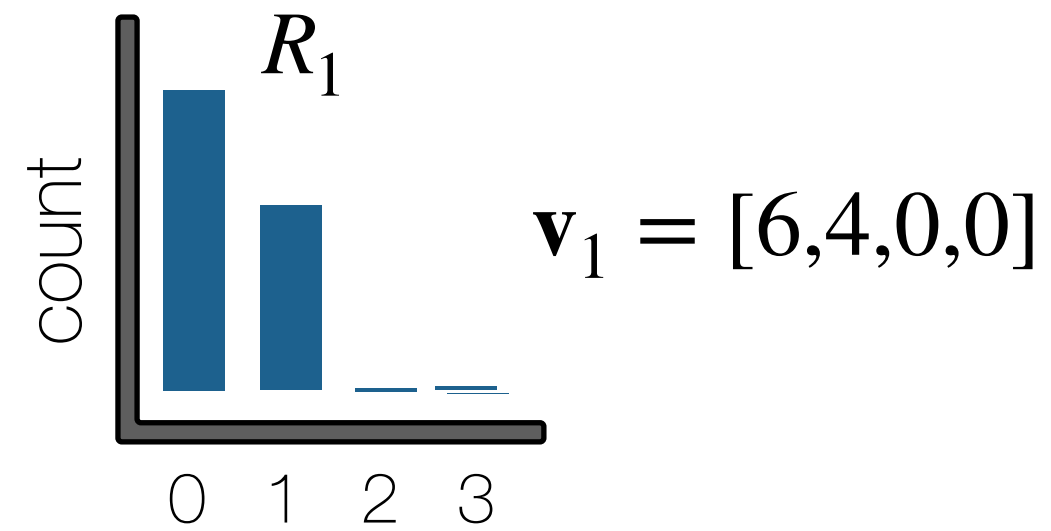


...

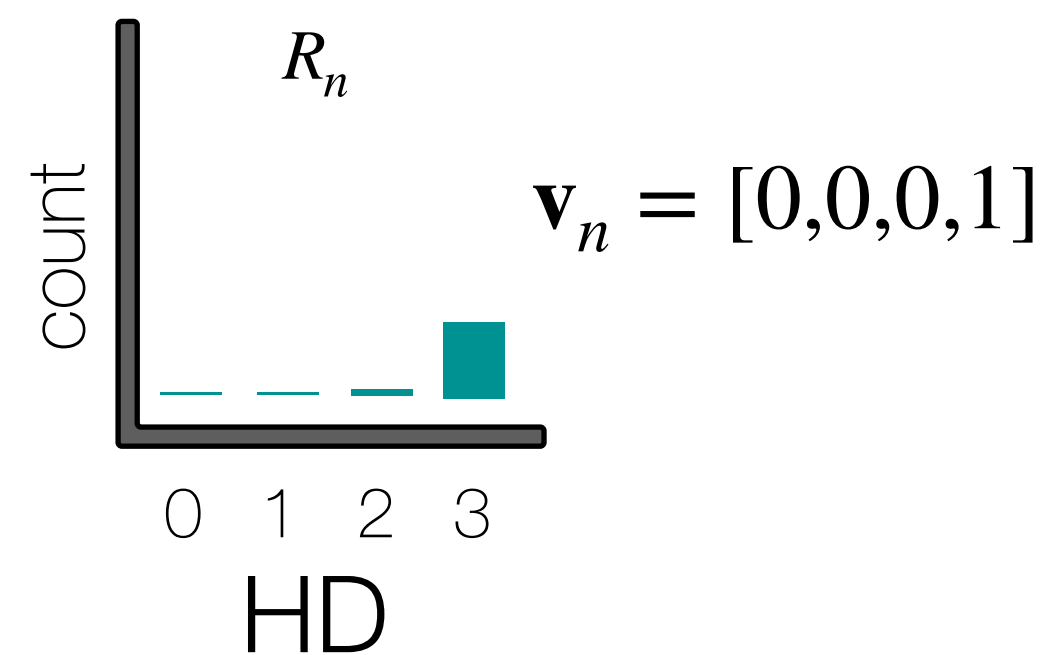


# Maximum likelihood estimation of distances

## Hamming distance histograms



...

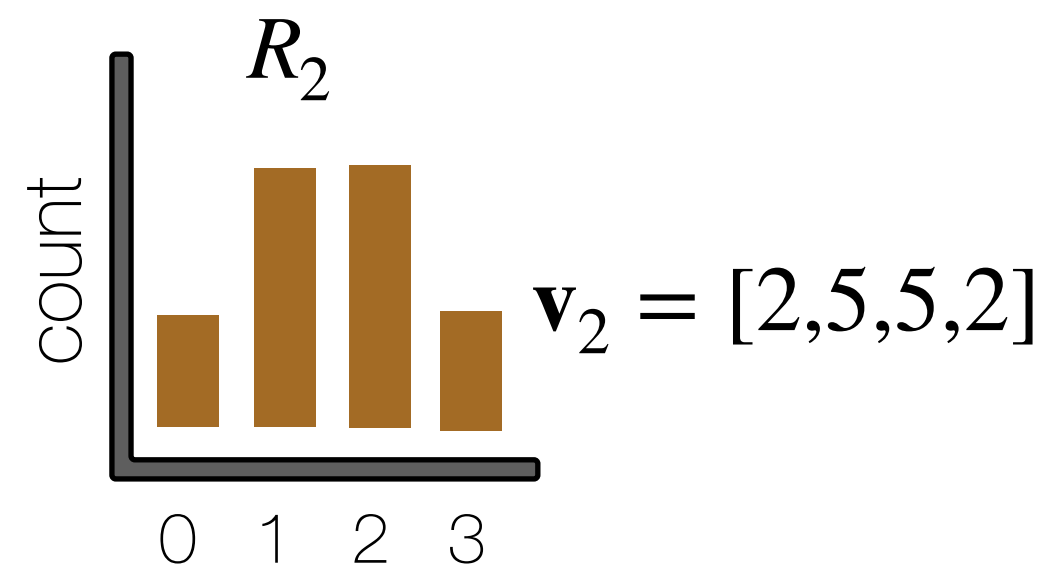
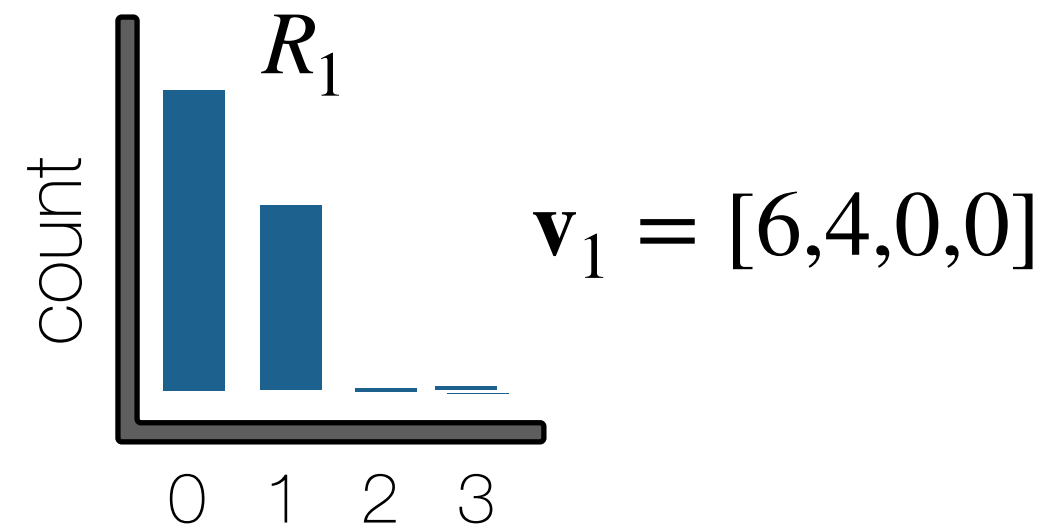


Optimize  $-\log \mathcal{L}_i$  w.r.t.  $D$  for each hitting reference  $R_i$ :

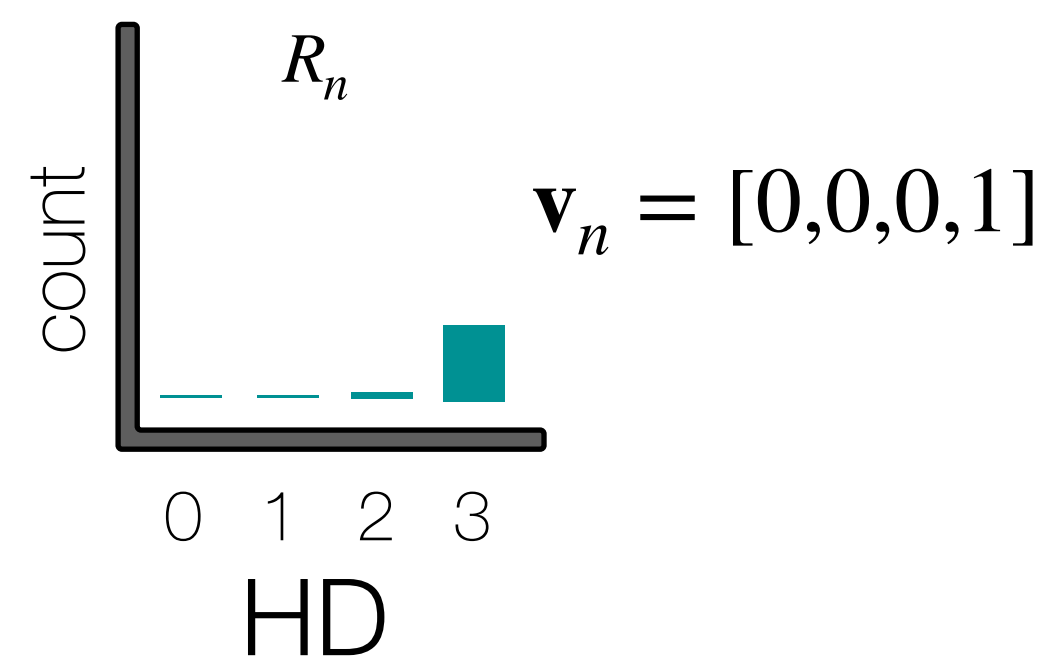
$$\arg \max_D u_i \log P_{miss}(D; k, h, \delta) + \sum_{d=0}^{\delta} v_{i,d} (d \log(D) + (k - d) \log(1 - D))$$

# Maximum likelihood estimation of distances

## Hamming distance histograms



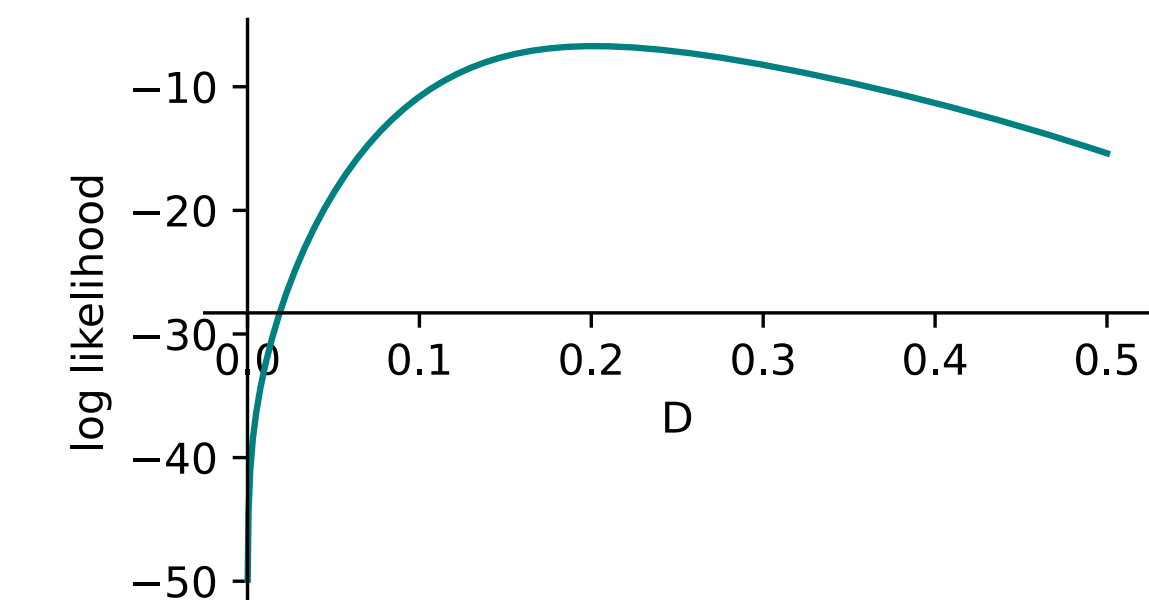
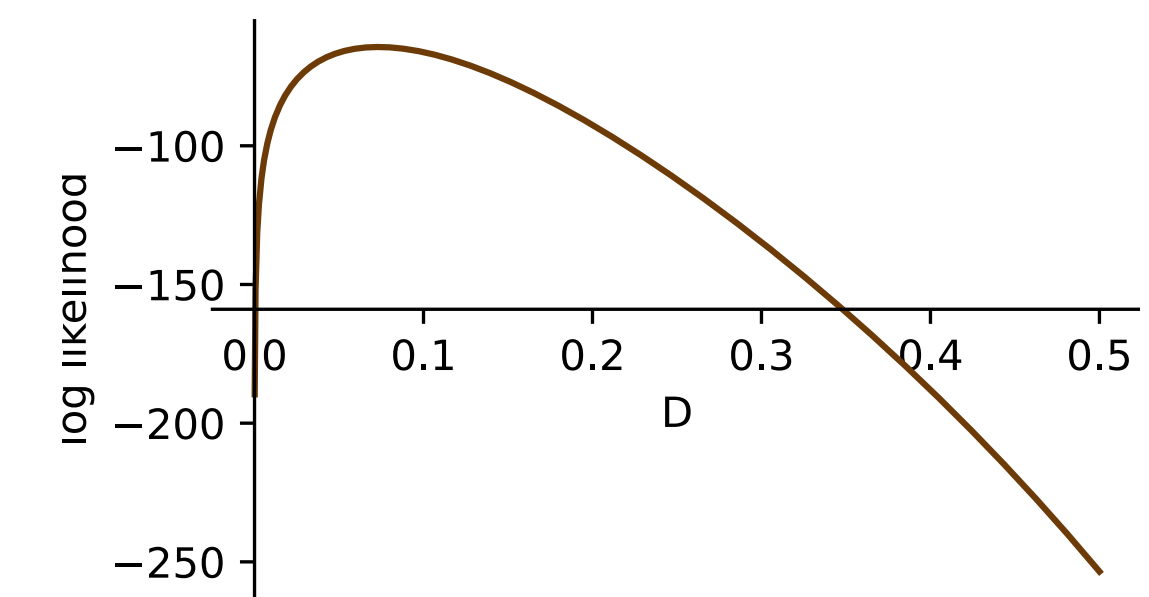
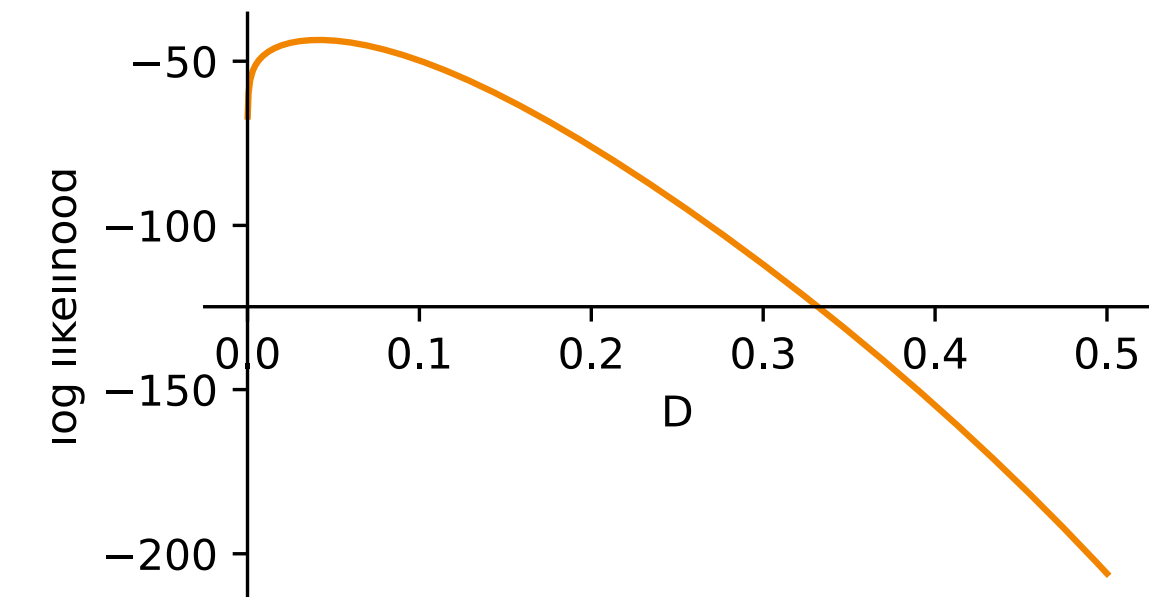
...



single variable & convex with a sensible choice of parameters

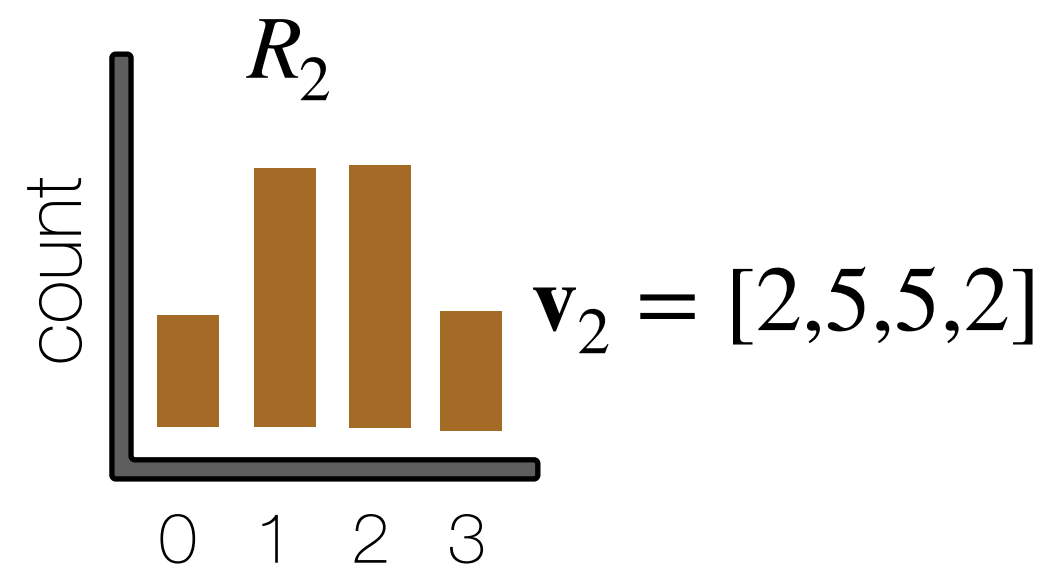
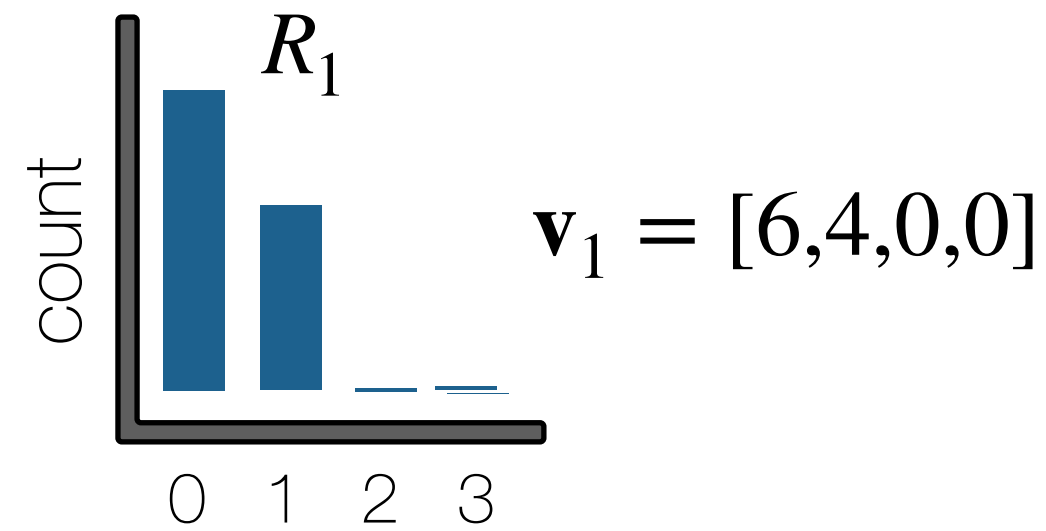
Optimize  $-\log \mathcal{L}_i$  w.r.t.  $D$  for each hitting reference  $R_i$ :

$$\arg \max_D u_i \log P_{miss}(D; k, h, \delta) + \sum_{d=0}^{\delta} v_{i,d} (d \log(D) + (k - d) \log(1 - D))$$

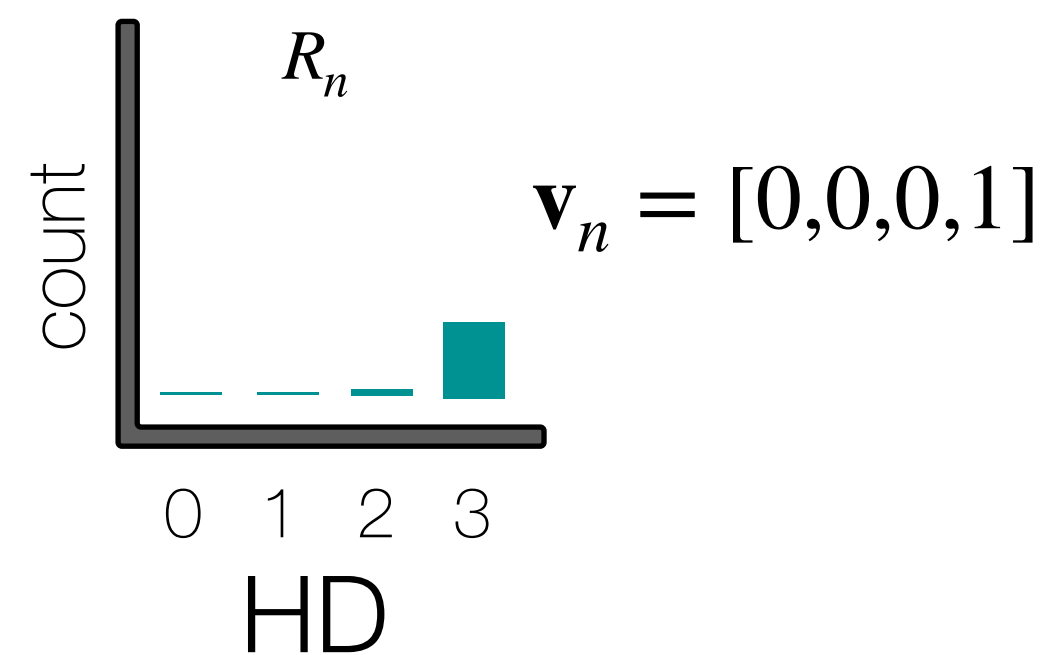


# Maximum likelihood estimation of distances

## Hamming distance histograms



...

