# Memory-bound and taxonomy-aware k-mer selection for large reference databases
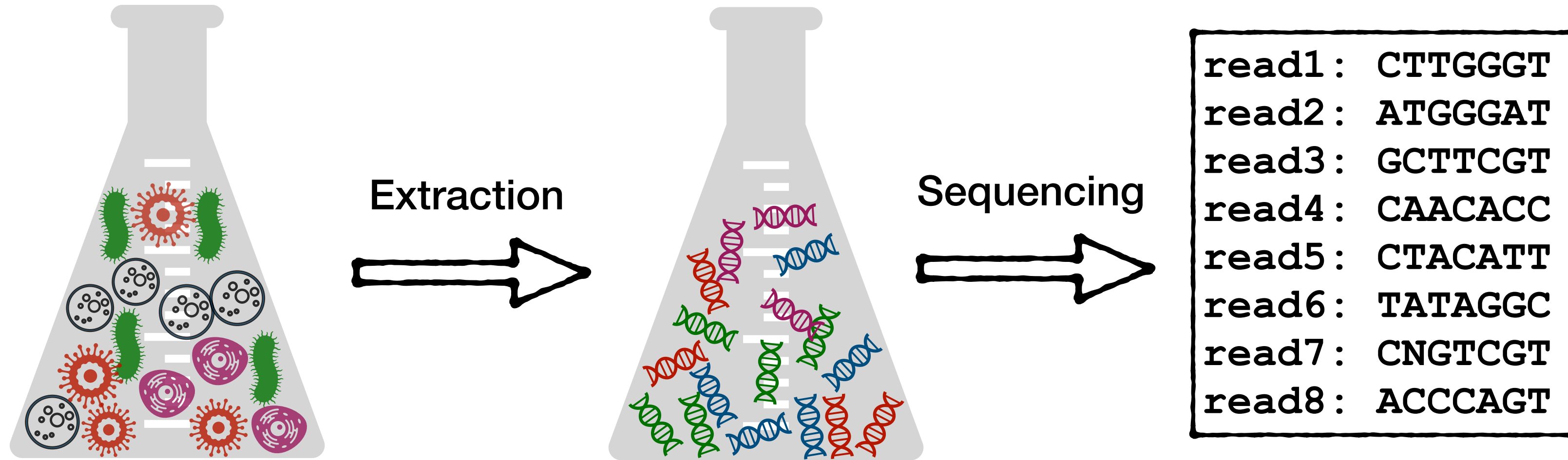