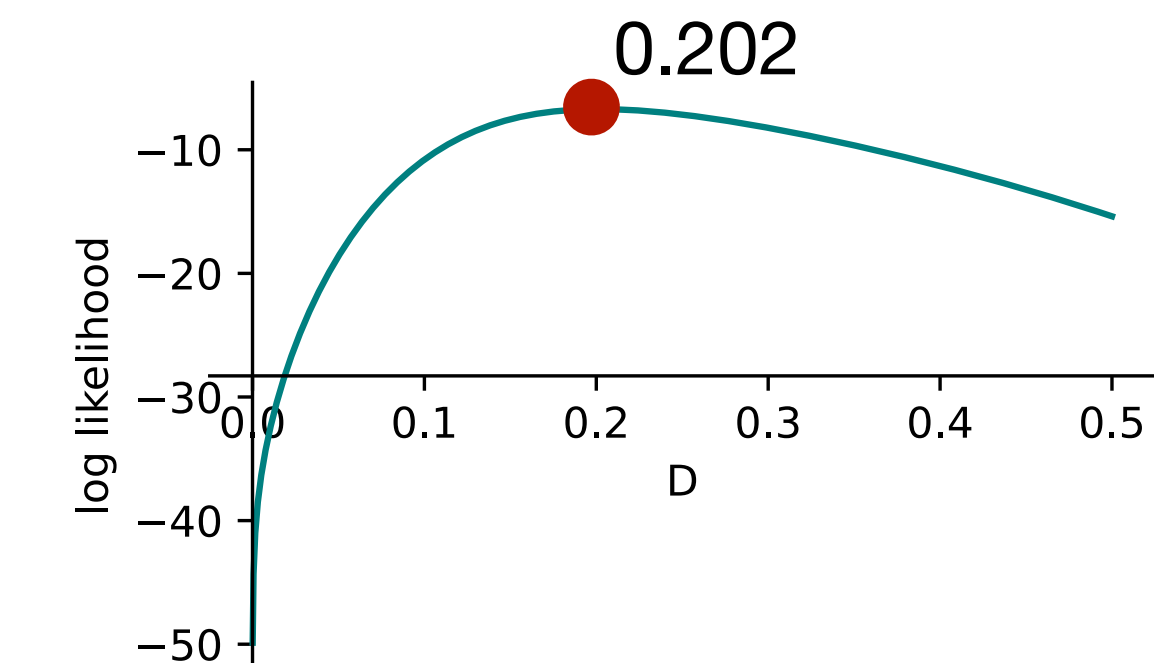
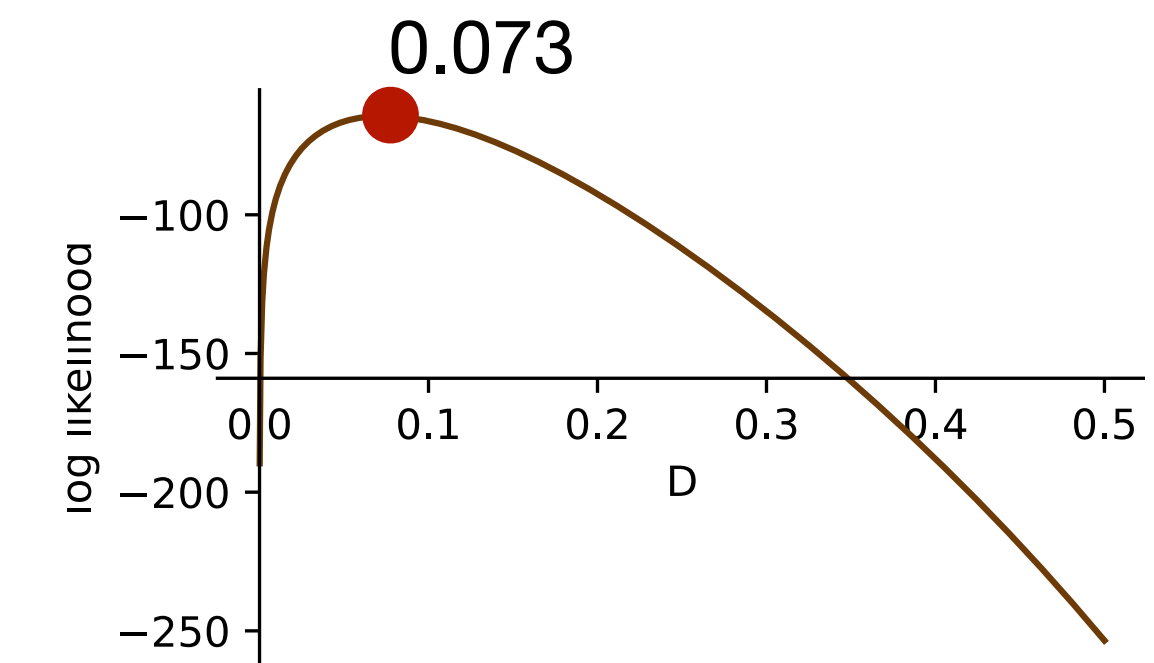
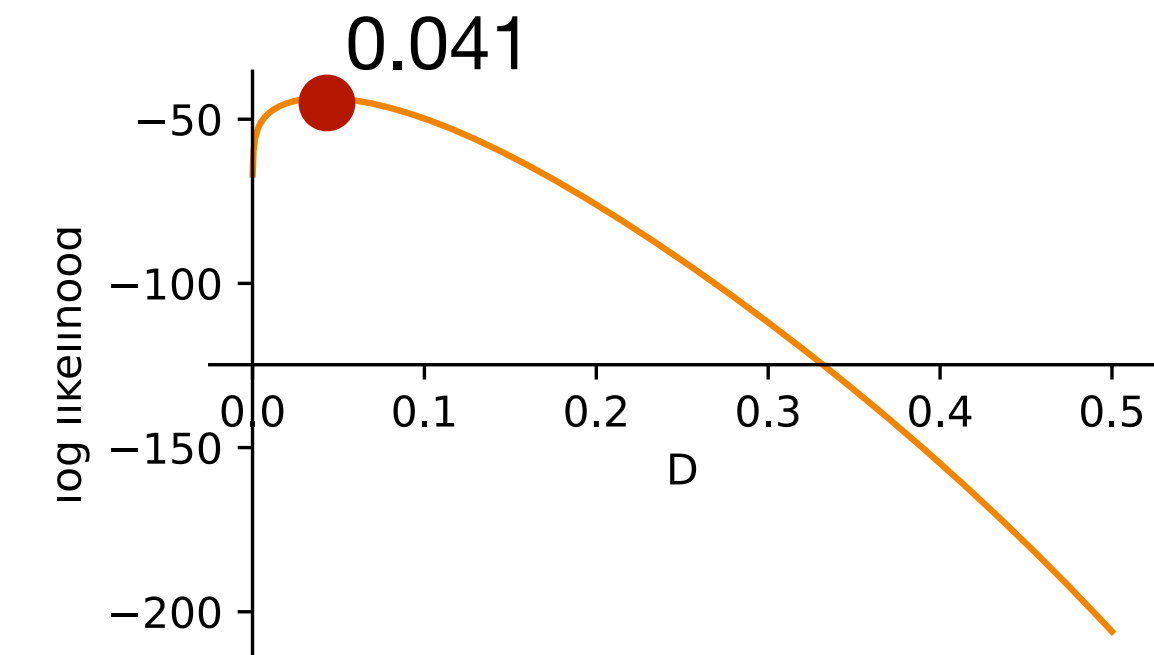


single variable & convex with a sensible choice of parameters

Optimize  $-\log \mathcal{L}_i$  w.r.t.  $D$  for each hitting reference  $R_i$ :

$$\arg \max_D u_i \log P_{miss}(D; k, h, \delta) + \sum_{d=0}^{\delta} v_{i,d} (d \log(D) + (k - d) \log(1 - D))$$

using Brent's Method



**krepp estimates read-level distances accurately**

# krepp estimates read-level distances accurately

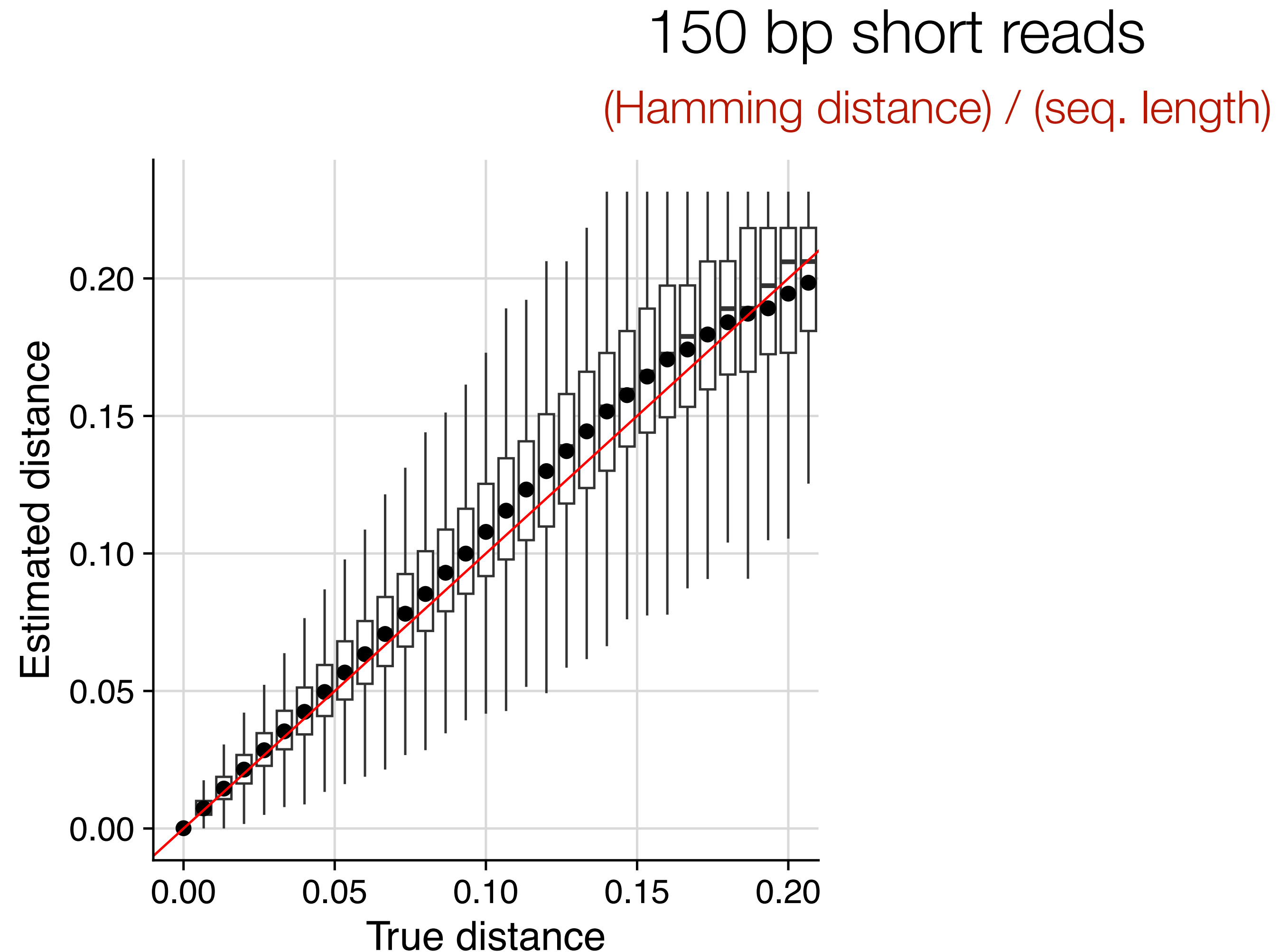
default: 29-mer  
minimizers of 35-mers

150 bp short reads  
(Hamming distance) / (seq. length)

# krepp estimates read-level distances accurately

default: 29-mer  
minimizers of 35-mers

- Simulated genomes and known coordinates
- **Highly accurate** (despite some noise)



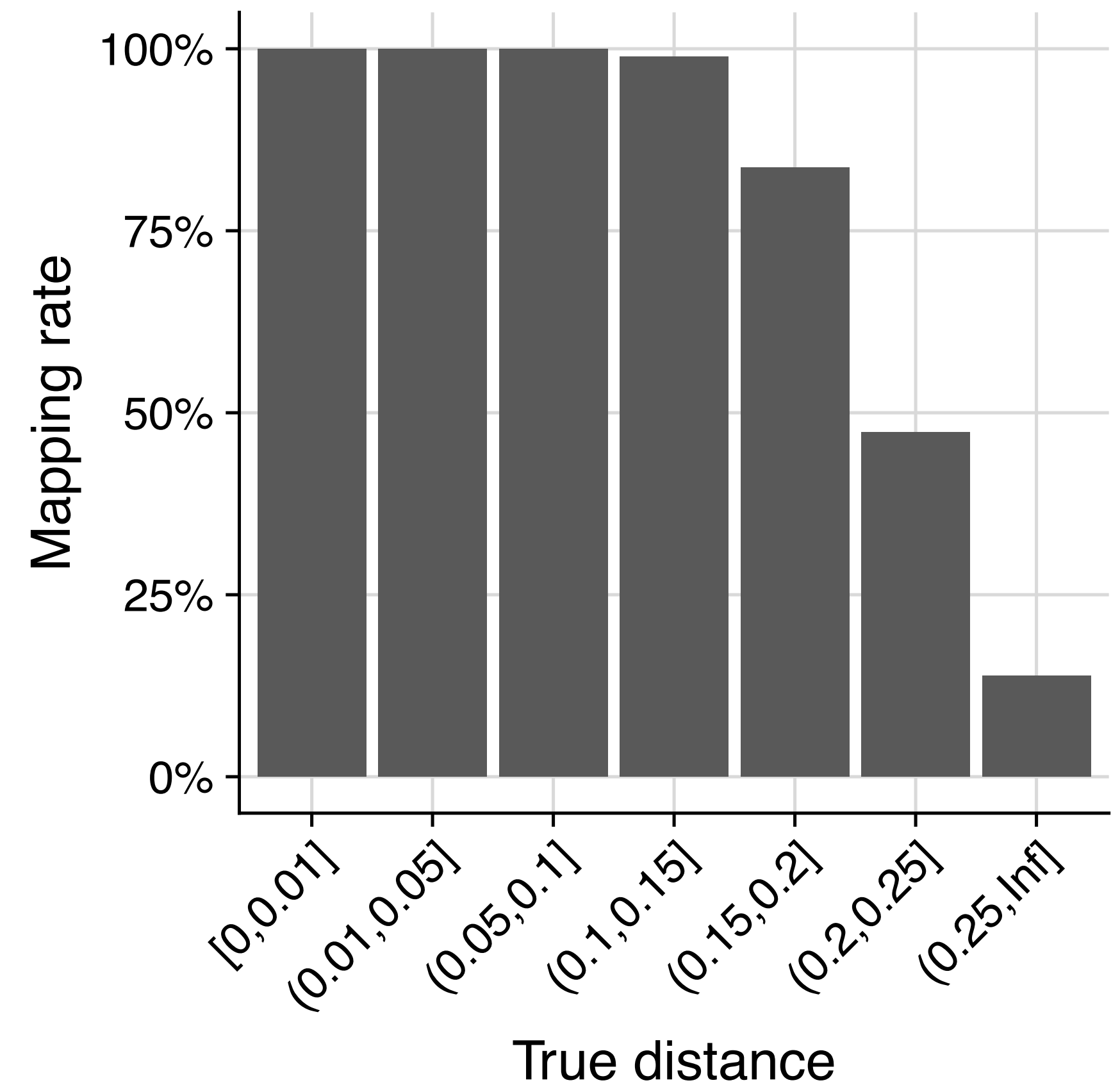
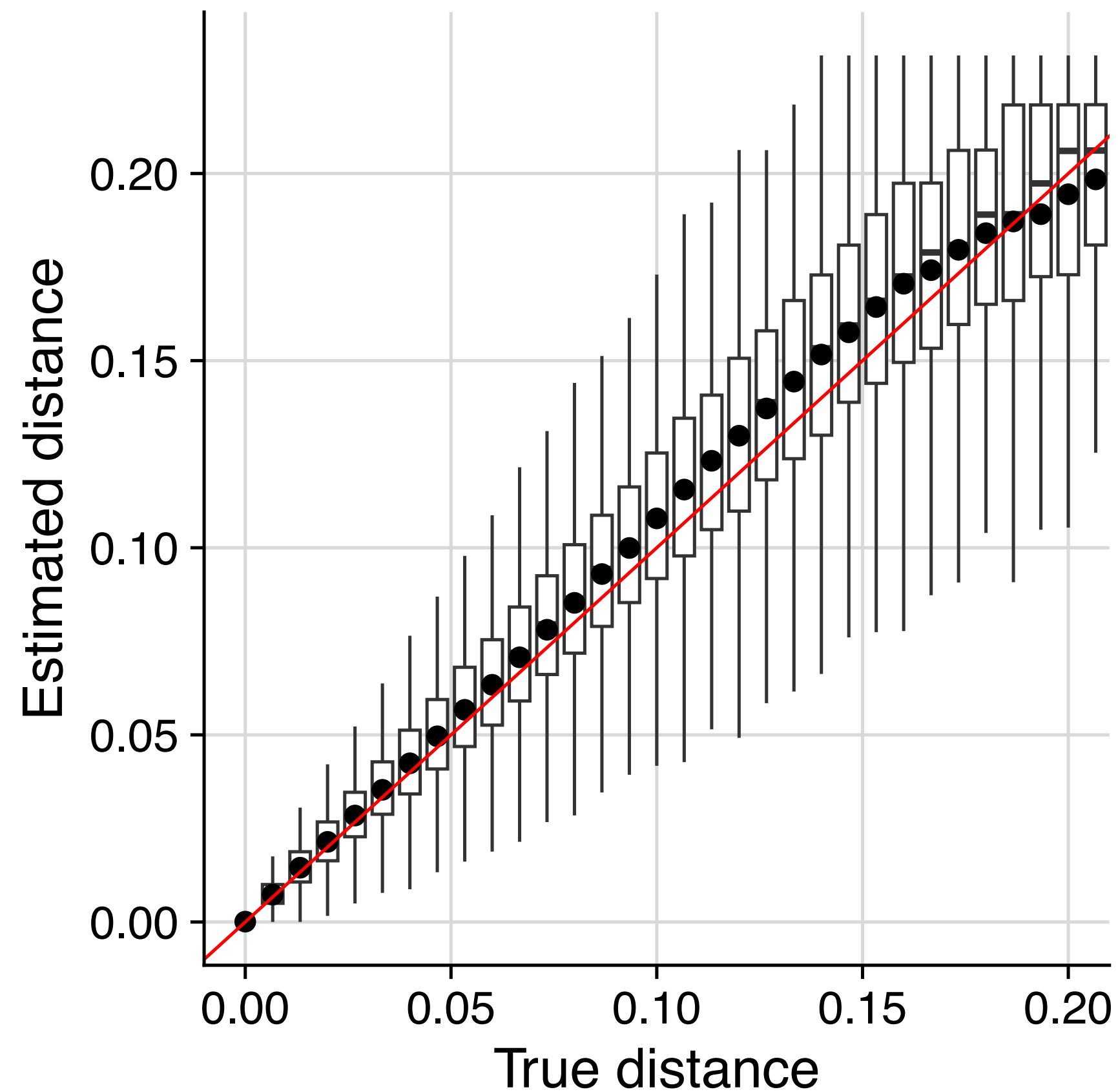
# krepp estimates read-level distances accurately

default: 29-mer  
minimizers of 35-mers

- Simulated genomes and known coordinates
- **Highly accurate** (despite some noise)

150 bp short reads

(Hamming distance) / (seq. length)



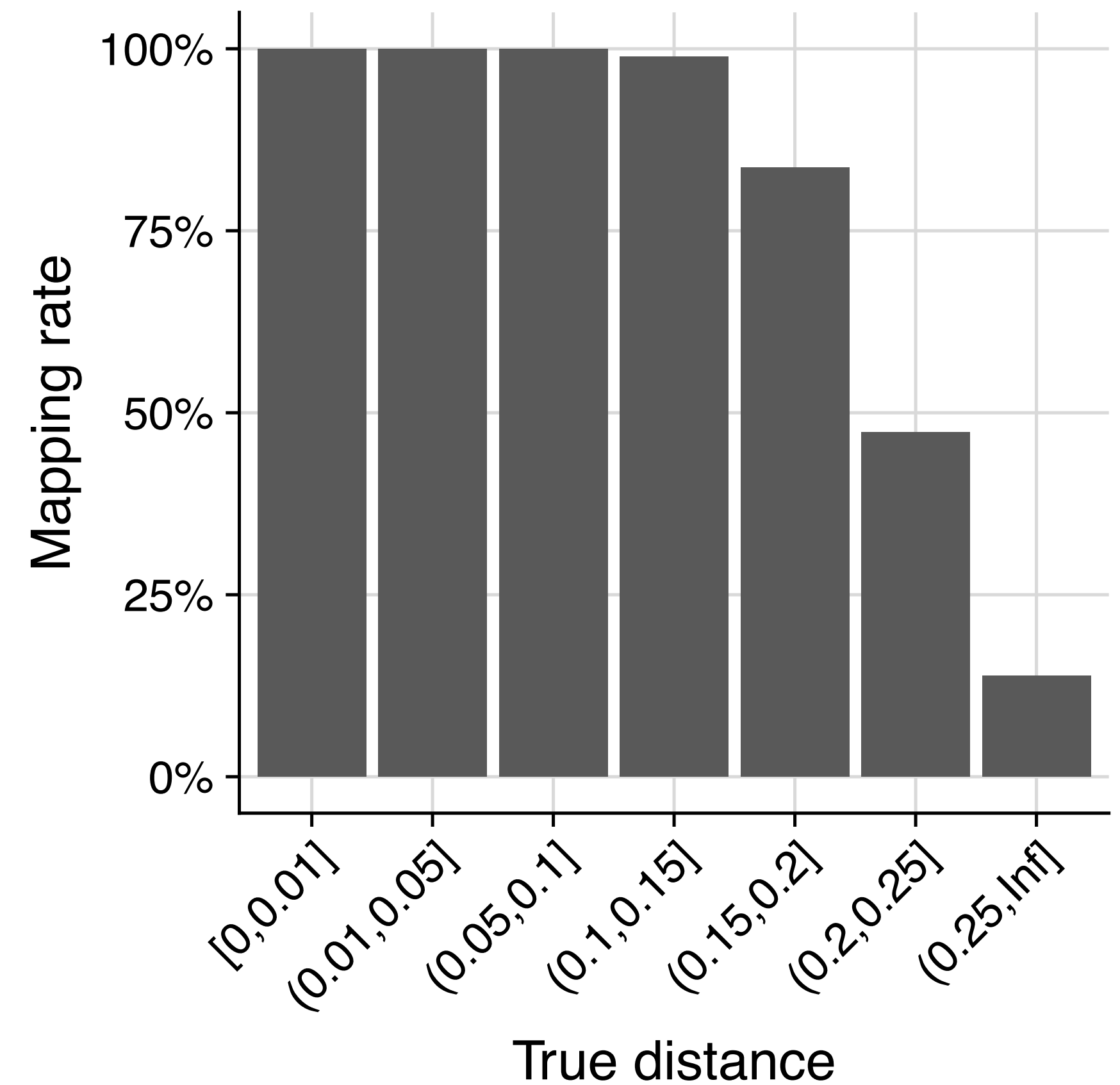
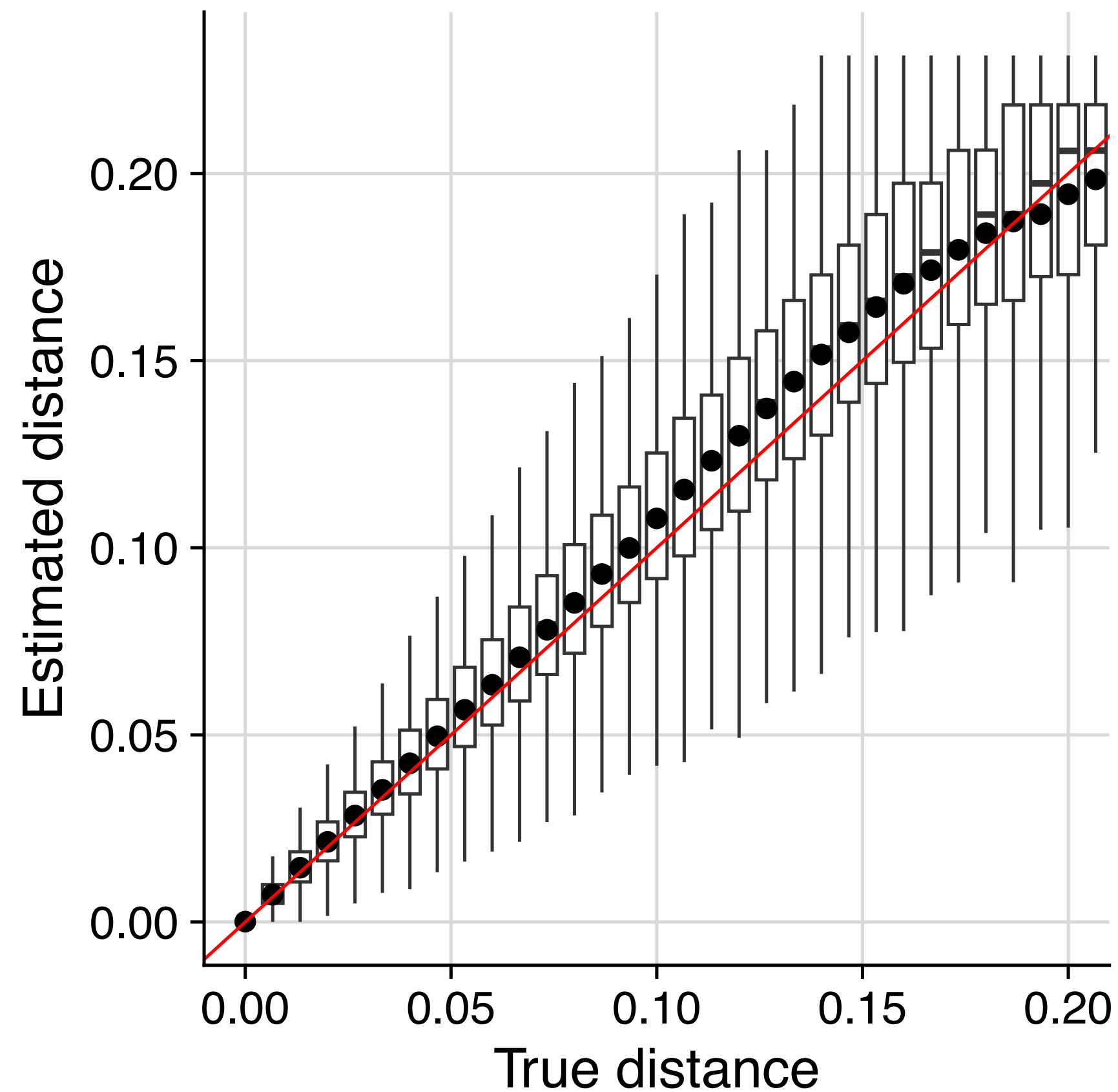
# krepp estimates read-level distances accurately

default: 29-mer  
minimizers of 35-mers

- Simulated genomes and known coordinates
- **Highly accurate** (despite some noise)
- **Slight overestimation** bias for high distances

150 bp short reads

(Hamming distance) / (seq. length)



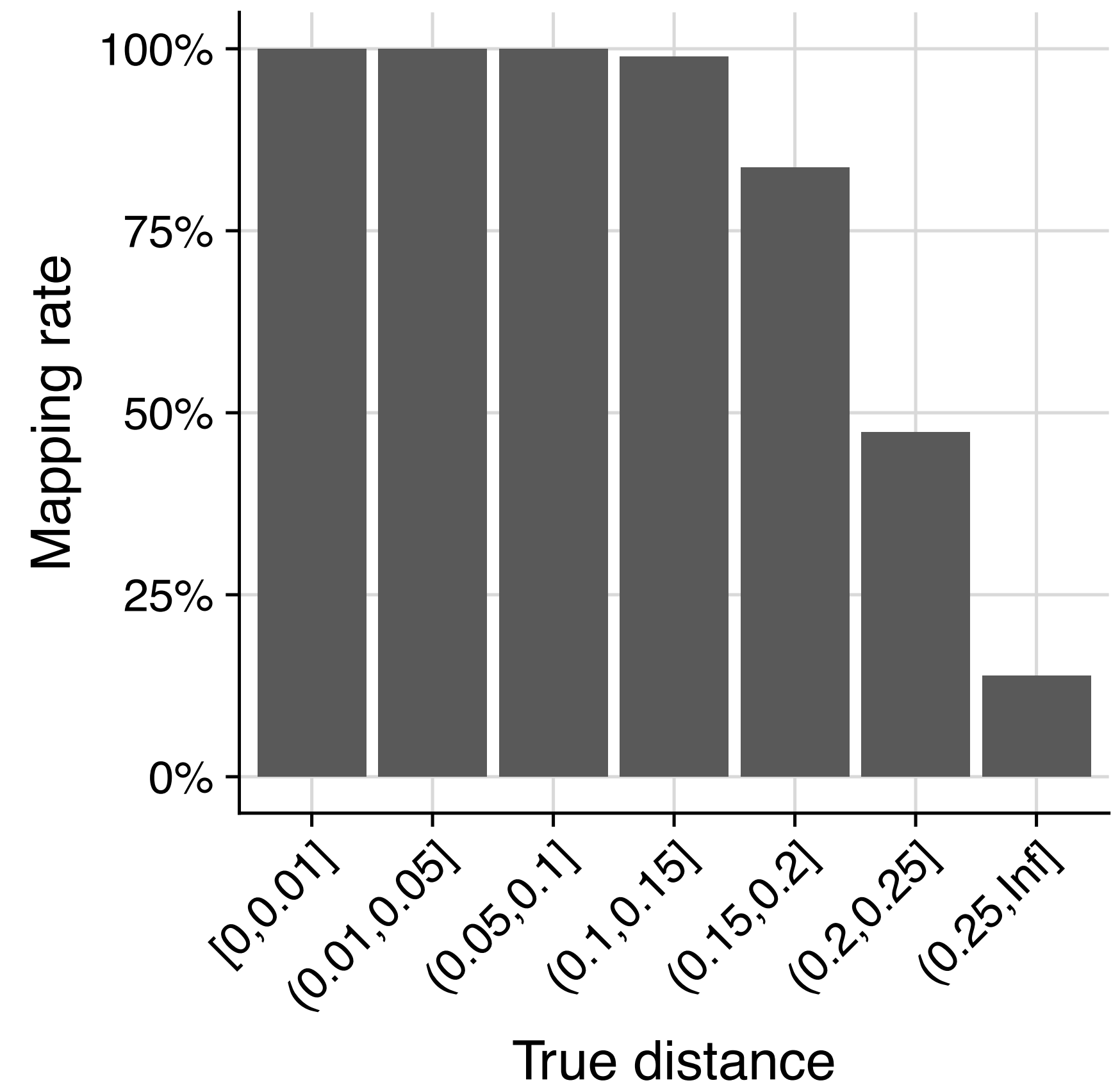
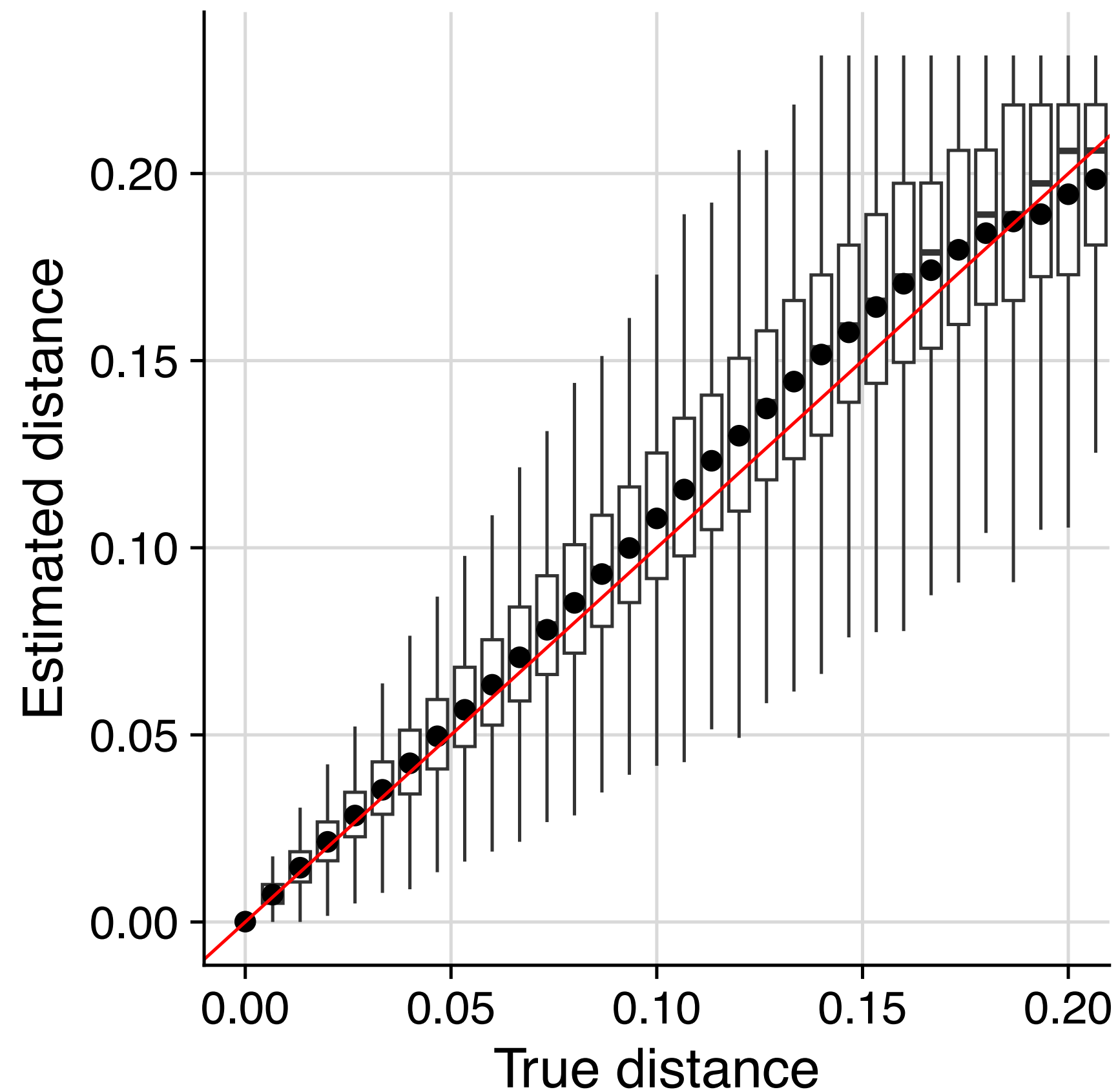
# krepp estimates read-level distances accurately

default: 29-mer  
minimizers of 35-mers

- Simulated genomes and known coordinates
- **Highly accurate** (despite some noise)
- **Slight overestimation** bias for high distances
- subsampling  $\uparrow$   
bias  $\downarrow$ , noise  $\uparrow$   
mapping rate  $\downarrow$

150 bp short reads

(Hamming distance) / (seq. length)

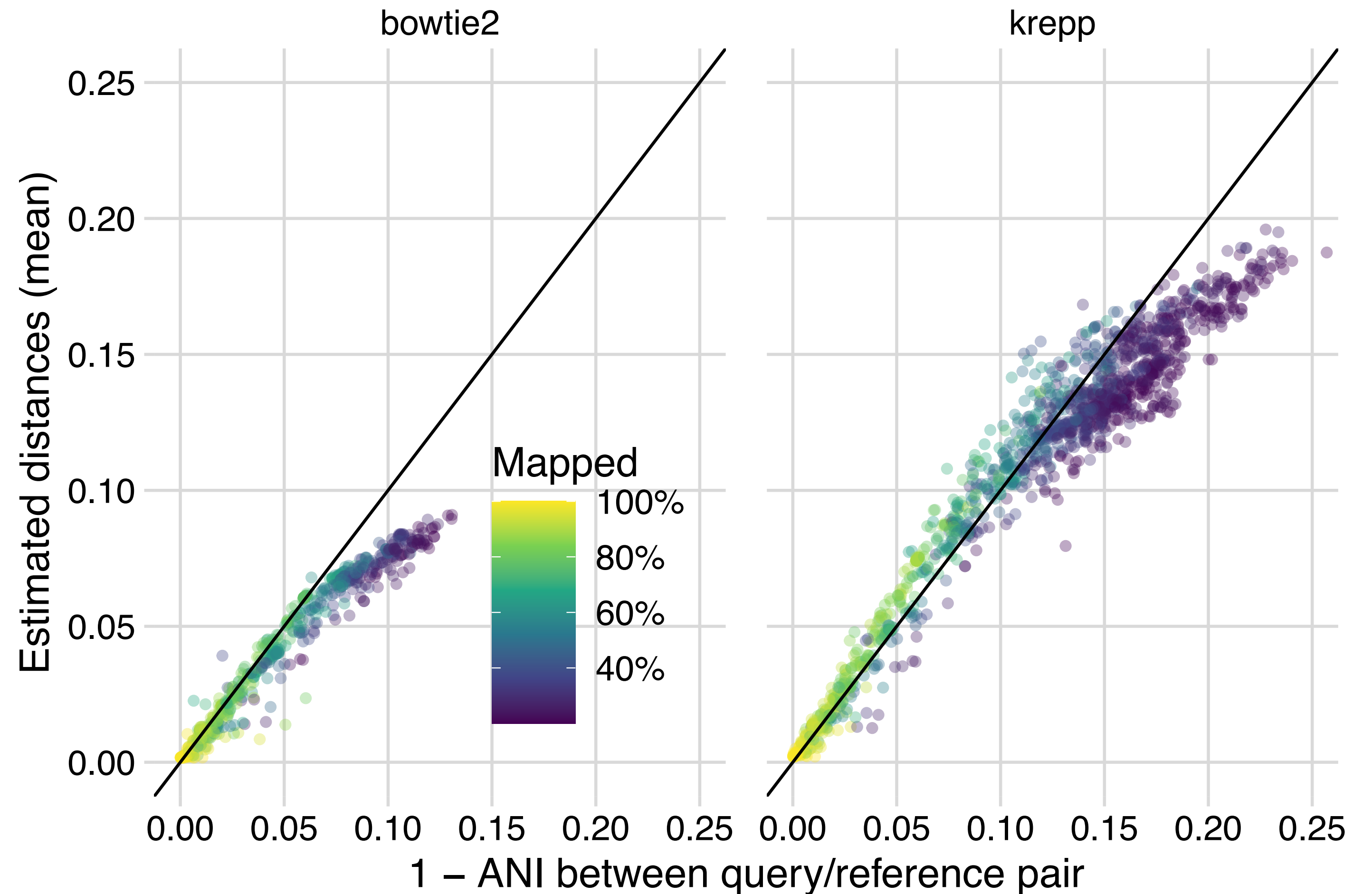


**krepp matches genome-wide nucleotide identity on average**

# krepp matches genome-wide nucleotide identity on average

Reference: Web of Life (v2)  
16,000 microbial genomes

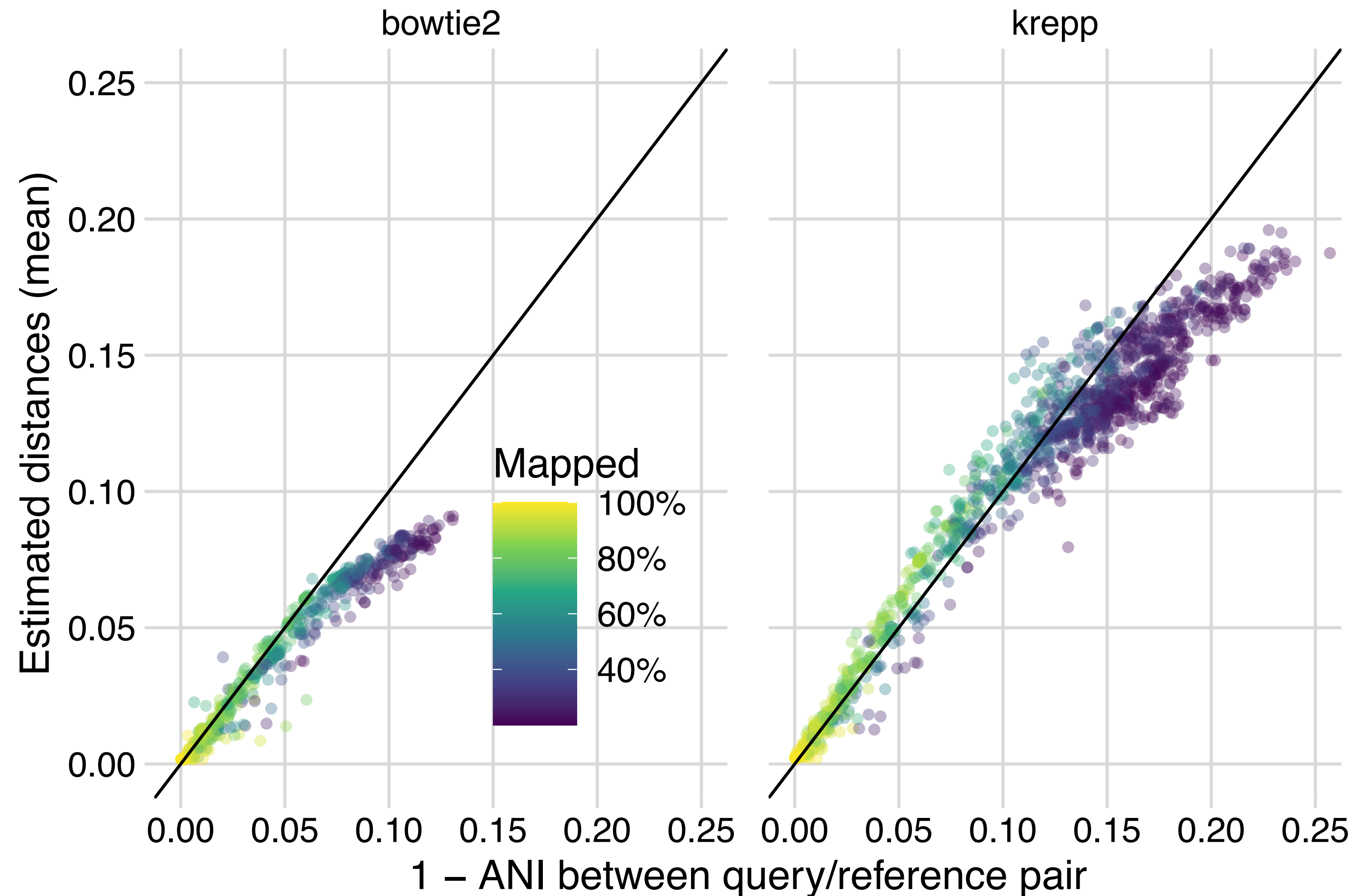
- Real query/reference genomes w/ simulated short reads



# krepp matches genome-wide nucleotide identity on average

Reference: Web of Life (v2)  
16,000 microbial genomes

- Real query/reference genomes w/ simulated short reads

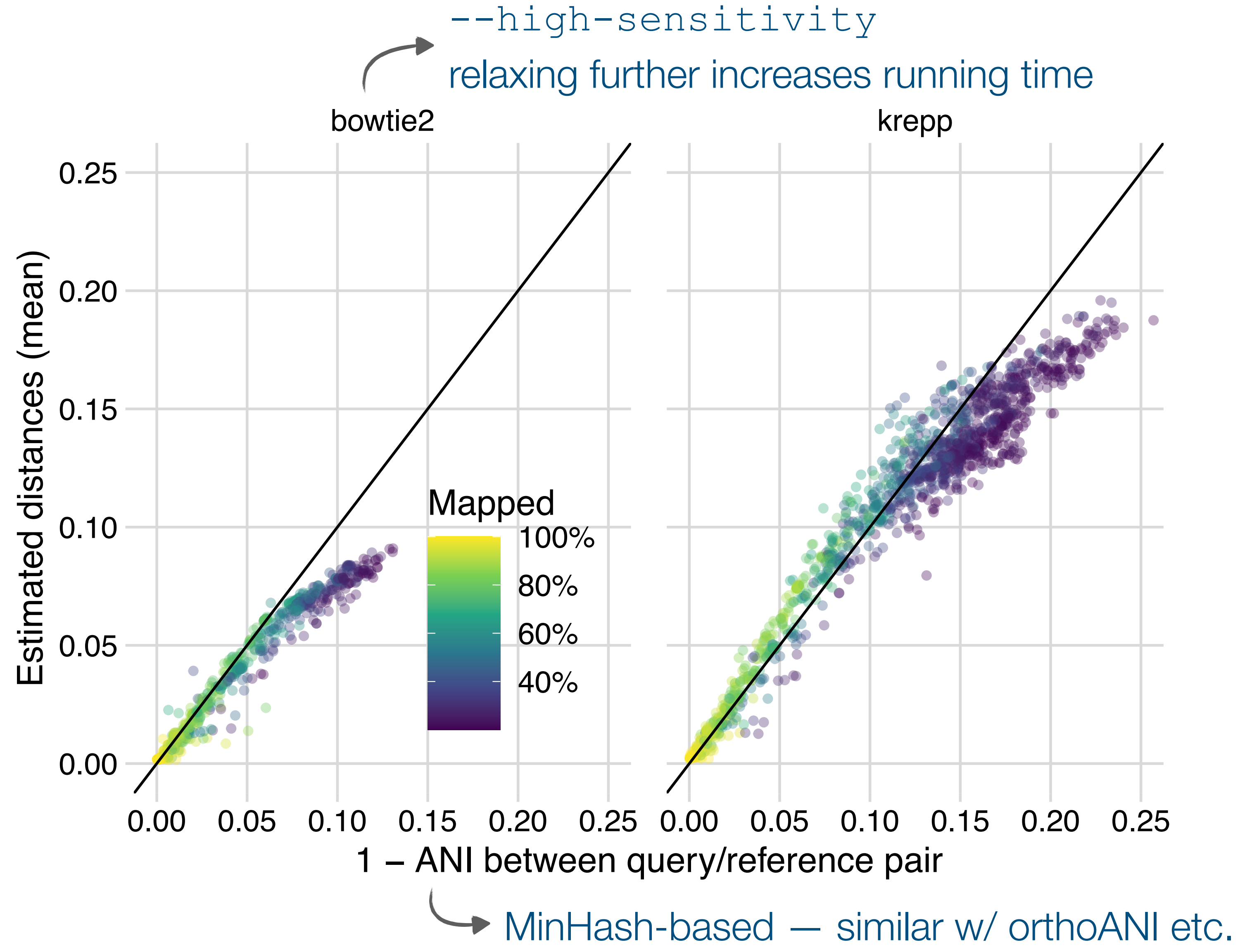


↪ MinHash-based — similar w/ orthoANI etc.

# krepp matches genome-wide nucleotide identity on average

Reference: Web of Life (v2)  
16,000 microbial genomes

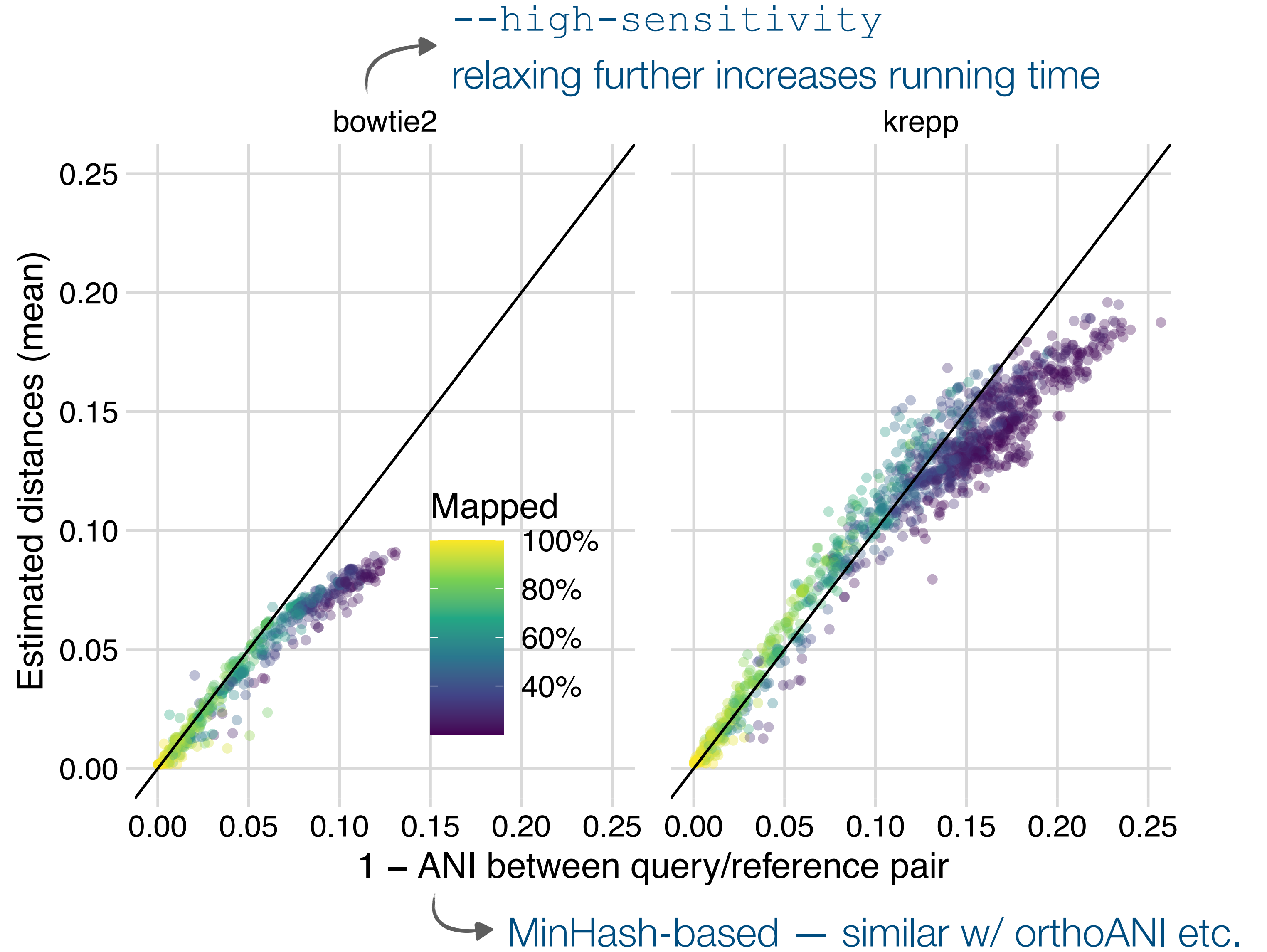
- Real query/reference genomes w/ simulated short reads



# krepp matches genome-wide nucleotide identity on average

Reference: Web of Life (v2)  
16,000 microbial genomes

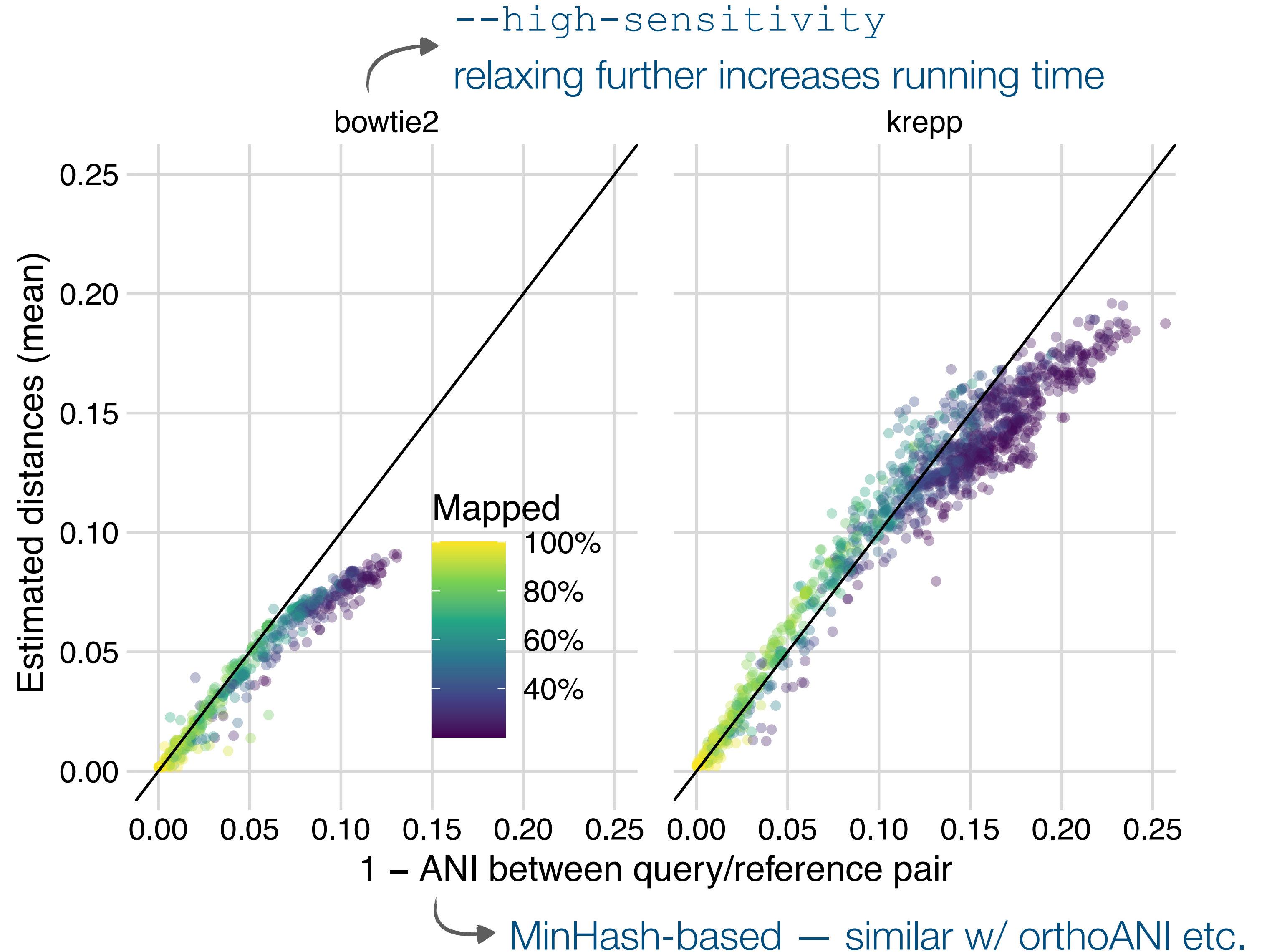
- Real query/reference genomes w/ simulated short reads
- *krepp* extends to more distant reference genomes accurately



# krepp matches genome-wide nucleotide identity on average

Reference: Web of Life (v2)  
16,000 microbial genomes

- Real query/reference genomes w/ simulated short reads
- *krepp* extends to more distant reference genomes accurately
- ANI↓: mapping rate ↓  
bias towards similar regions ↑



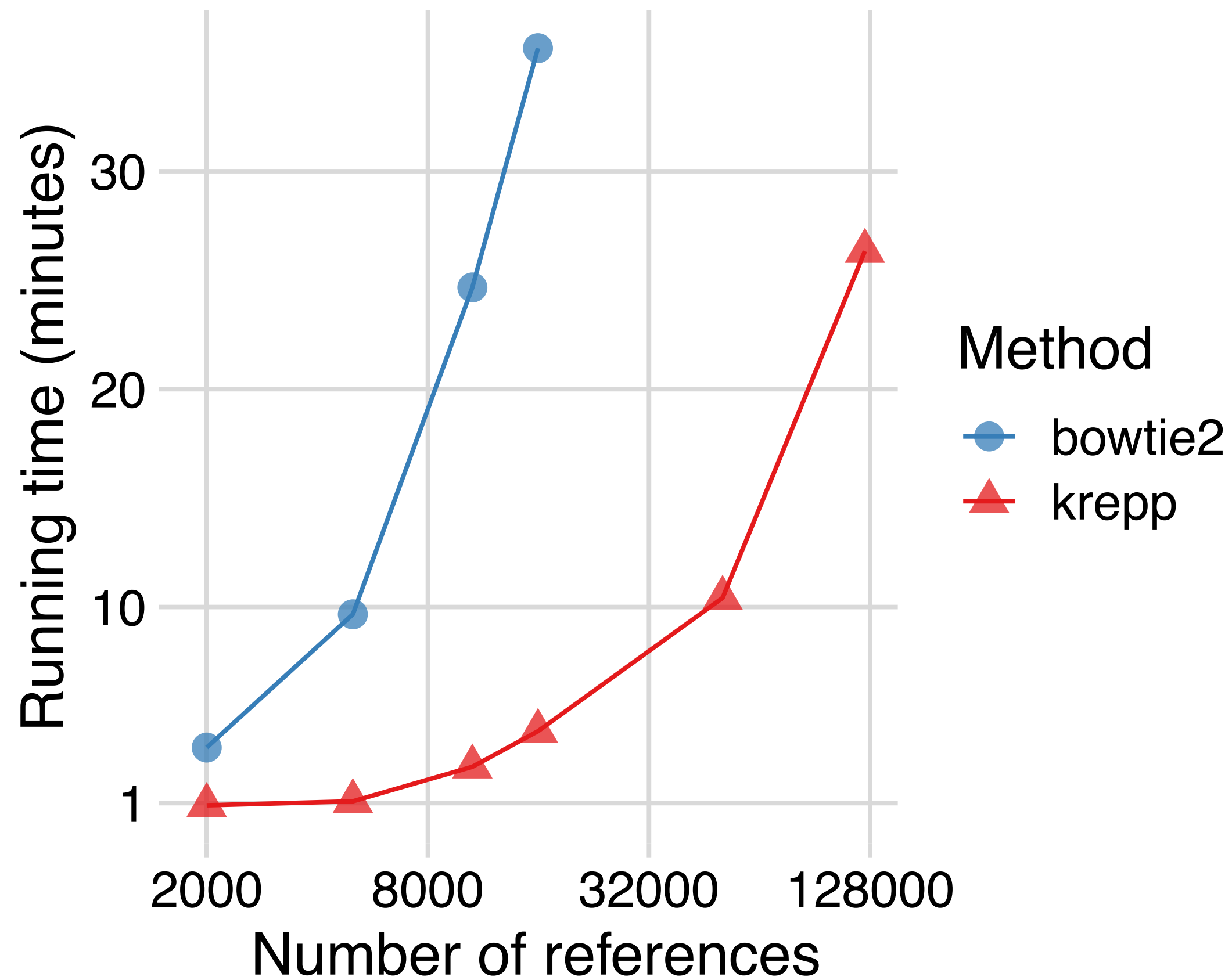
# **Scalability:**

**Avoiding the more difficult problem & effective parallelization**

# Scalability:

## Avoiding the more difficult problem & effective parallelization

Mapping 10M reads (16 threads):

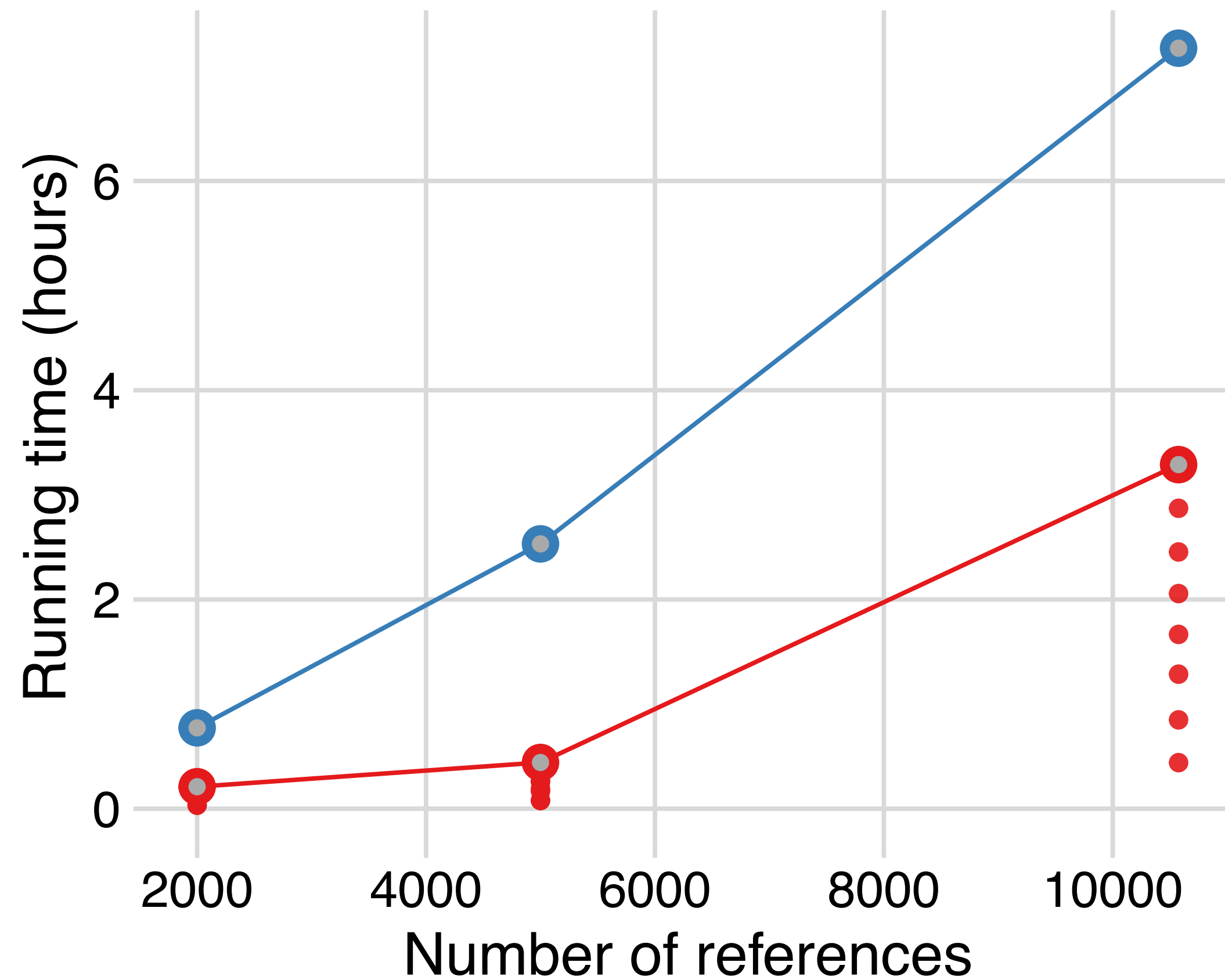
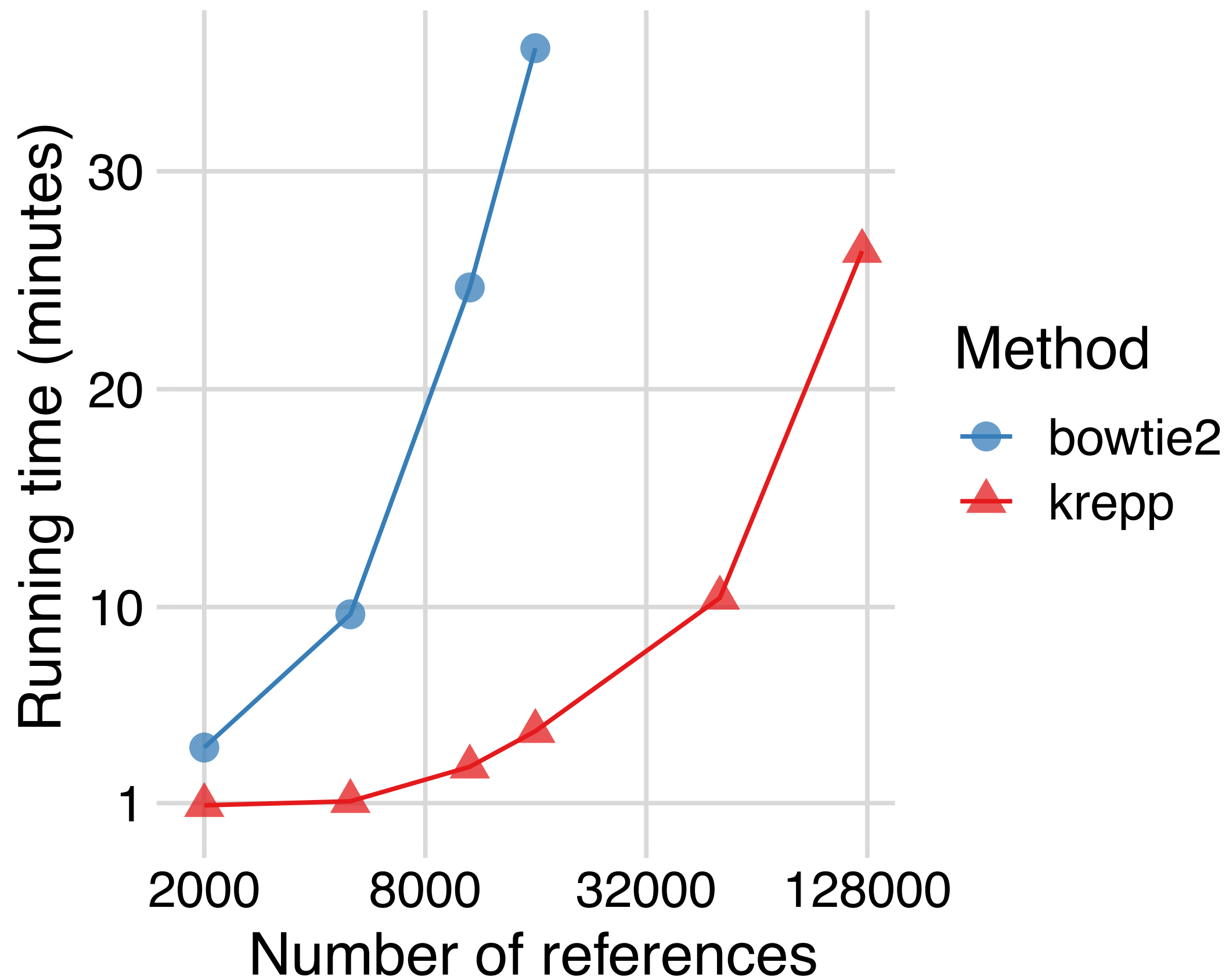


# Scalability:

## Avoiding the more difficult problem & effective parallelization

Mapping 10M reads (16 threads):

Indexing microbial genomes (32 threads):

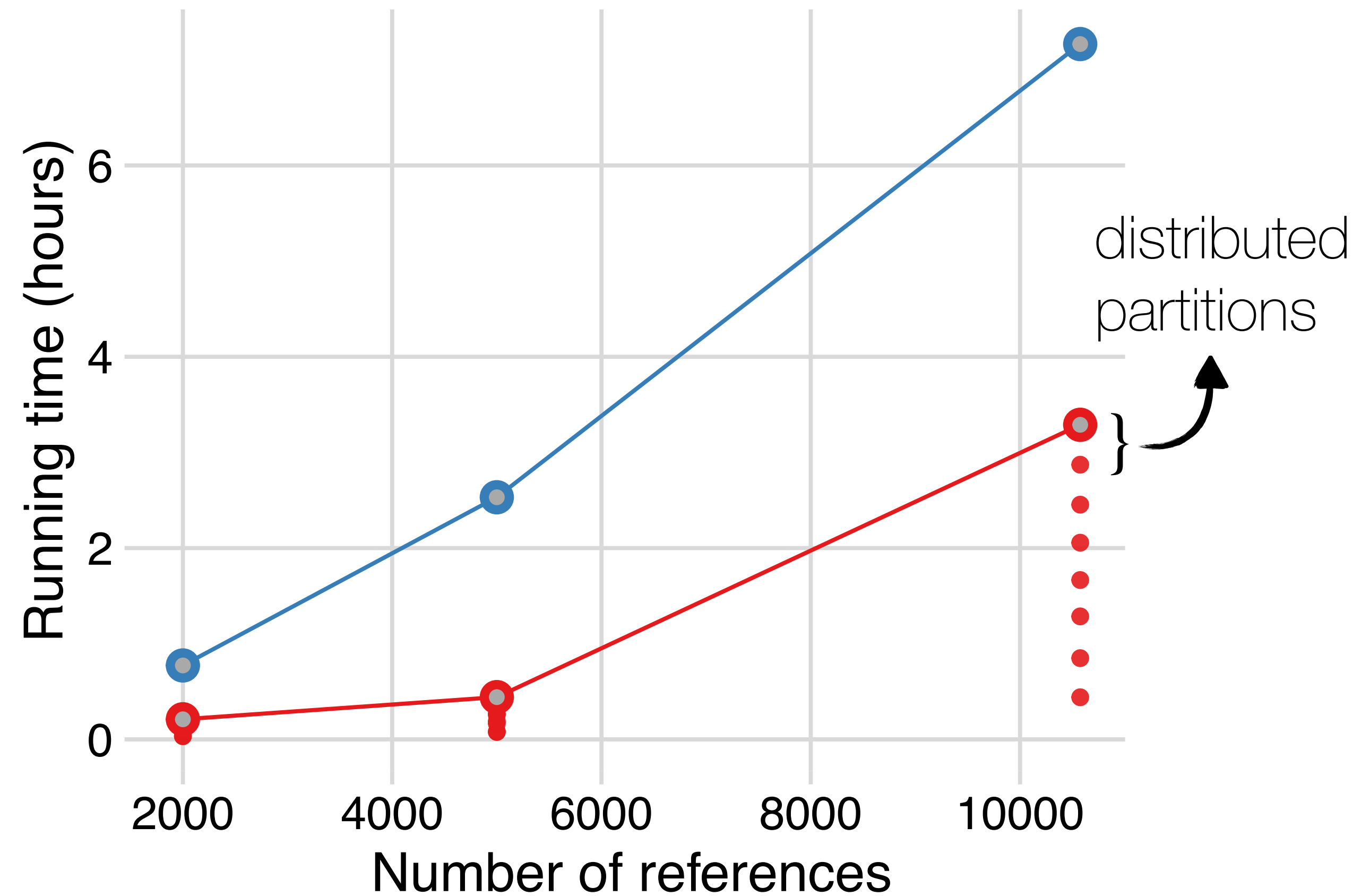
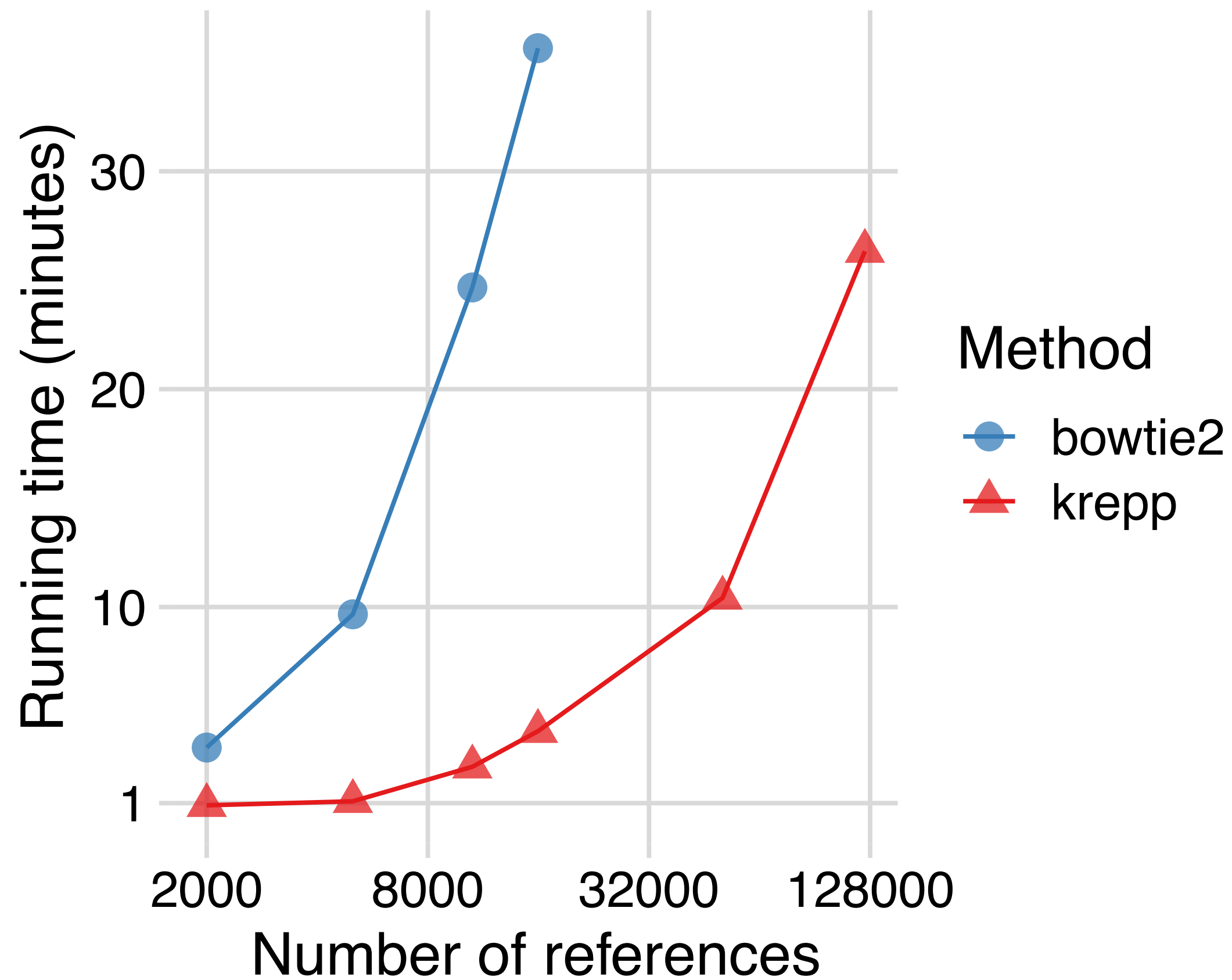


# Scalability:

## Avoiding the more difficult problem & effective parallelization

Mapping 10M reads (16 threads):

Indexing microbial genomes (32 threads):

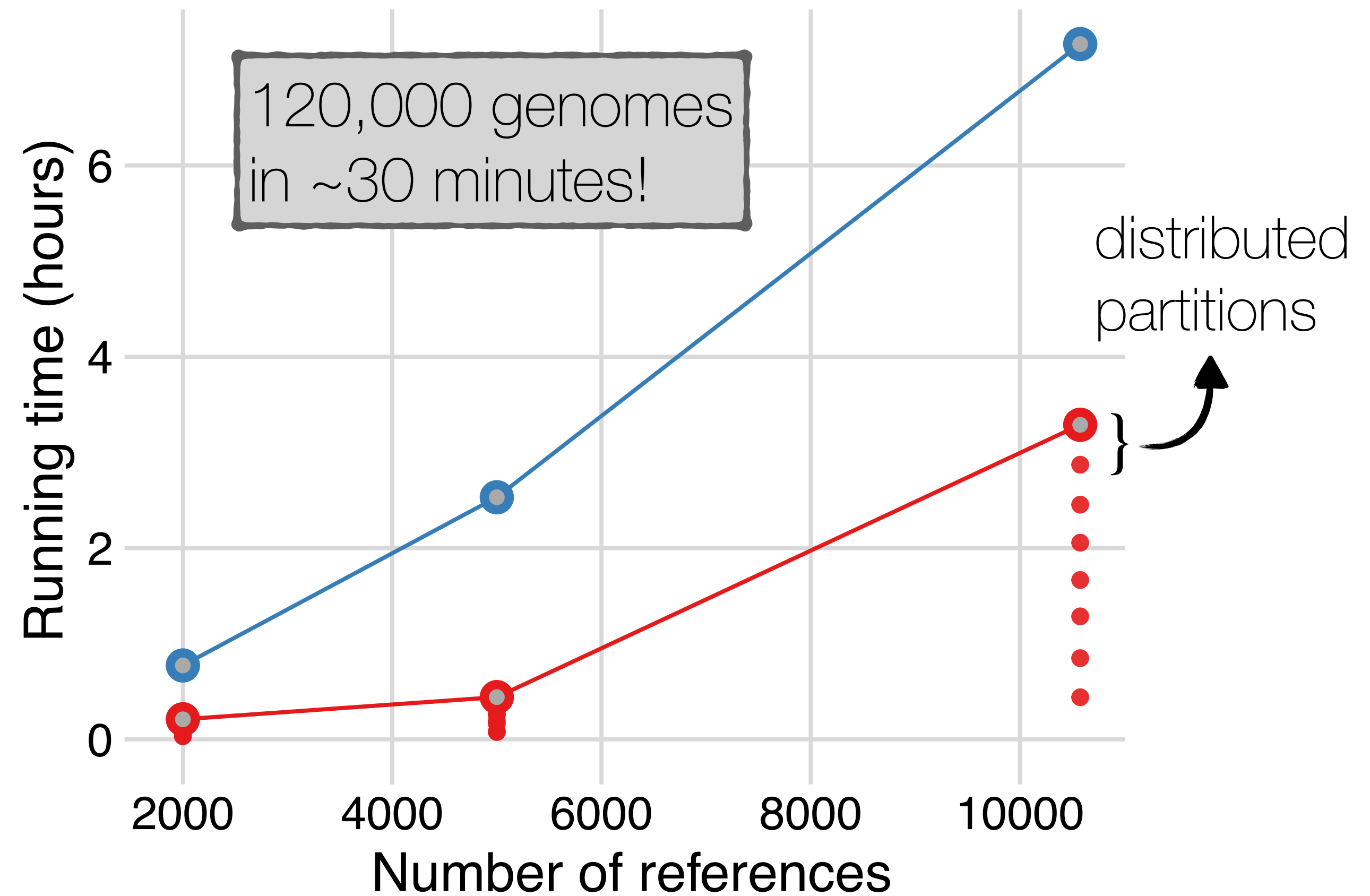
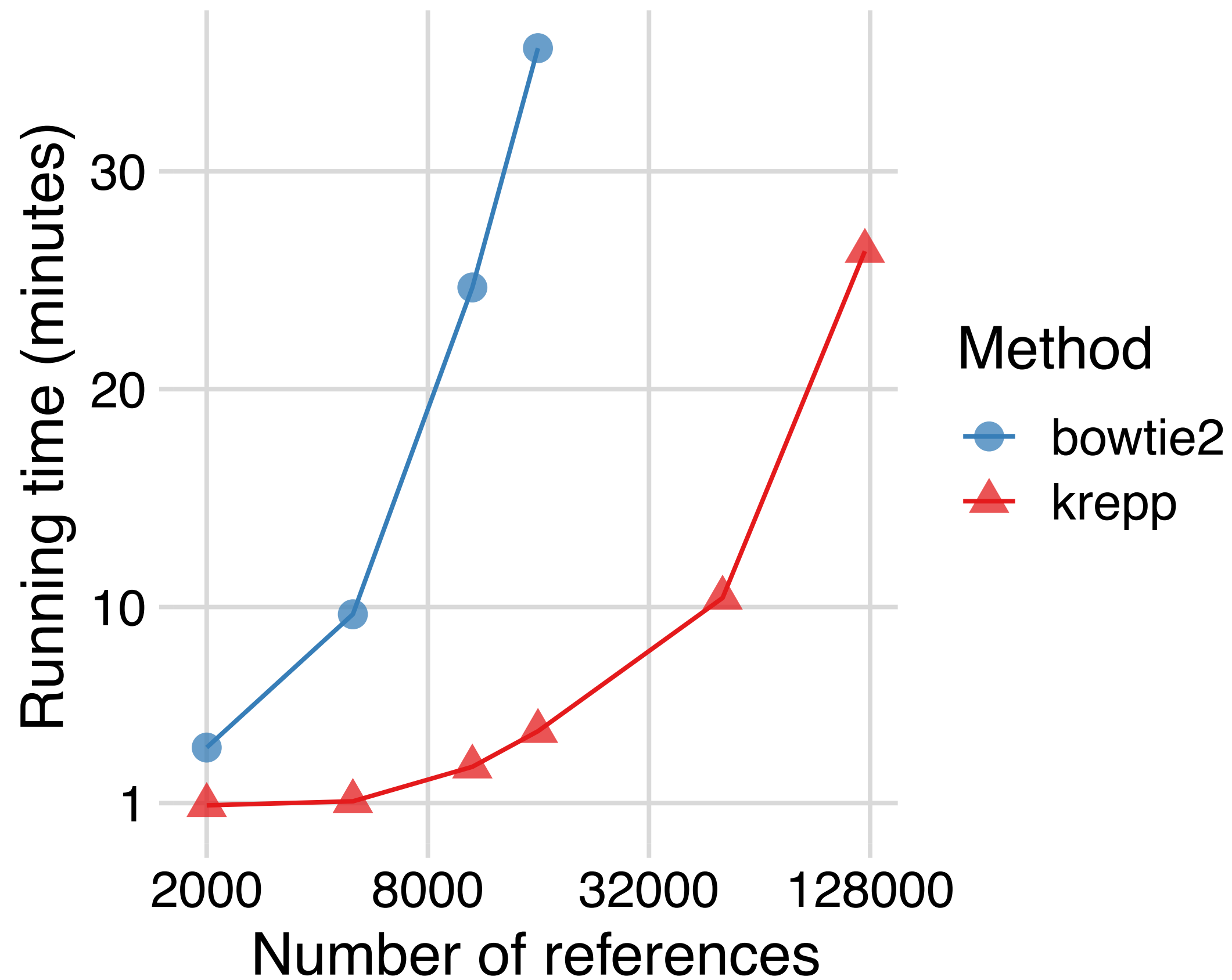


# Scalability:

## Avoiding the more difficult problem & effective parallelization

Mapping 10M reads (16 threads):

Indexing microbial genomes (32 threads):

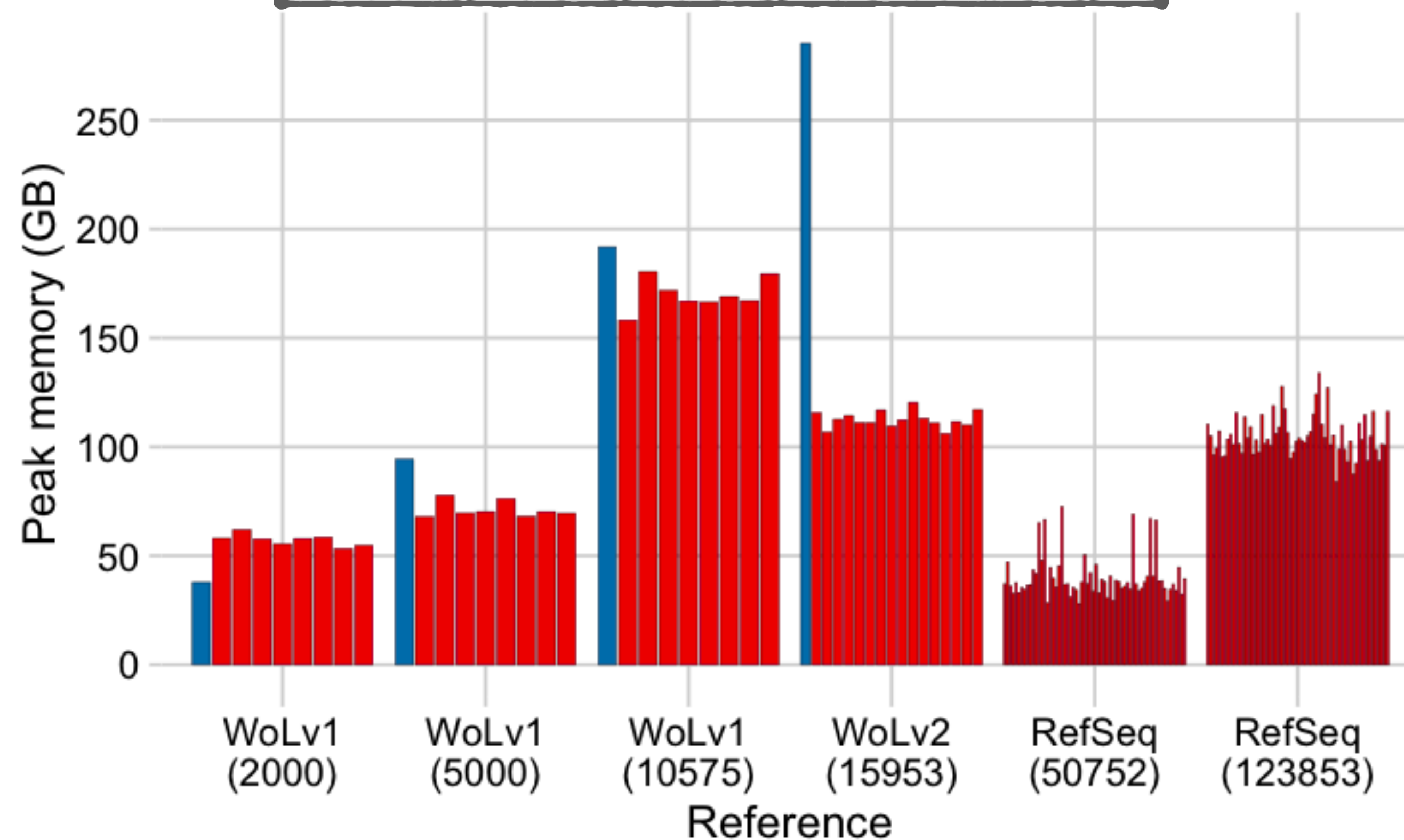


# Scalability:

**krepp can be distributed and has flexible memory requirements**

Indexing microbial genomes (32 threads):

adjusting partitioning based on the input size & available memory...



# **What about longer sequences?**

## **Taxonomic binning of contigs in CAMI-II**

# What about longer sequences?

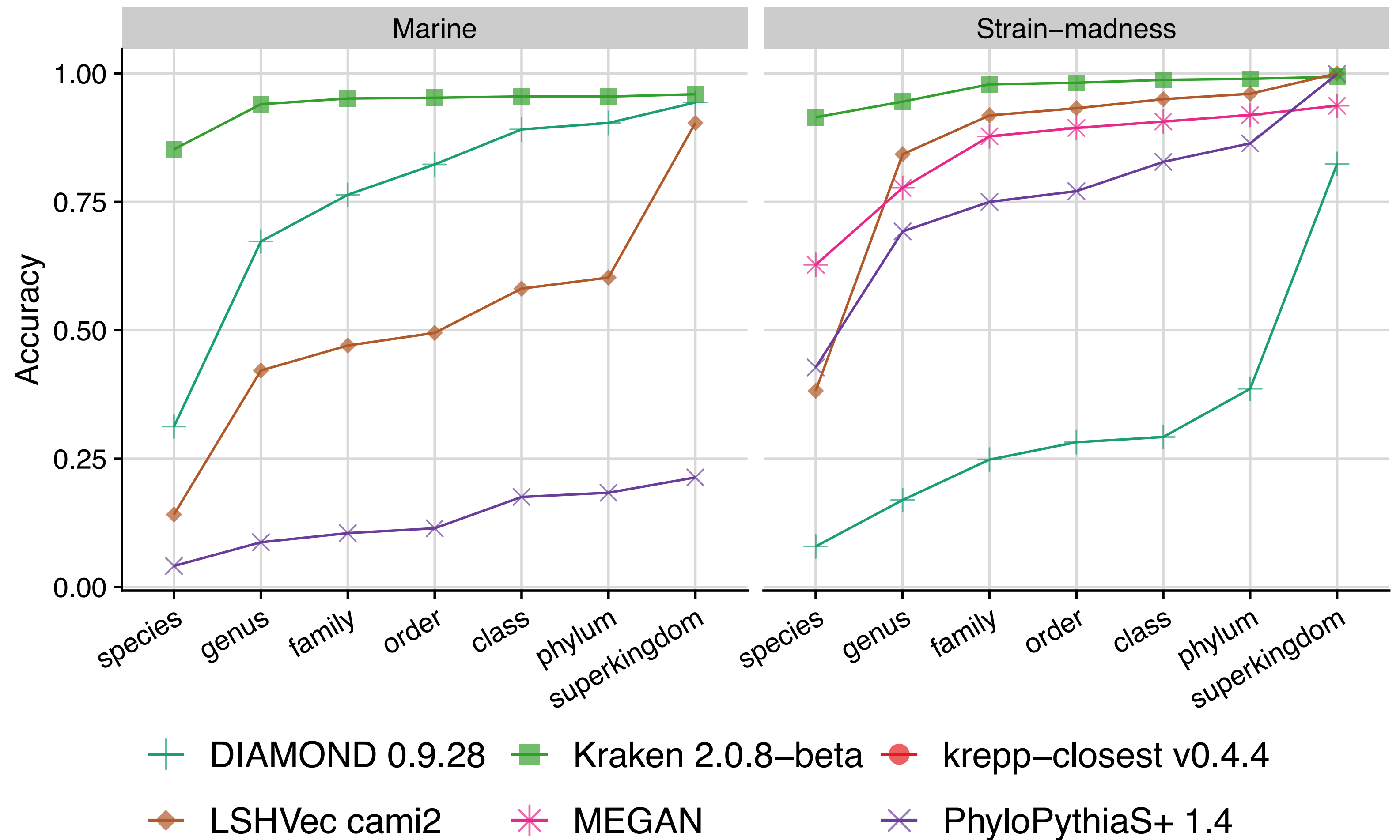
## Taxonomic binning of contigs in CAMI-II

- **CAMI-II:** Binning contigs to taxonomic groups

# What about longer sequences?

## Taxonomic binning of contigs in CAMI-II

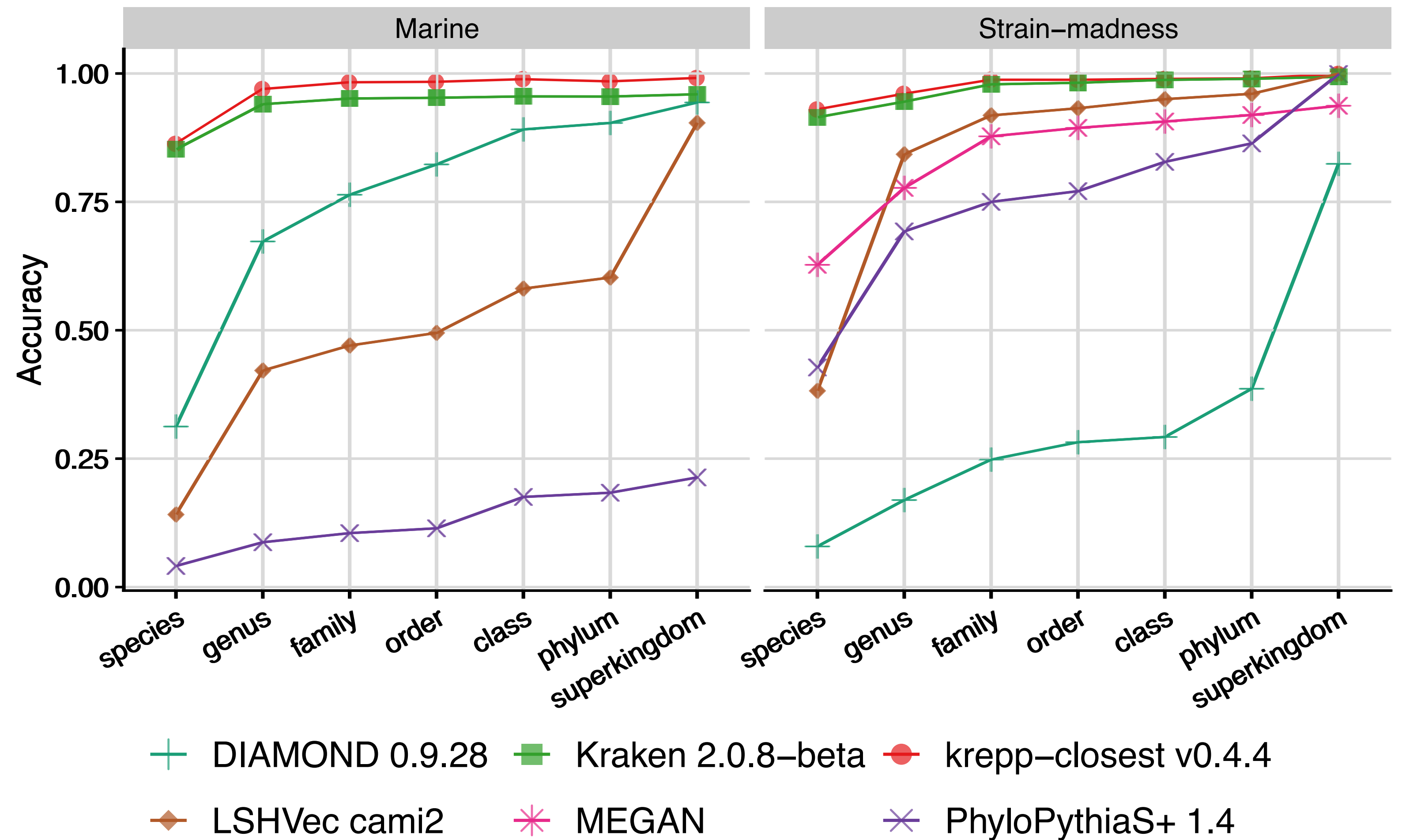
- **CAMI-II:** Binning contigs to taxonomic groups
- Simulated data mimicking two different environments



# What about longer sequences?

## Taxonomic binning of contigs in CAMI-II

- **CAMI-II:** Binning contigs to taxonomic groups
- Simulated data mimicking two different environments
- *krepp* assigns each contig to the **closest genome's group**



# **Dealing with uncertainty: statistically distinguishability**

# Dealing with uncertainty: statistically distinguishability

- short reads — **low signal**
- high distances — fewer matching  $k$ -mers
- small differences may not be statistically meaningful
  - ▶ **test distinguishability**


# Dealing with uncertainty: statistically distinguishability


- short reads — **low signal**
- high distances — fewer matching  $k$ -mers
- small differences may not be statistically meaningful
  - ▶ **test distinguishability**

## likelihood-ratio test

with the closest reference:

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

  $D$ : alternative distance

  $i^*$ : closest reference

# Dealing with uncertainty: statistically distinguishability

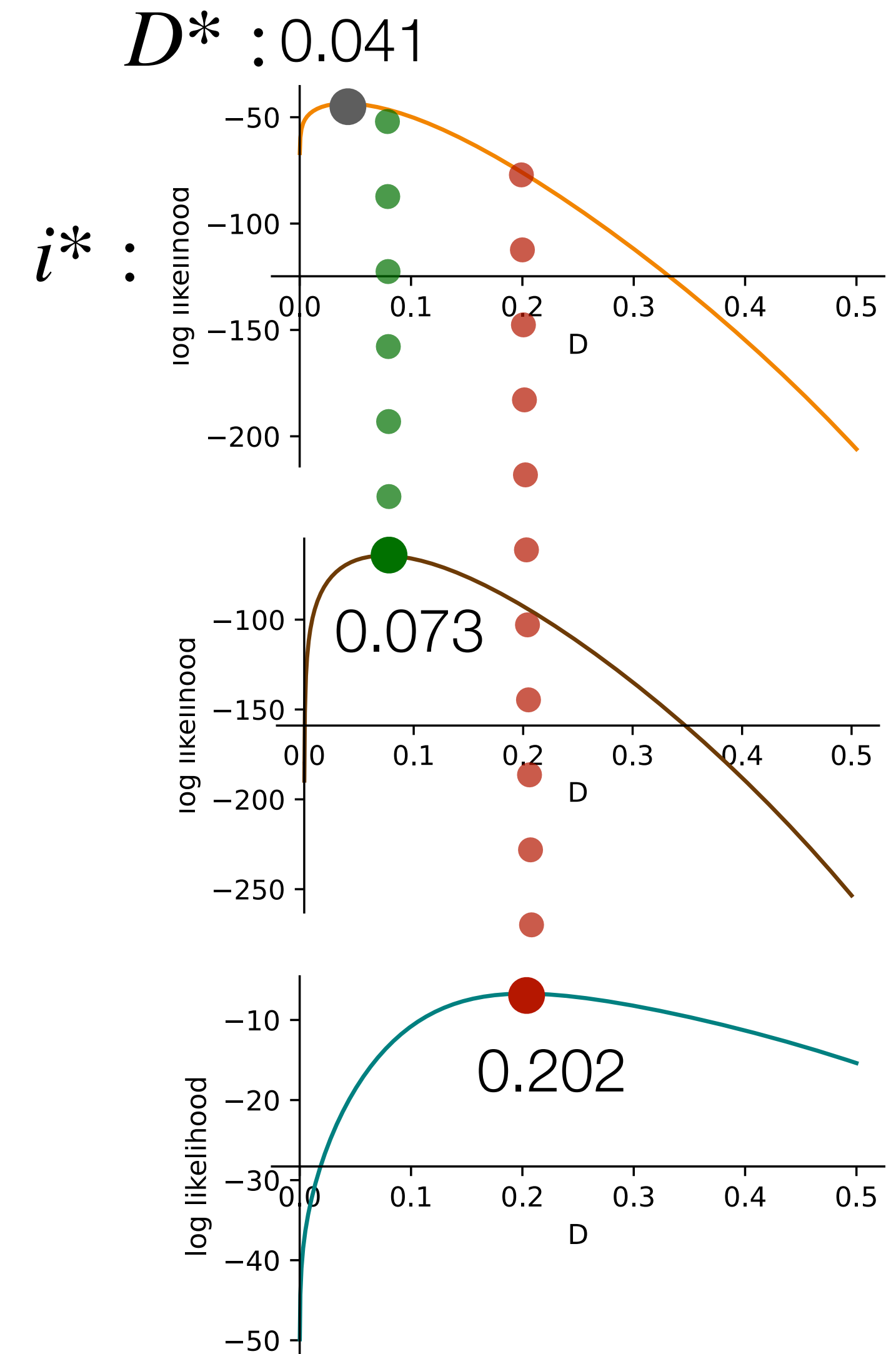
- short reads — **low signal**
- high distances — fewer matching  $k$ -mers
- small differences may not be statistically meaningful
  - ▶ **test distinguishability**

**likelihood-ratio test**  
with the closest reference:

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

↗  $D$ : alternative distance

↘  $i^*$ : closest reference



# Dealing with uncertainty: statistically distinguishability

- short reads — **low signal**
- high distances — fewer matching  $k$ -mers
- small differences may not be statistically meaningful
  - ▶ **test distinguishability**

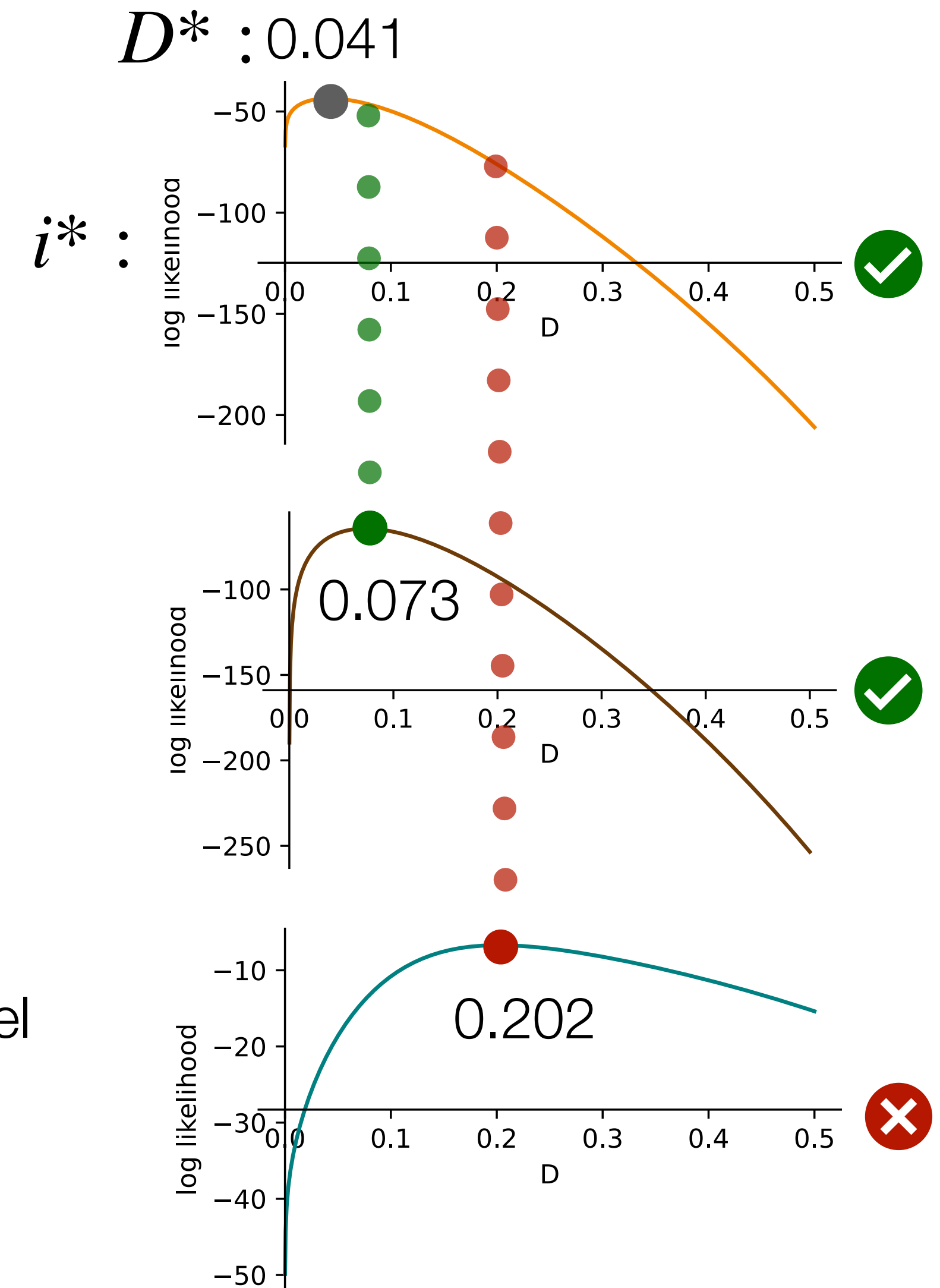
**likelihood-ratio test**  
with the closest reference:

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

↗  $D$ : alternative distance  
↘  $i^*$ : closest reference

$$\lambda_{LR} \sim \chi^2$$

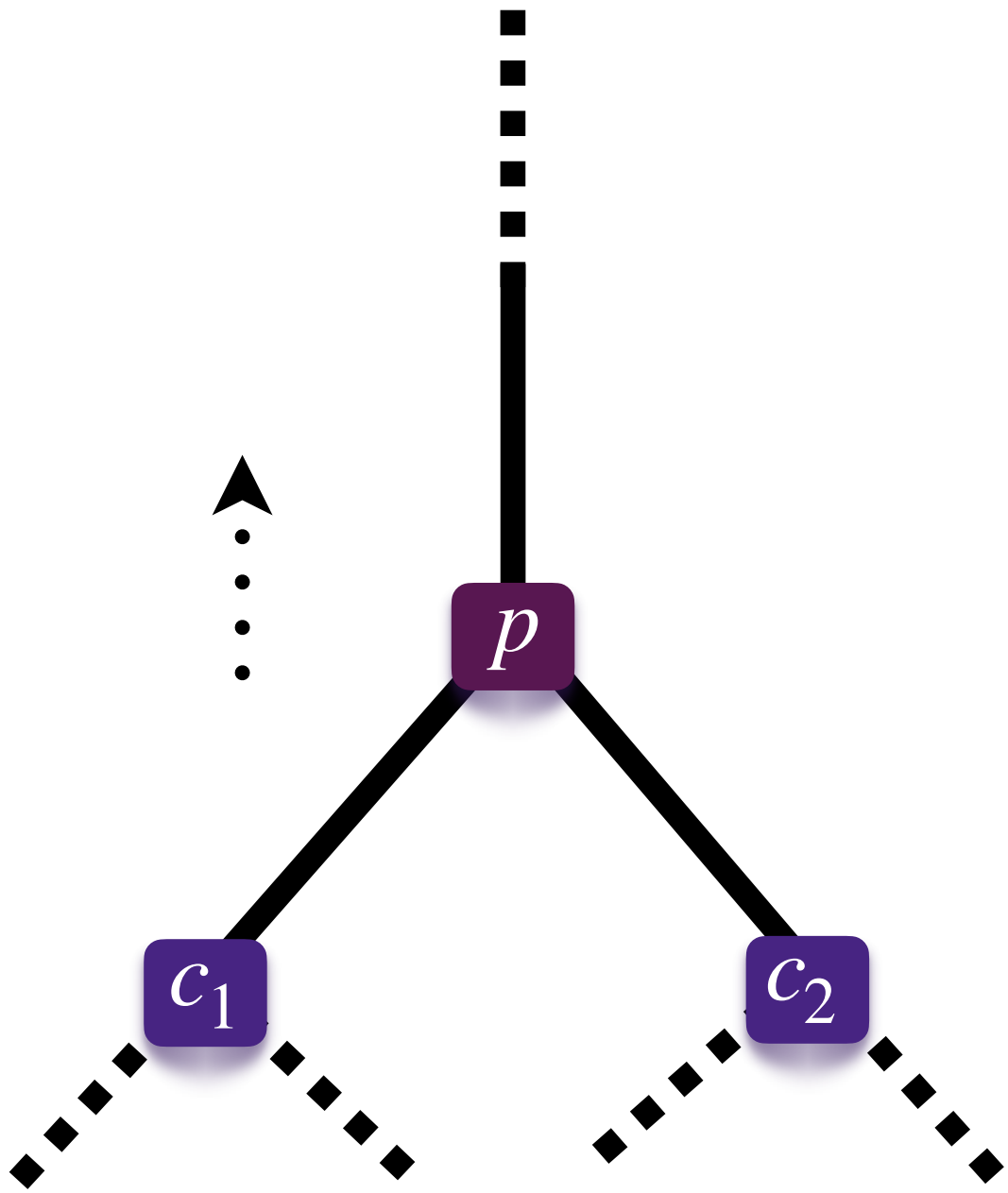
- ▶ select a significance level  
(default:  $\alpha=90\%$ )



# **Defining clade distances & branch length agnostic placement**

# Defining clade distances & branch length agnostic placement

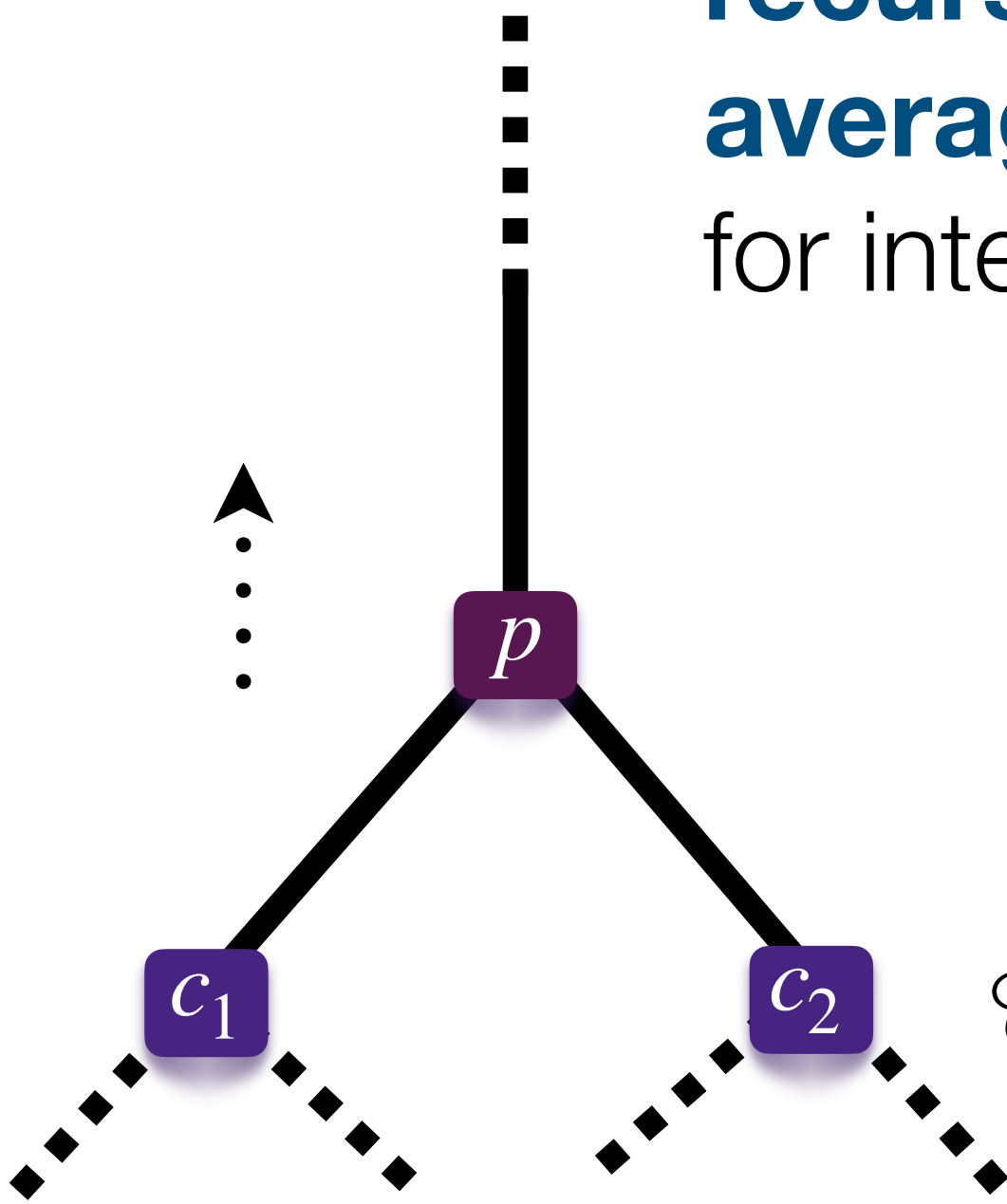
A notion of distance  
for internal nodes:



# Defining clade distances & branch length agnostic placement

A notion of distance  
for internal nodes:

recursively compute  
average HD histograms  
for internal nodes



$$\mathbf{v}_p = \frac{\sum_{c \in \mathcal{C}(p)} \mathbf{v}_c}{|\mathcal{C}(p)|}$$

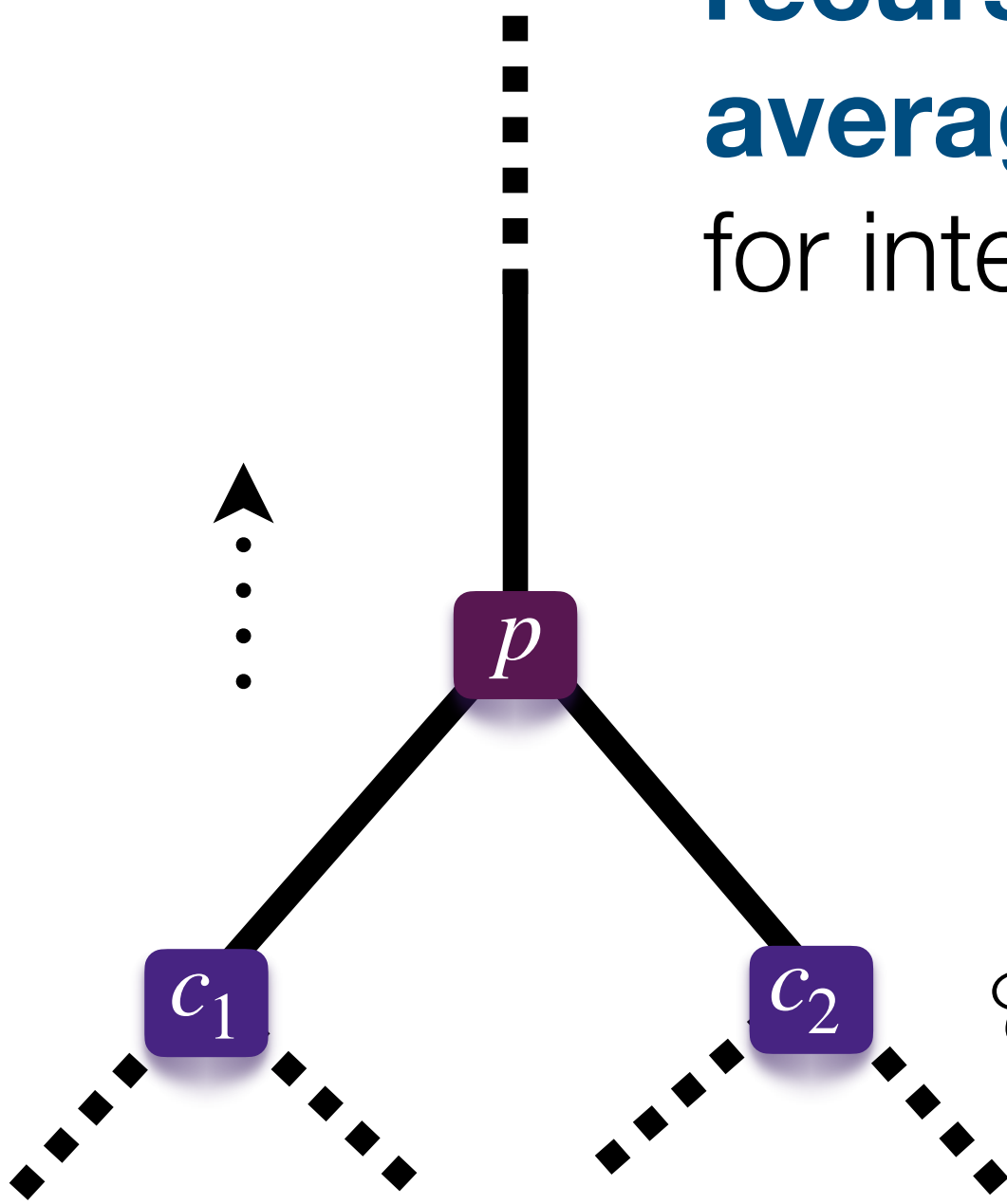
$\mathcal{C}(p)$ : set of children of  $p$ ,  $\{c_1, c_2\}$

# Defining clade distances & branch length agnostic placement

A notion of distance  
for internal nodes:

**recursively compute  
average** HD histograms  
for internal nodes

use the same likelihood model  
and log-likelihood ratio test



$$\mathbf{v}_p = \frac{\sum_{c \in \mathcal{C}(p)} \mathbf{v}_c}{|\mathcal{C}(p)|}$$

$\mathcal{C}(p)$ : set of children of  $p$ ,  $\{c_1, c_2\}$

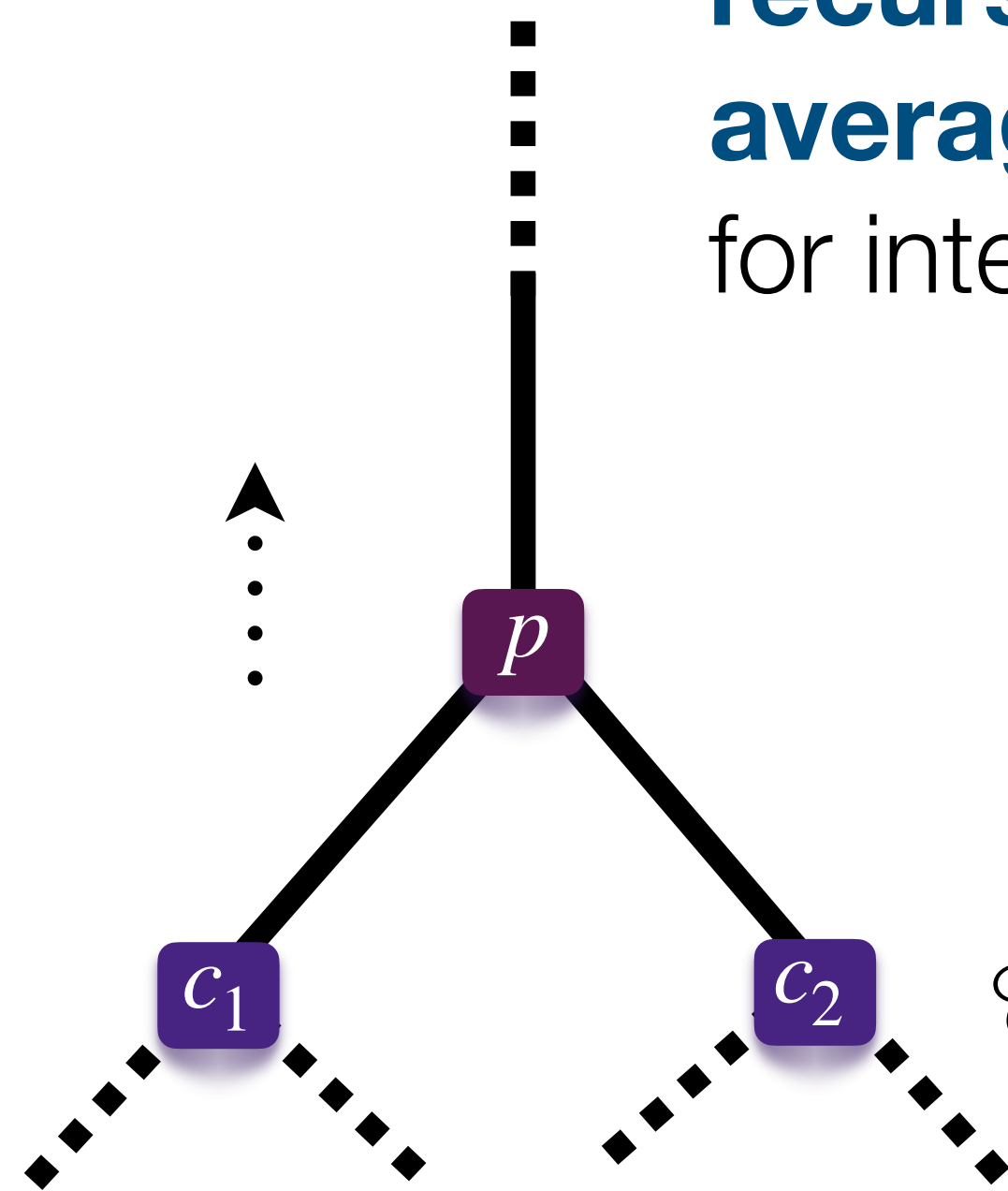
# Defining clade distances & branch length agnostic placement

A notion of distance  
for internal nodes:

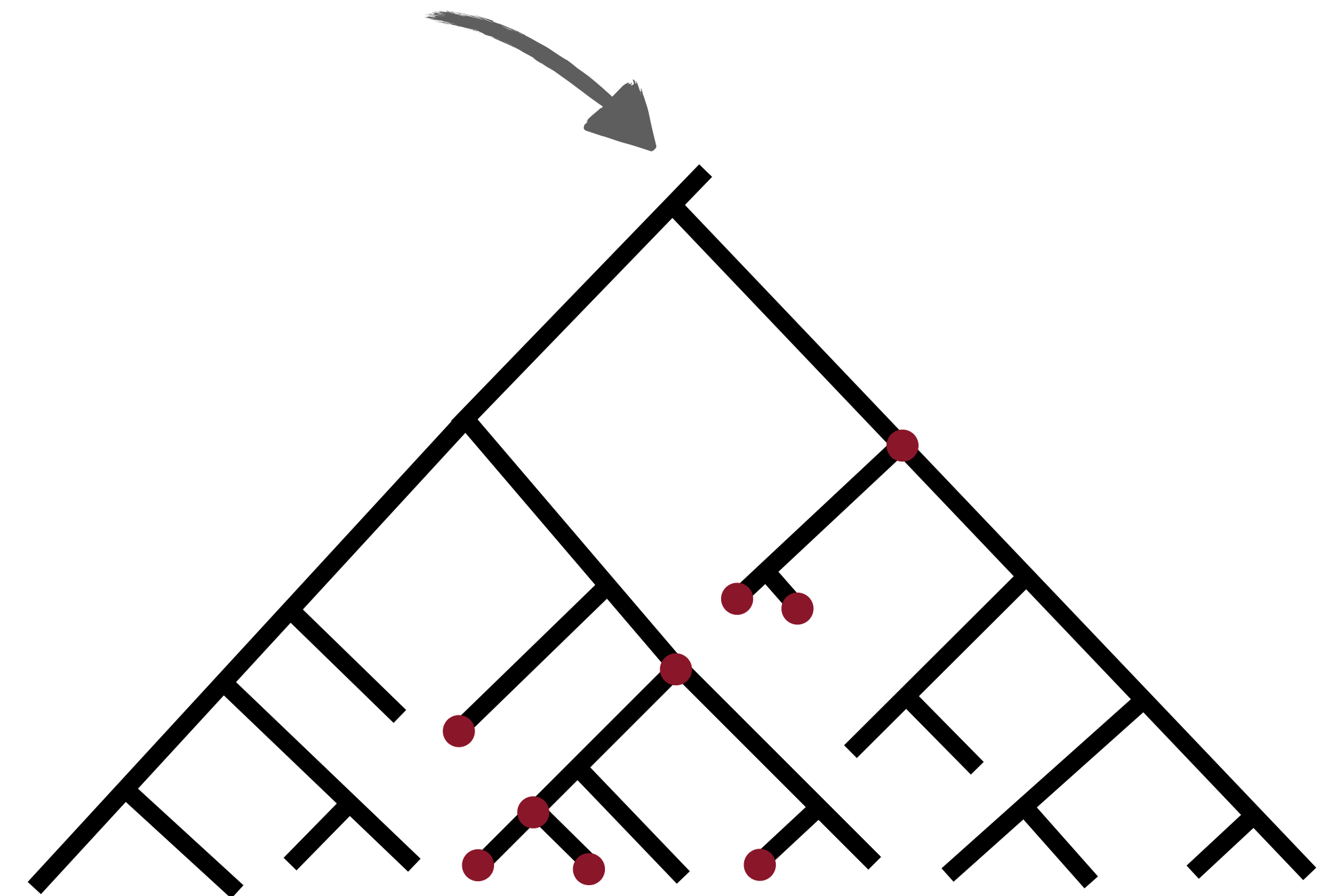
**recursively compute  
average** HD histograms  
for internal nodes

$$\mathbf{v}_p = \frac{\sum_{c \in \mathcal{C}(p)} \mathbf{v}_c}{|\mathcal{C}(p)|}$$

$\mathcal{C}(p)$ : set of children of  $p$ ,  $\{c_1, c_2\}$



use the same likelihood model  
and log-likelihood ratio test



● : indistinguishable w.r.t.  
the closest reference

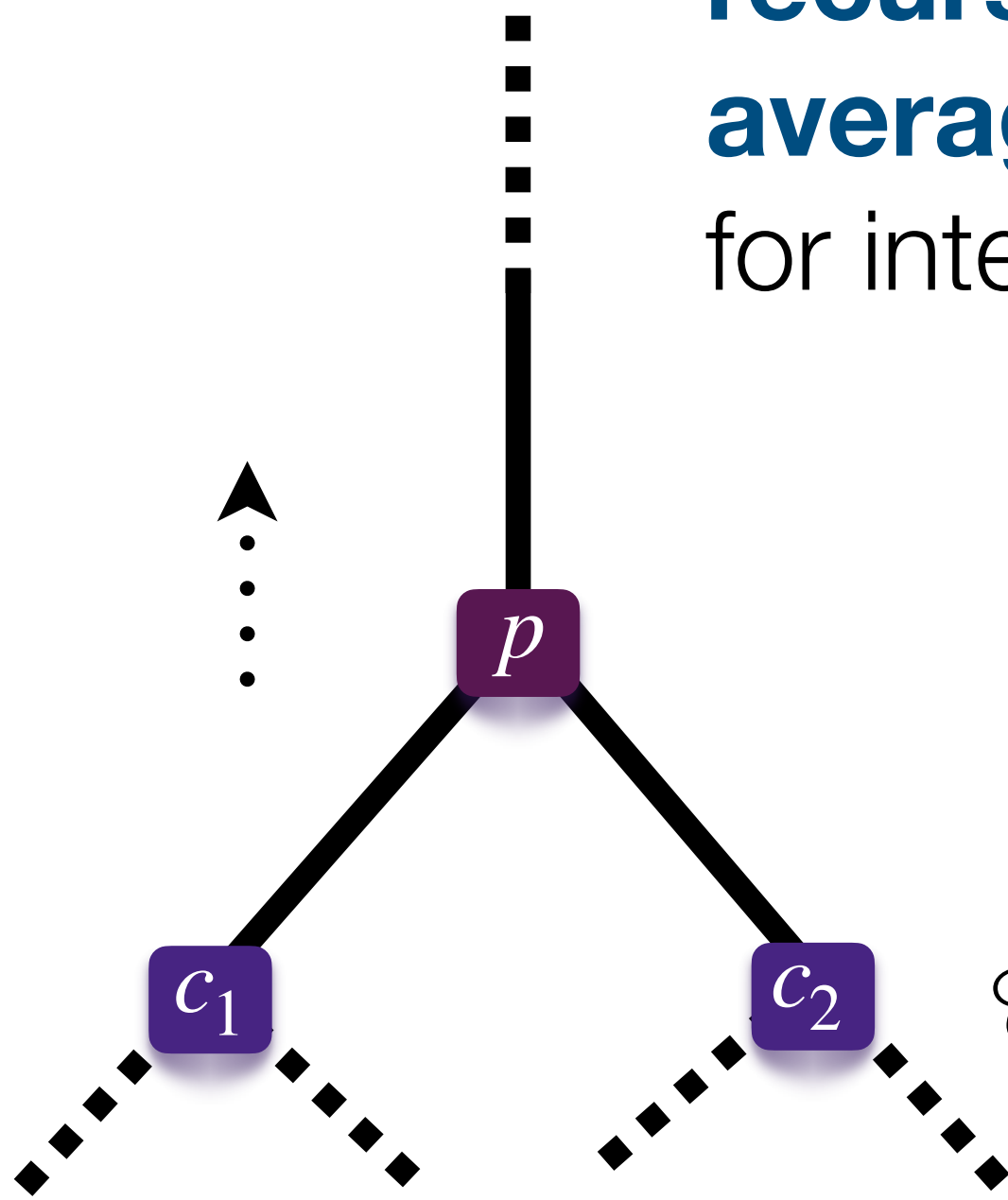
# Defining clade distances & branch length agnostic placement

A notion of distance  
for internal nodes:

**recursively compute  
average** HD histograms  
for internal nodes

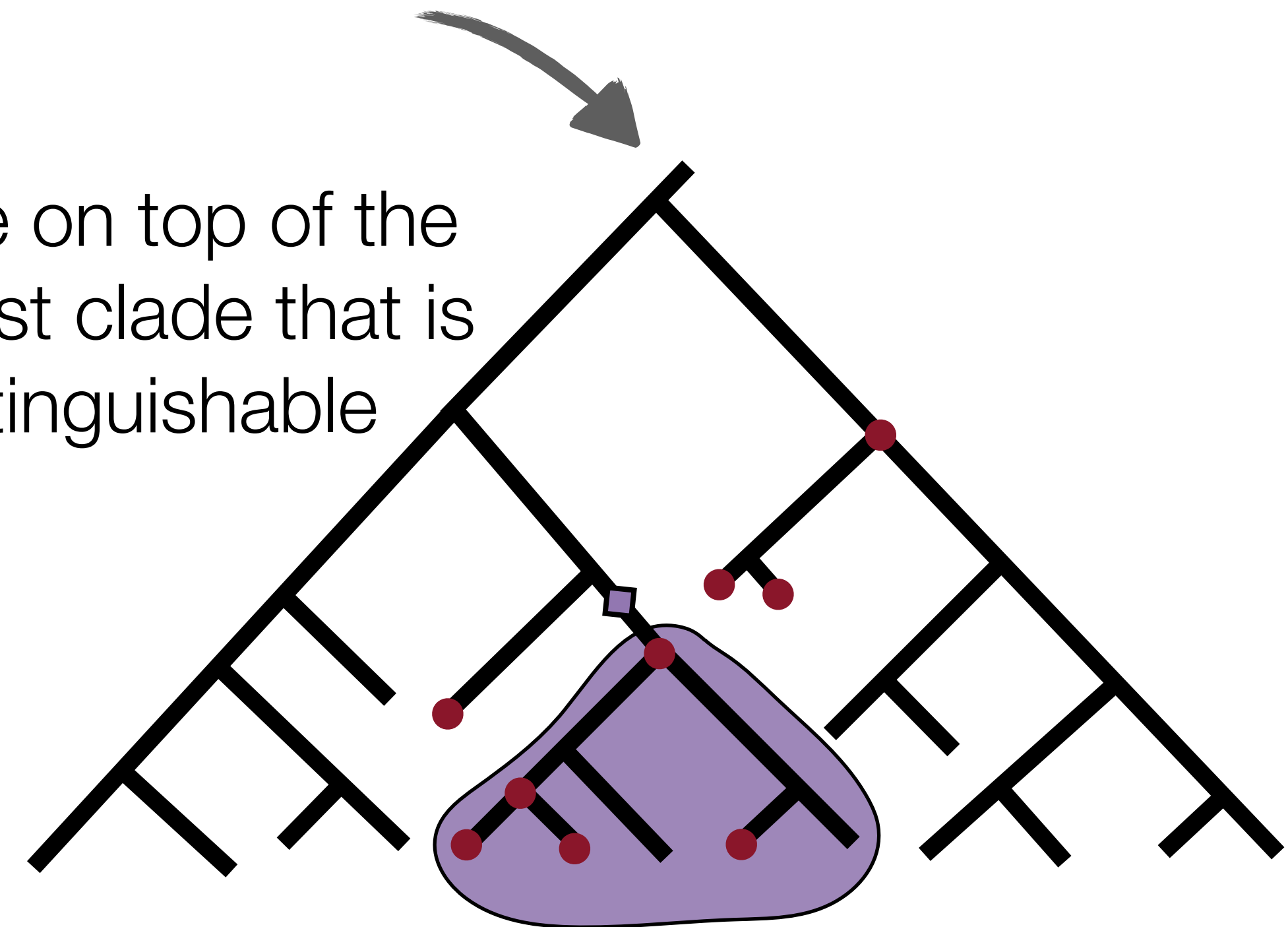
$$\mathbf{v}_p = \frac{\sum_{c \in \mathcal{C}(p)} \mathbf{v}_c}{|\mathcal{C}(p)|}$$

$\mathcal{C}(p)$ : set of children of  $p$ ,  $\{c_1, c_2\}$



use the same likelihood model  
and log-likelihood ratio test

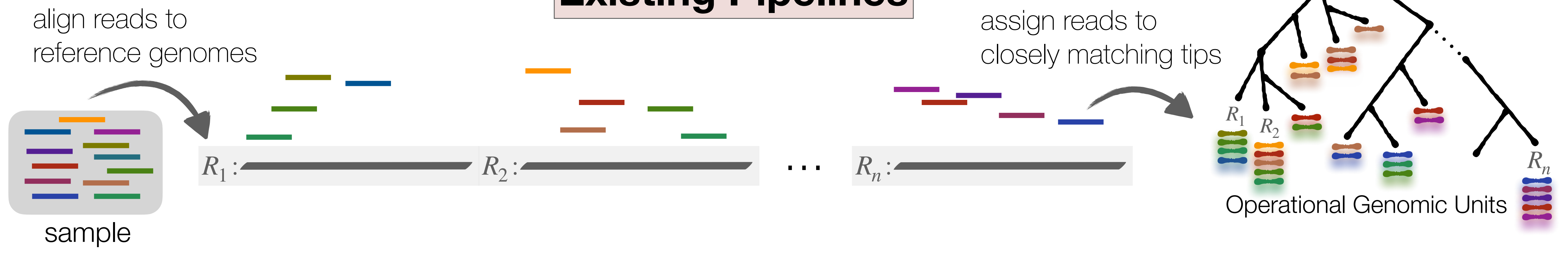
place on top of the  
largest clade that is  
indistinguishable



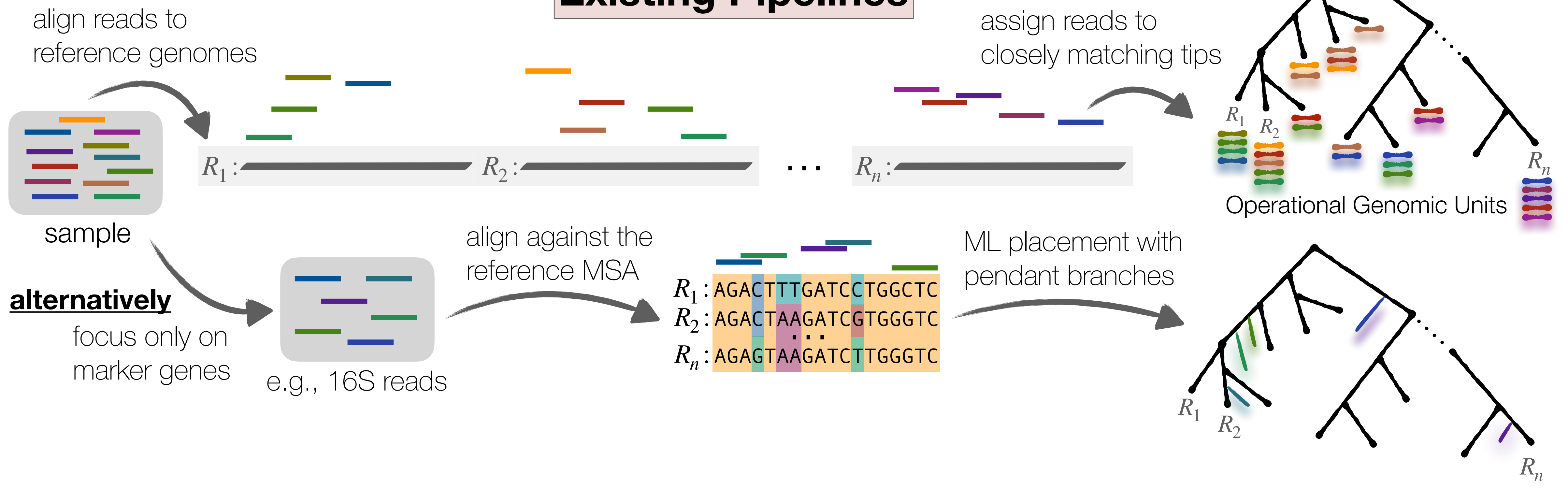
● : indistinguishable w.r.t.  
the closest reference

# Existing Pipelines

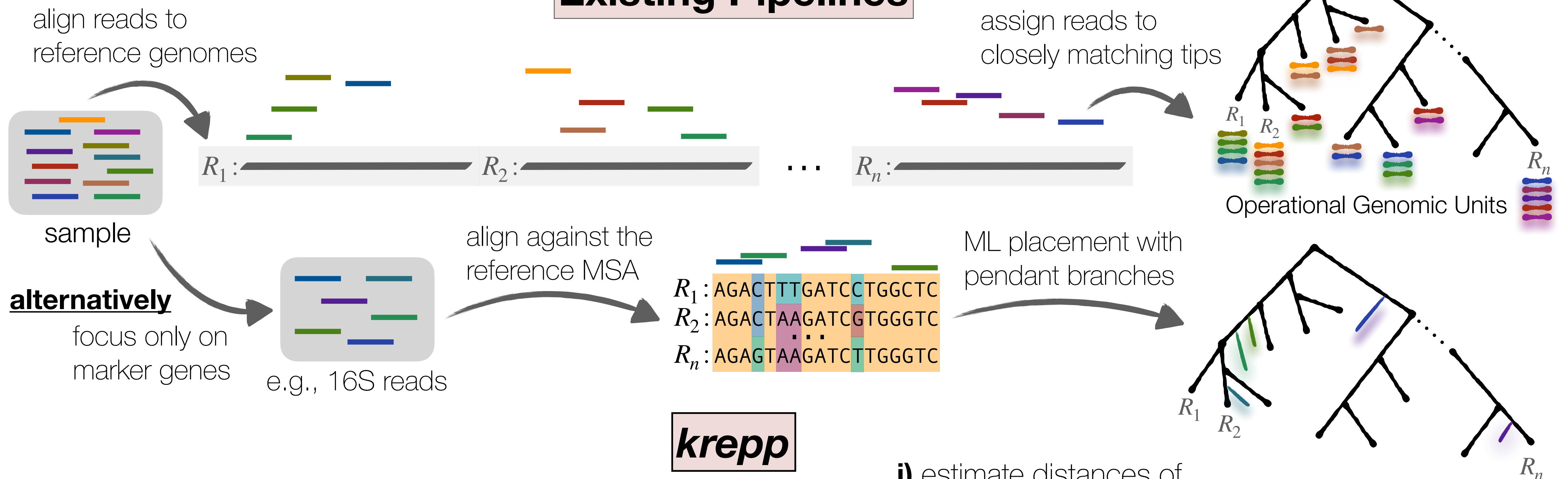
# Existing Pipelines



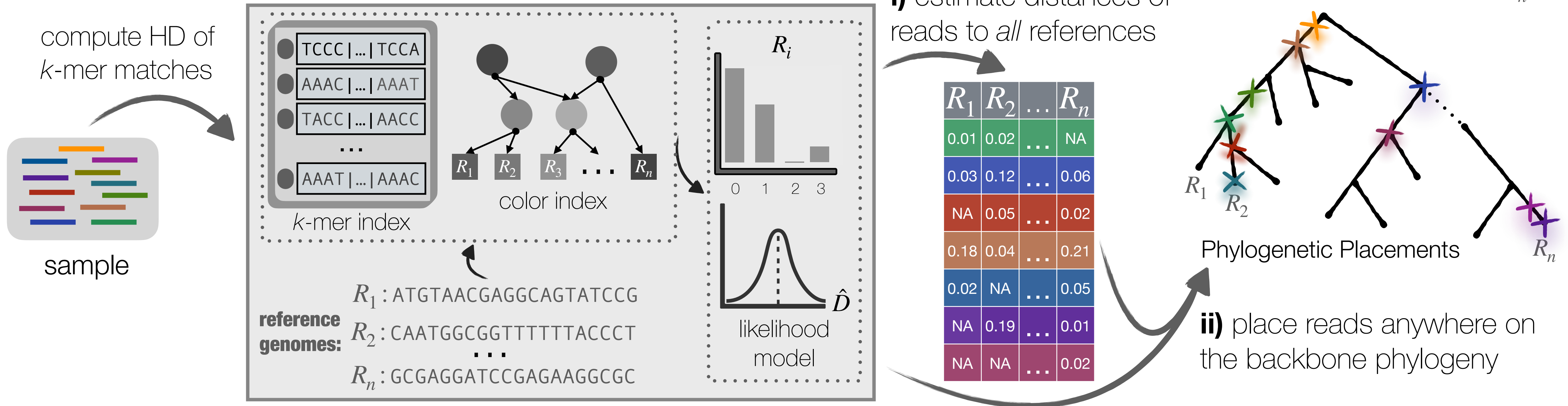
# Existing Pipelines



# Existing Pipelines



# krepp



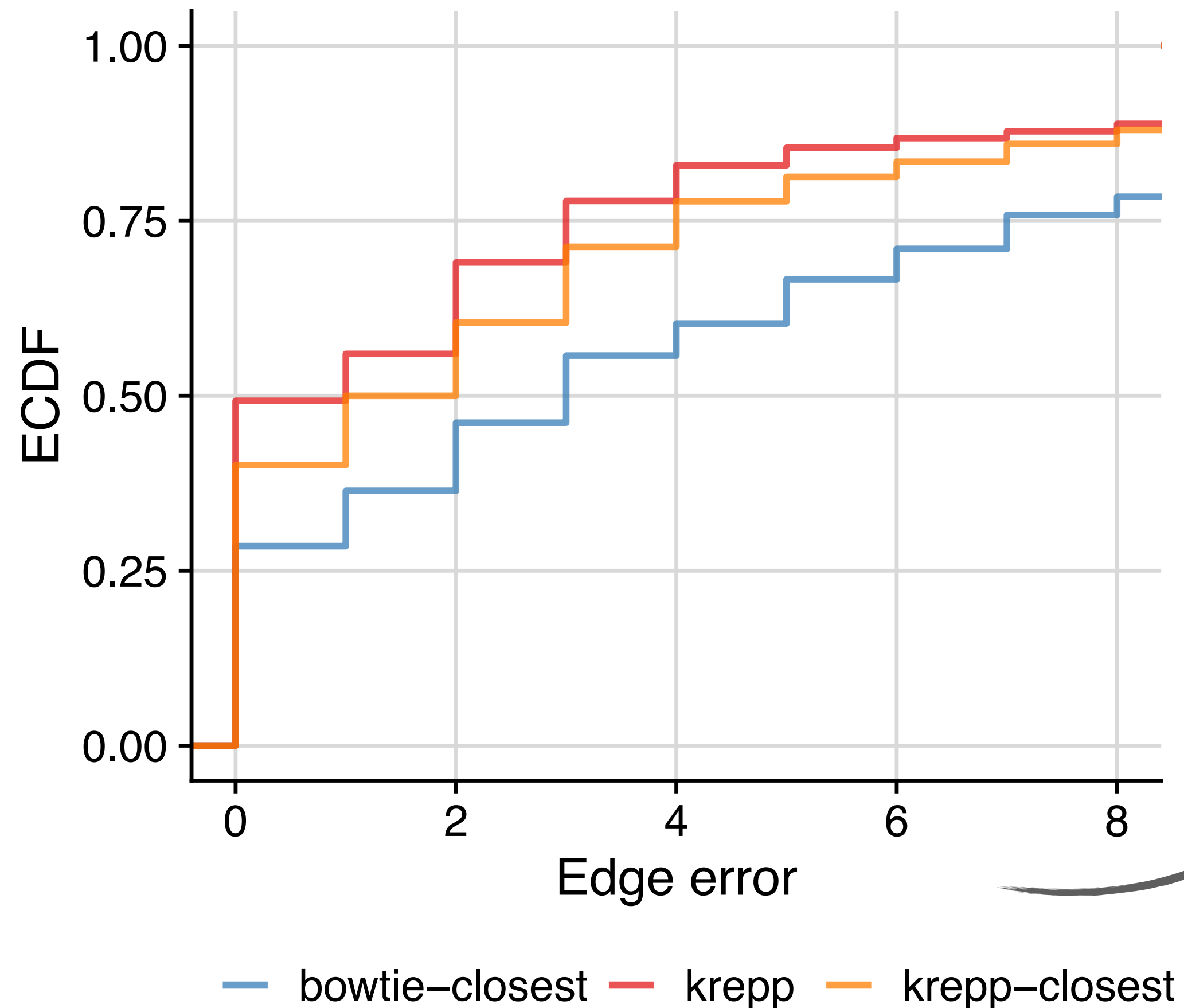
# **krepp's heuristic improves closest-tip placement**

# **krepp's heuristic improves closest-tip placement**

- Leave all out: 100 diverse taxa
- WoL-v2 tree with 15,952 leaves

# krepp's heuristic improves closest-tip placement

- Leave all out: 100 diverse taxa
- WoL-v2 tree with 15,952 leaves
- Outperforms baselines:  
on the closest, on the LCA, etc.
- >80% of all reads within four  
edges of the correct node



how many edges  
away is the  
placement from  
the correct edge?

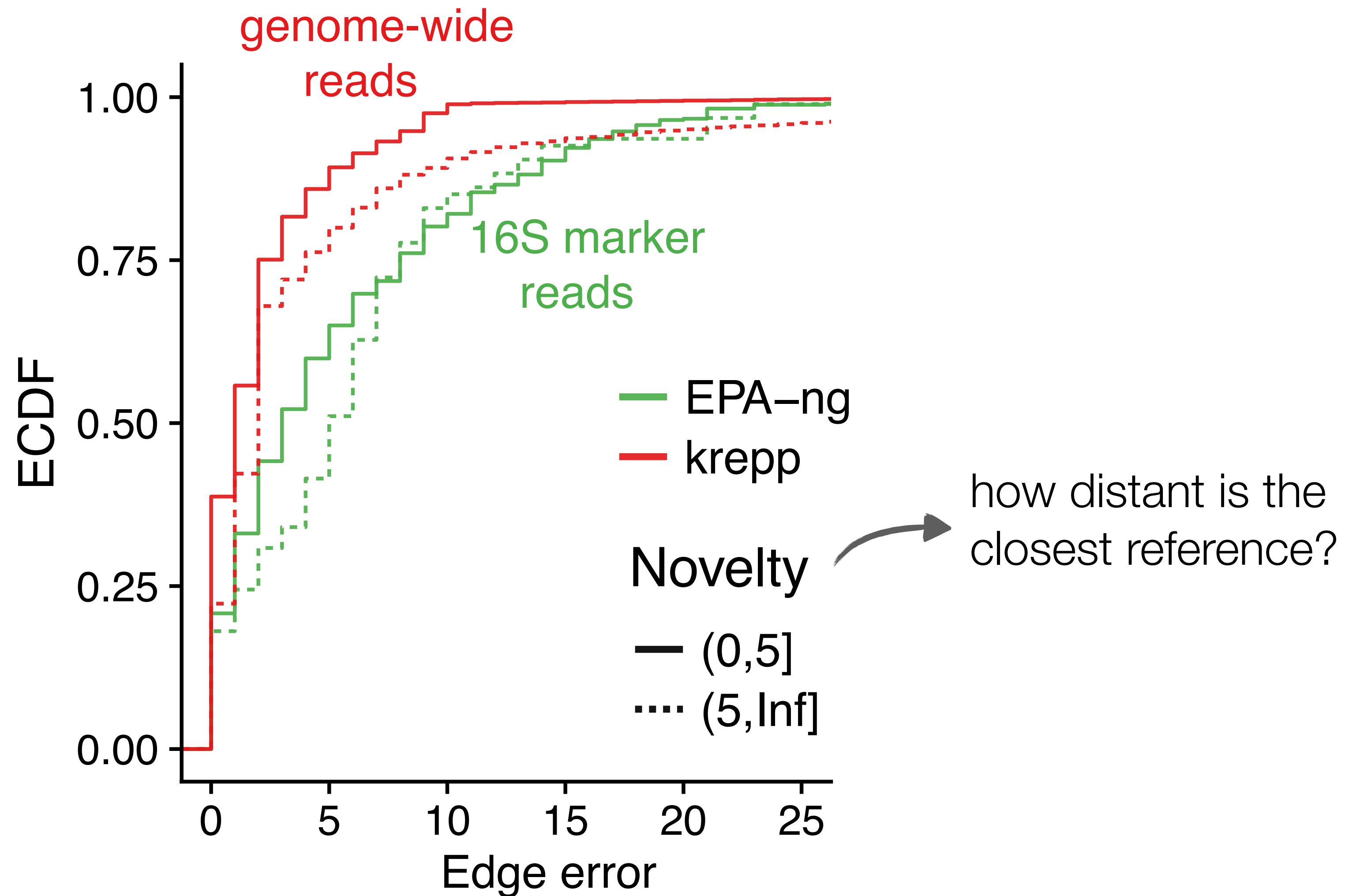
**krepp places genome-wide reads more accurately than  
ML-based placement of 16S reads**

# **krepp places genome-wide reads more accurately than ML-based placement of 16S reads**

- Leave all out — 100 diverse taxa
- WoL-v1 tree — 10,575 leaves

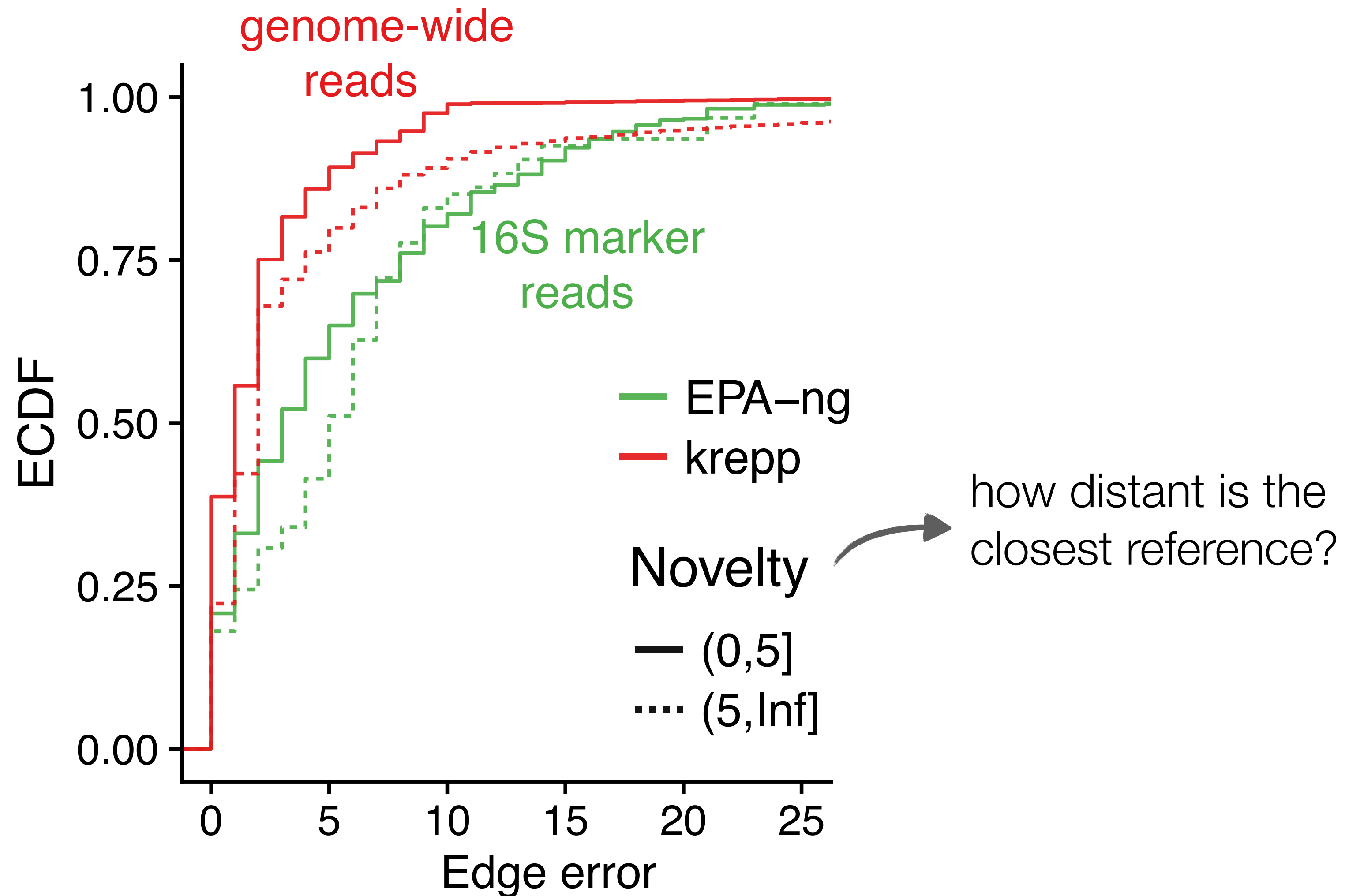
# krepp places genome-wide reads more accurately than ML-based placement of 16S reads

- Leave all out — 100 diverse taxa
- WoL-v1 tree — 10,575 leaves



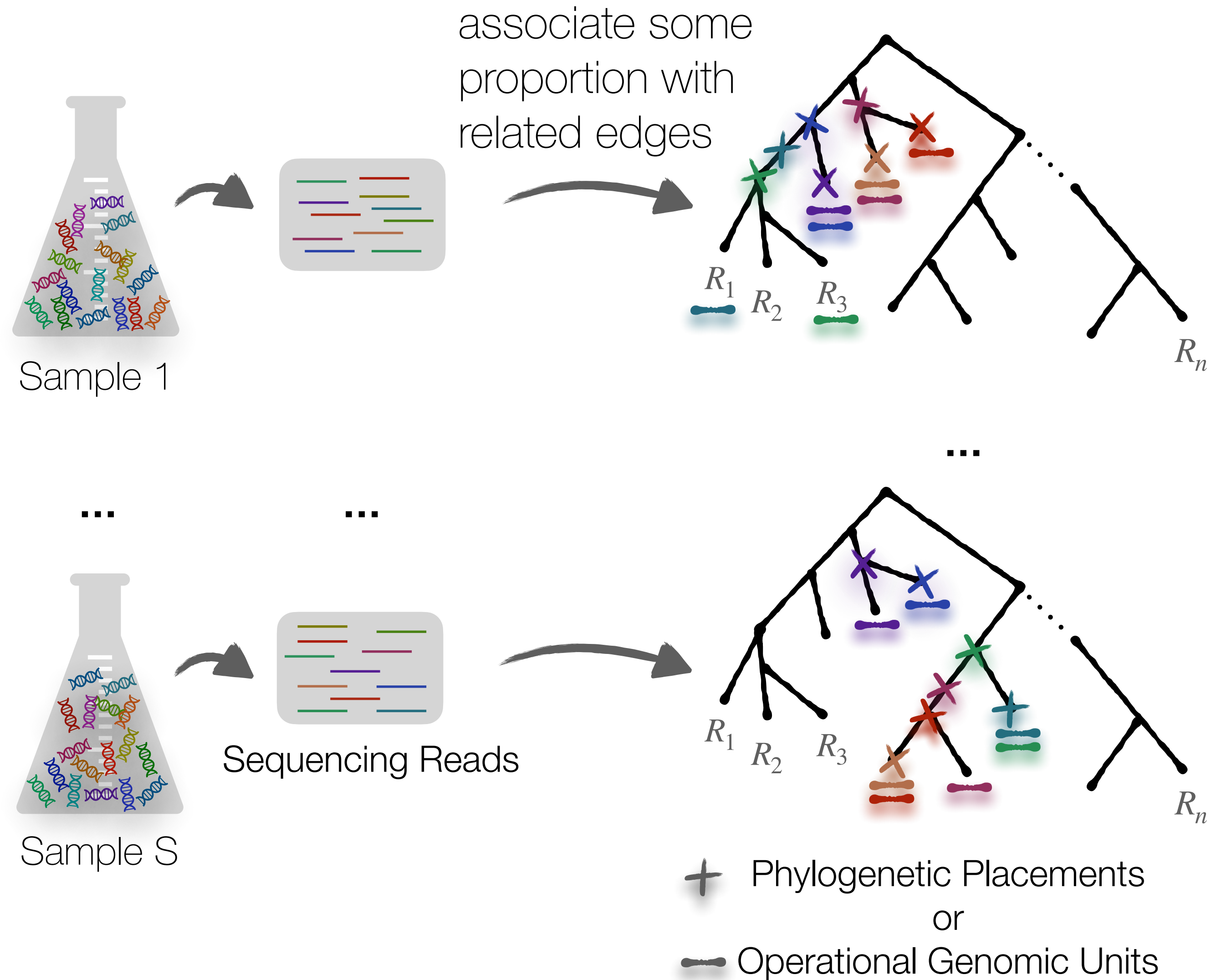
# krepp places genome-wide reads more accurately than ML-based placement of 16S reads

- Leave all out — 100 diverse taxa
- WoL-v1 tree — 10,575 leaves
- 2.4 versus 5.6 edges (average)
- *krepp* places 86% of all reads

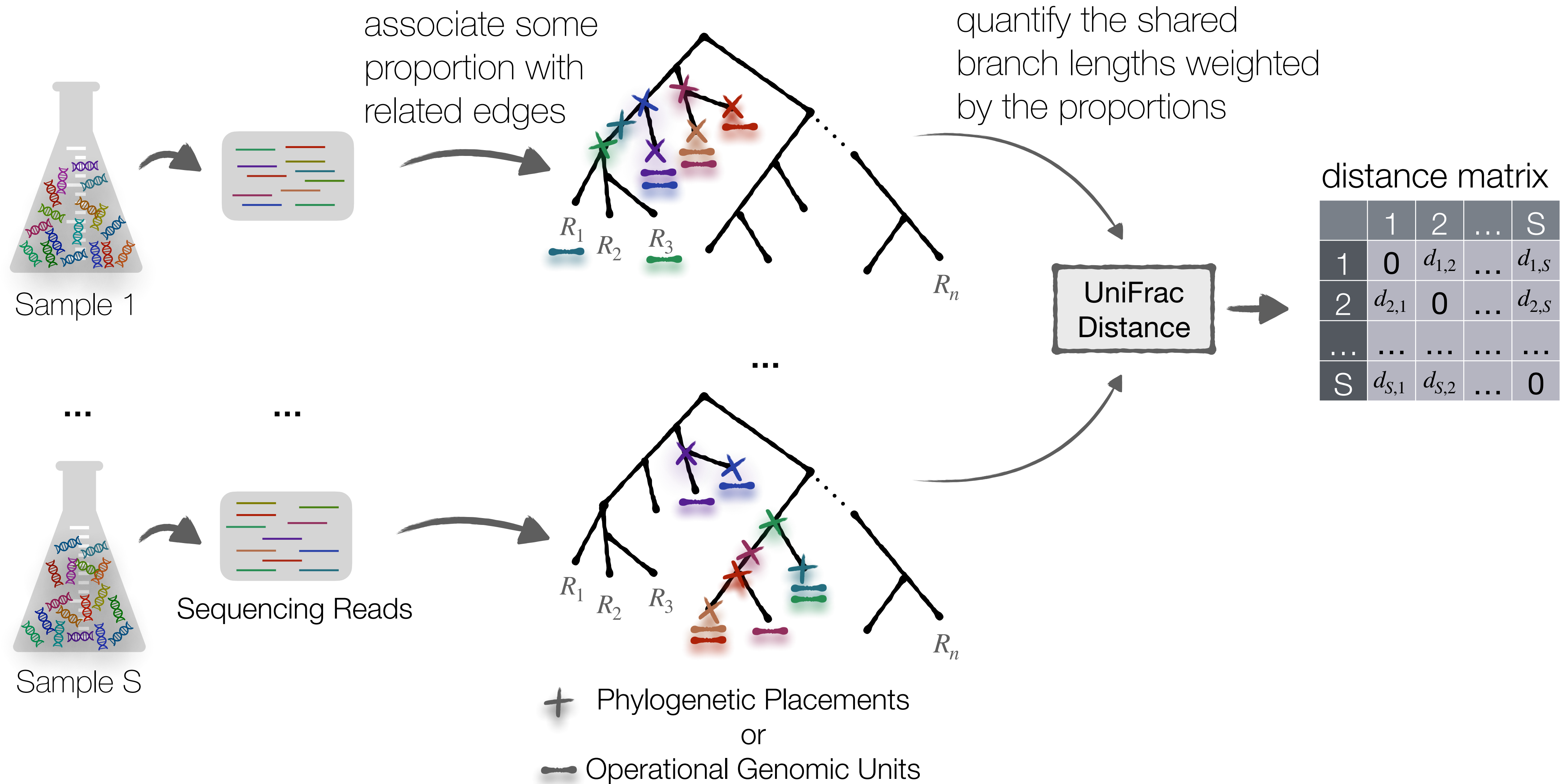


# **Characterization of metagenomic samples on a backbone phylogeny**

# Characterization of metagenomic samples on a backbone phylogeny

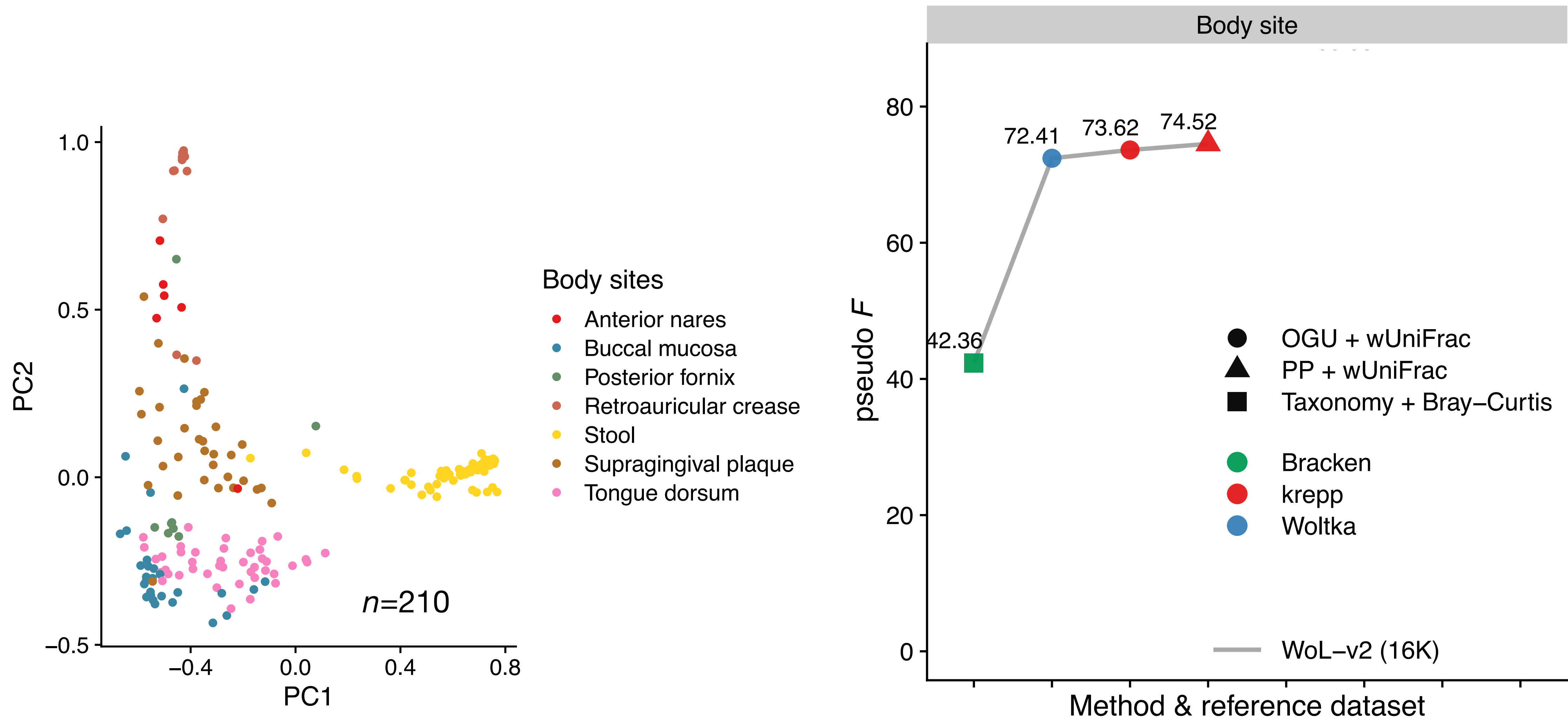


# Characterization of metagenomic samples on a backbone phylogeny



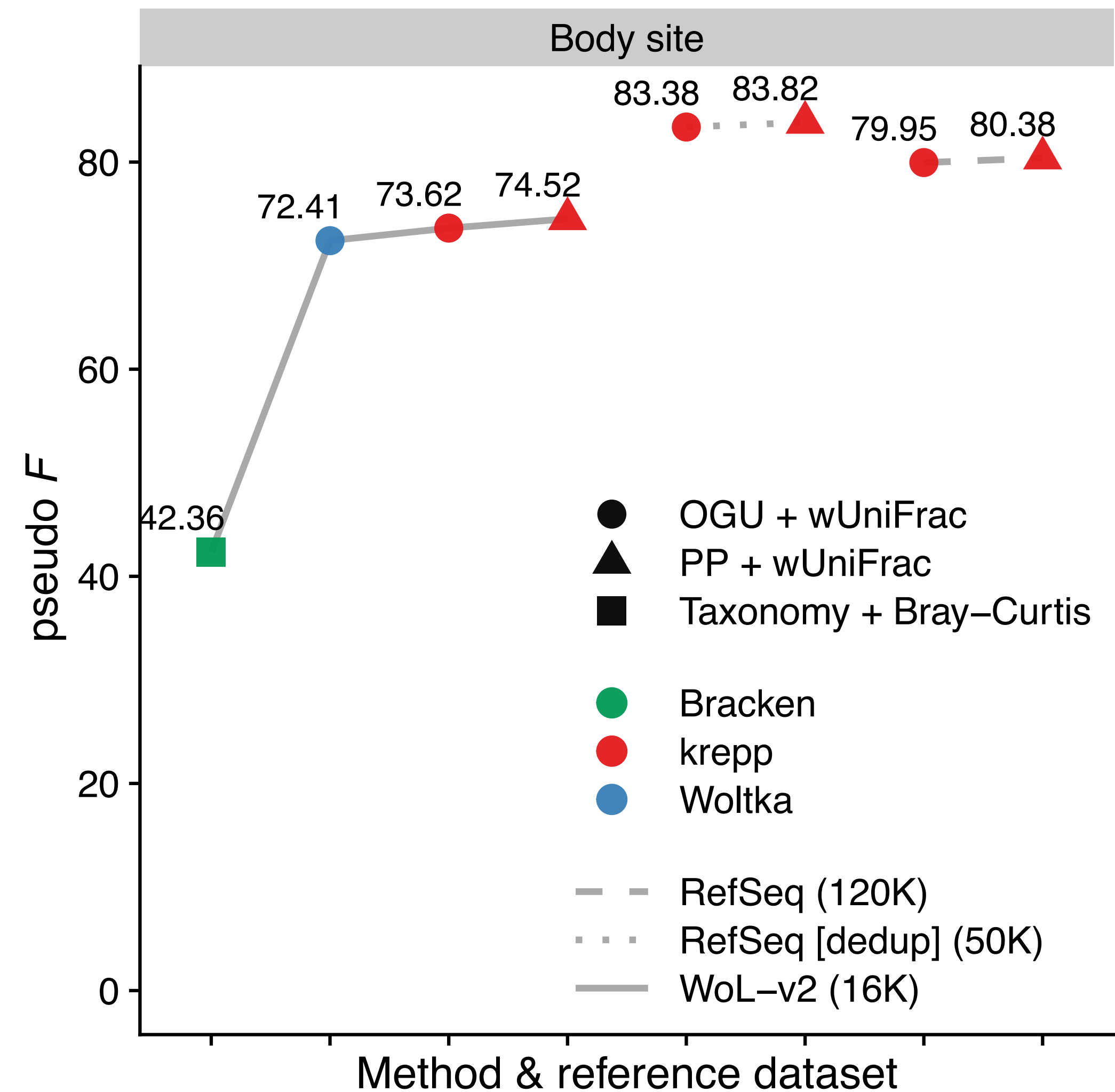
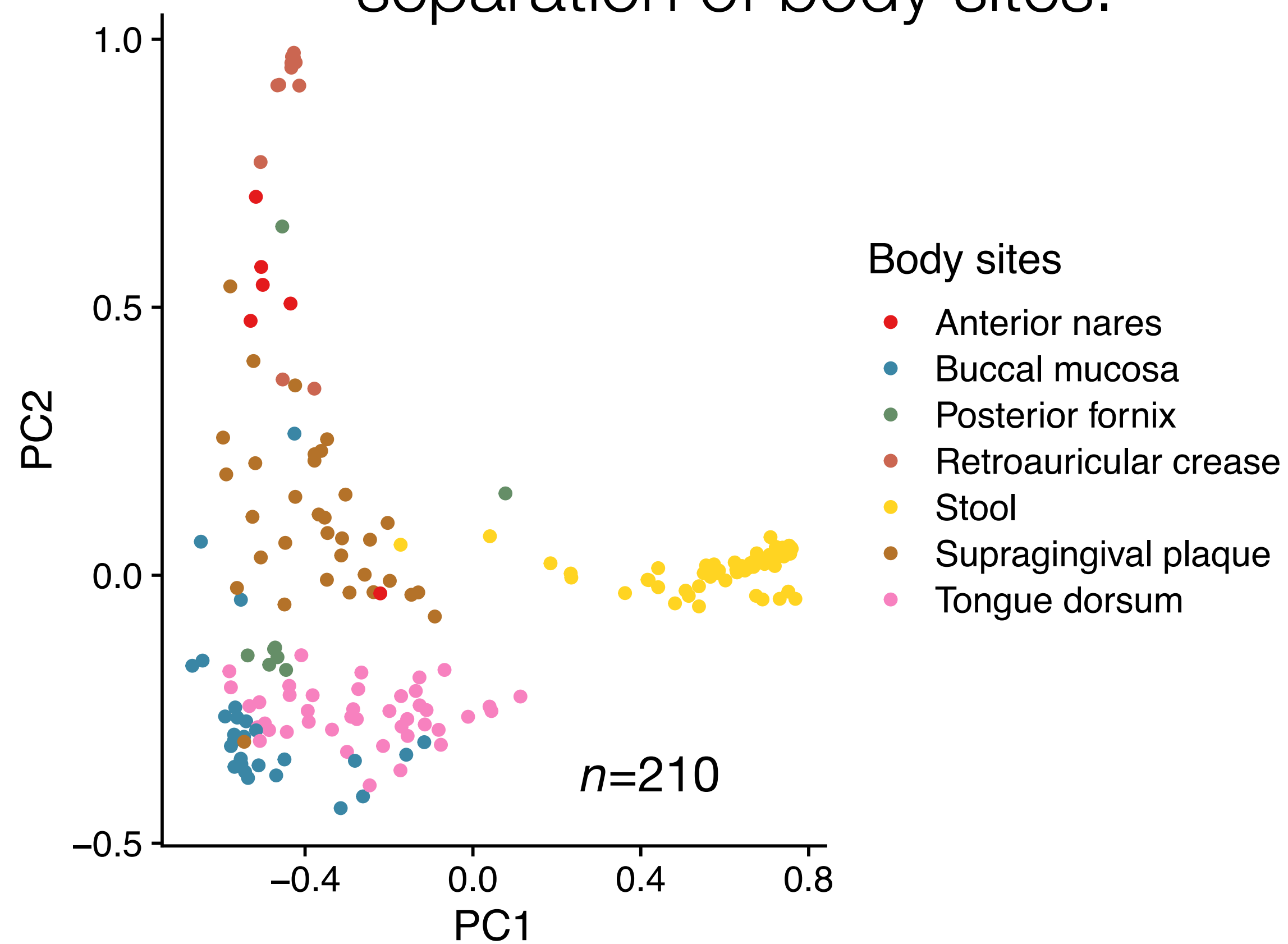


# Analyzing human microbiome with larger references



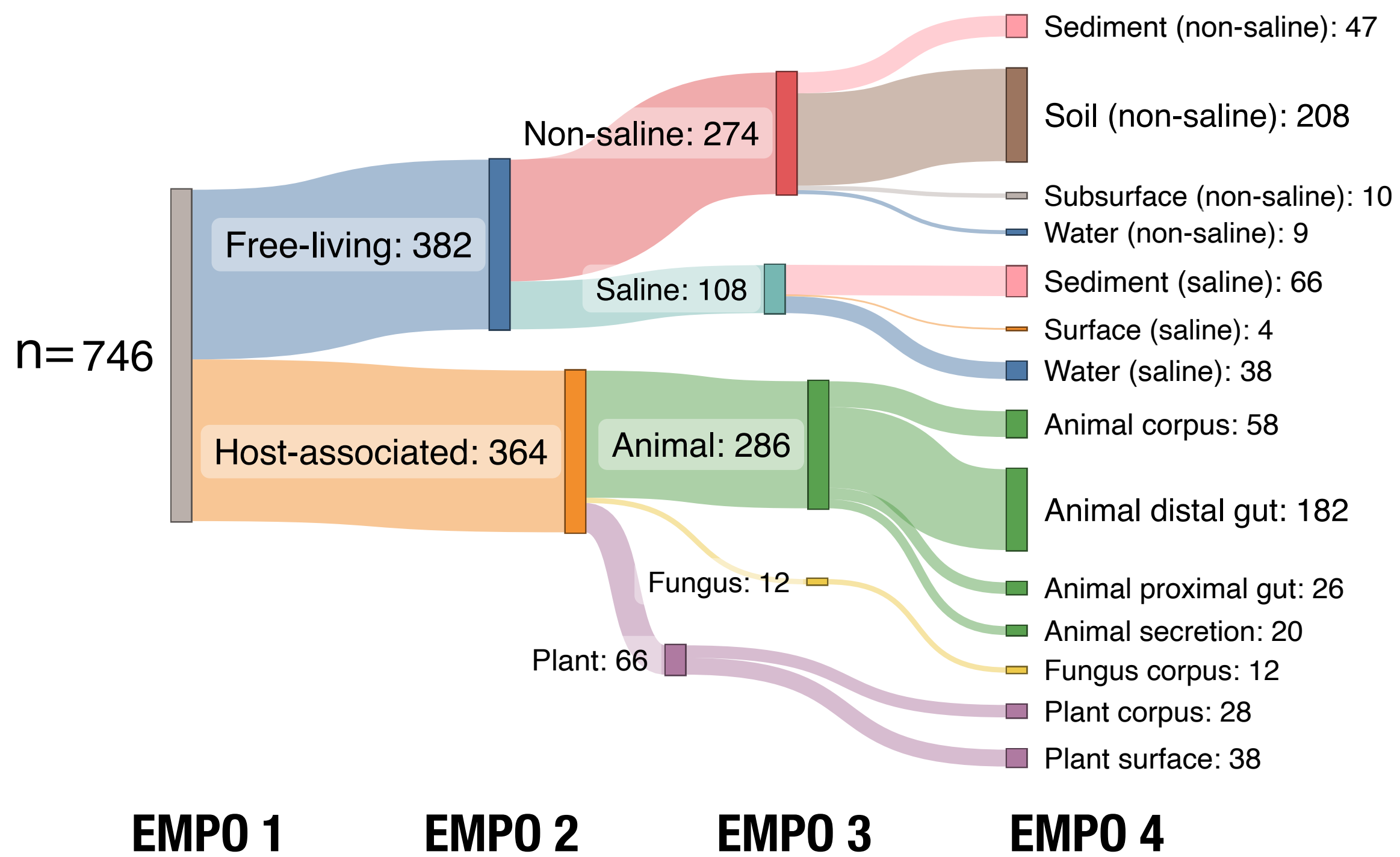
# Analyzing human microbiome with larger references

**Scaling** to large references further **improves** separation of body sites.

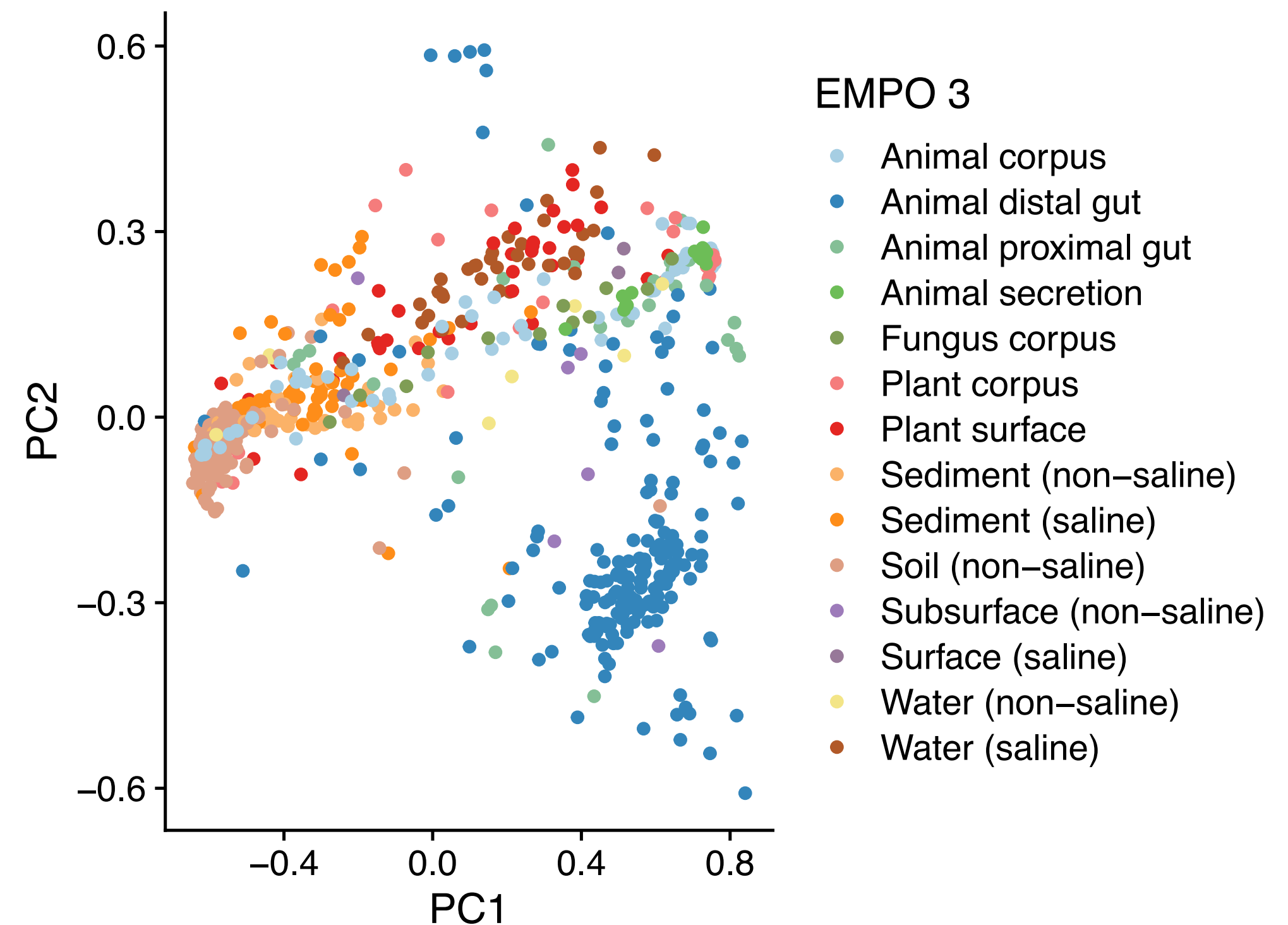


# Improved characterization of earth's microbiome

## Hierarchical categorization of earth microbiome samples:

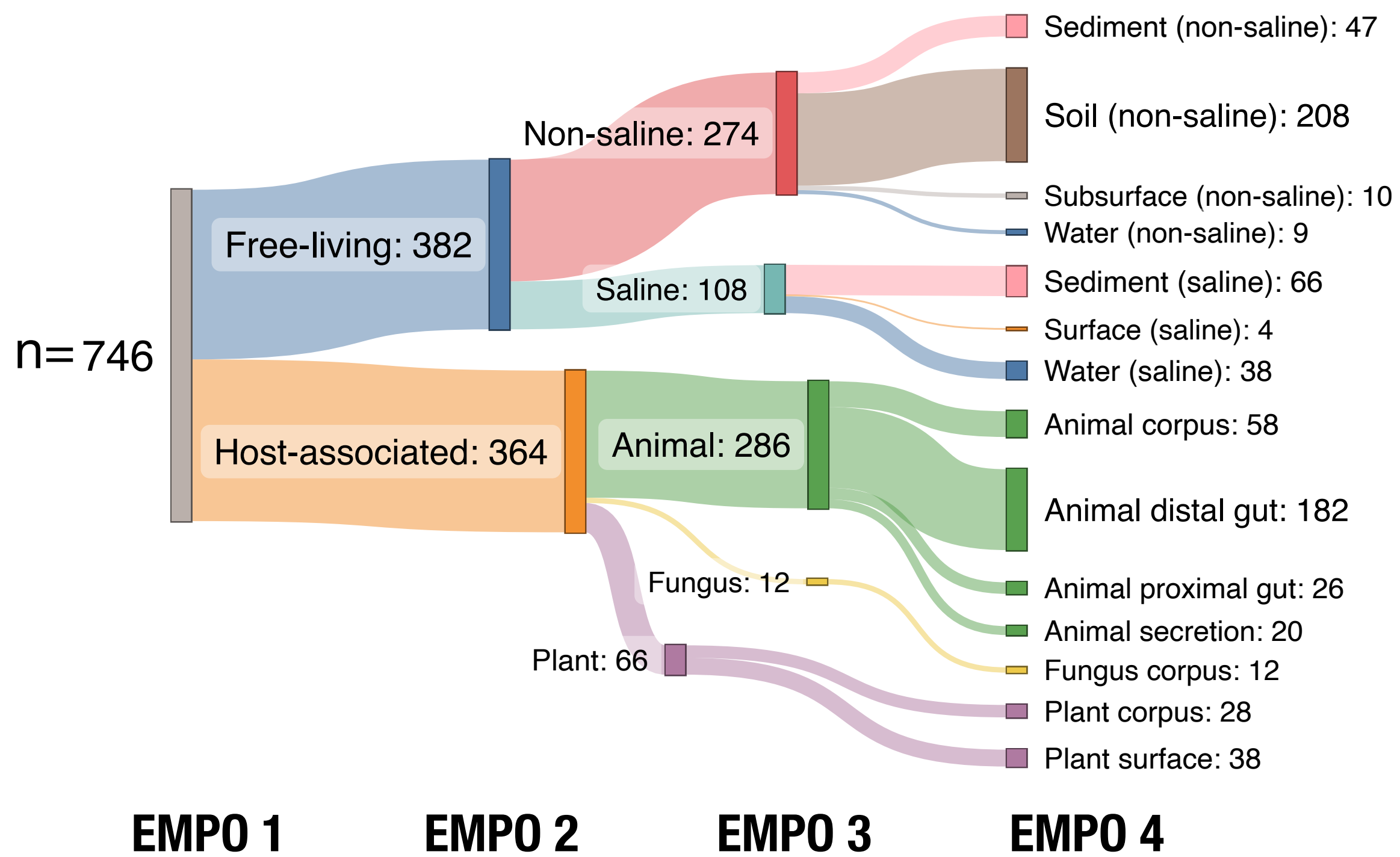


Reference: Web of Life (v1)  
11,000 microbial genomes

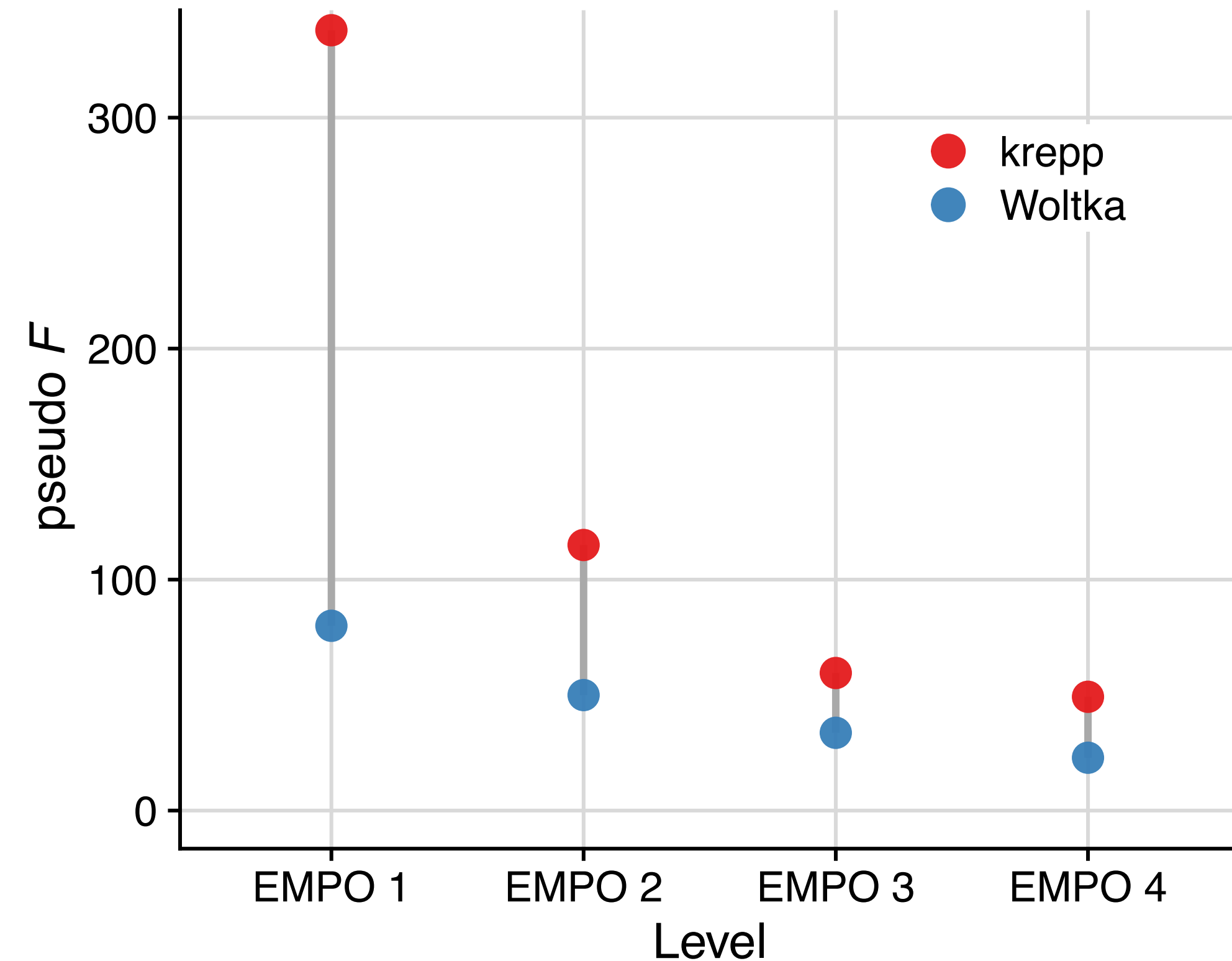


# Improved characterization of earth's microbiome

## Hierarchical categorization of earth microbiome samples:



Reference: Web of Life (v1)  
11,000 microbial genomes



# Conclusion

## krepp

- estimates read to genome distances **>10x faster** than alignment
- extends to more distant references and **maps novel reads**
- is the **only method** that can accurately **place genome-wide reads**
  - ▶ **One exception:** App-SpaM; cannot scale beyond 50-100 microbial genomes
- offers significantly **better characterization of metagenomic samples**



Siavash Mirarab



Siavash Mirarab

**Questions?**

# **Extra Slides**

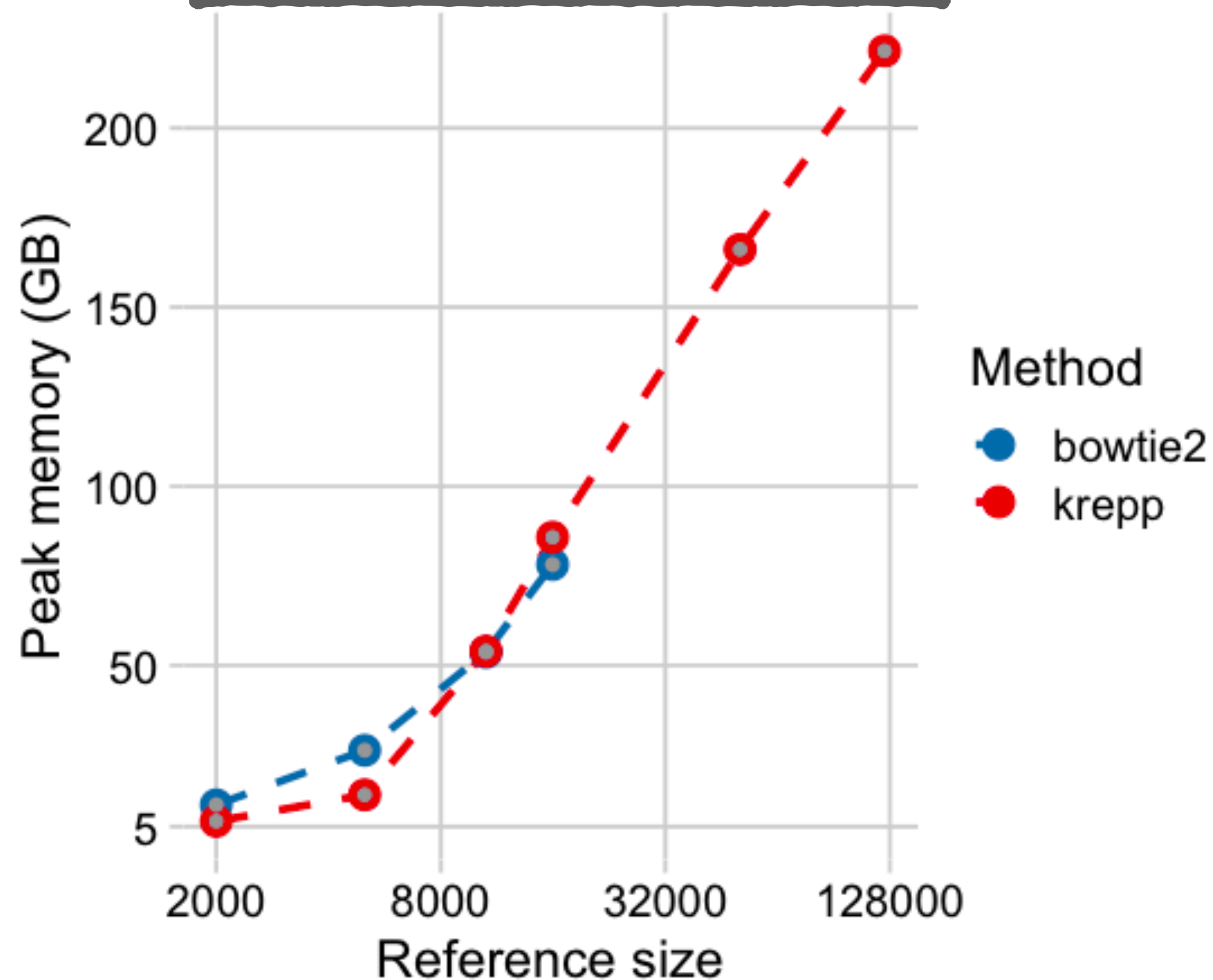
# Scalability:

## krepp can be distributed and has flexible memory requirements

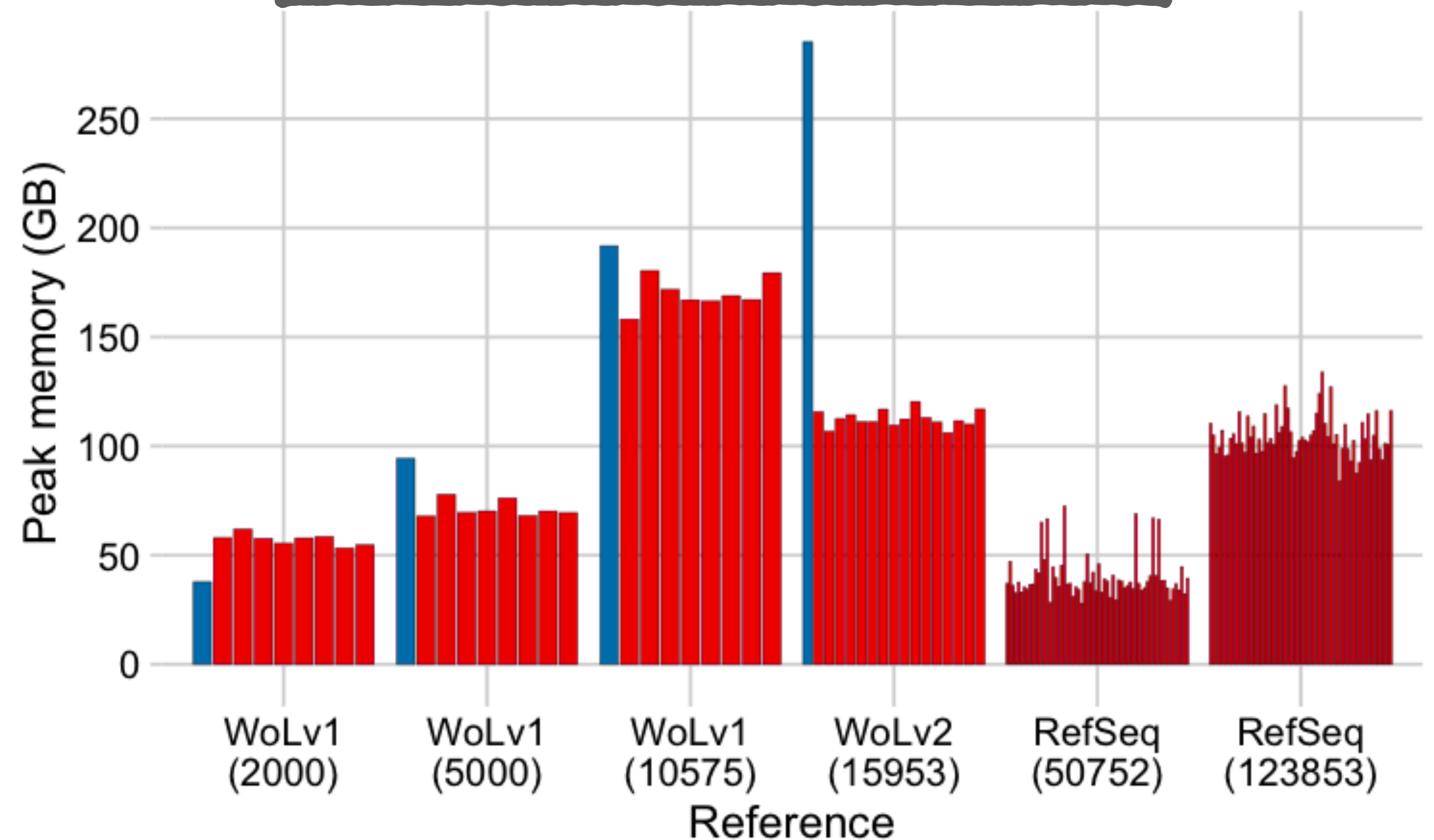
Mapping 10M reads (16 threads):

Indexing microbial genomes (32 threads):

reducing memory use  
w/ further subsampling...



adjusting partitioning based on the  
input size & available memory...



Let  $b_i$  be the length of the branch  $i$  and  $p_i^A$  and  $p_i^B$  are the taxa proportions descending from the branch  $i$  for community  $A$  and  $B$ , respectively.

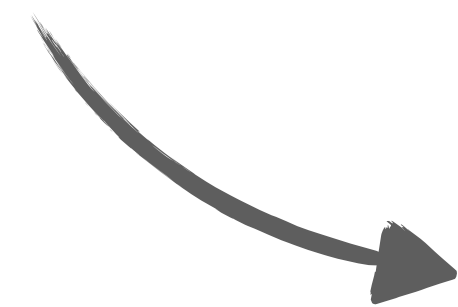
$$d(A, B) = \frac{\sum_i^n b_i |p_i^A - p_i^B|}{\sum_i^n b_i (p_i^A + p_i^B)}$$

$$SS_T = \frac{1}{N} \sum_{i=1}^{N-1} \sum_{j=i+1}^N d_{ij}^2 \quad SS_W = \frac{1}{n} \sum_{i=1}^{N-1} \sum_{j=i+1}^N d_{ij}^2 \delta_{ij}$$

$$SS_A = SS_T - SS_W$$

$$F = \frac{\left( \frac{SS_A}{p-1} \right)}{\left( \frac{SS_W}{N-p} \right)}$$

Multiple permutations  
to get a  $p$ -value!



$N$  is the number of groups,  $p$  is the number of objects in each group.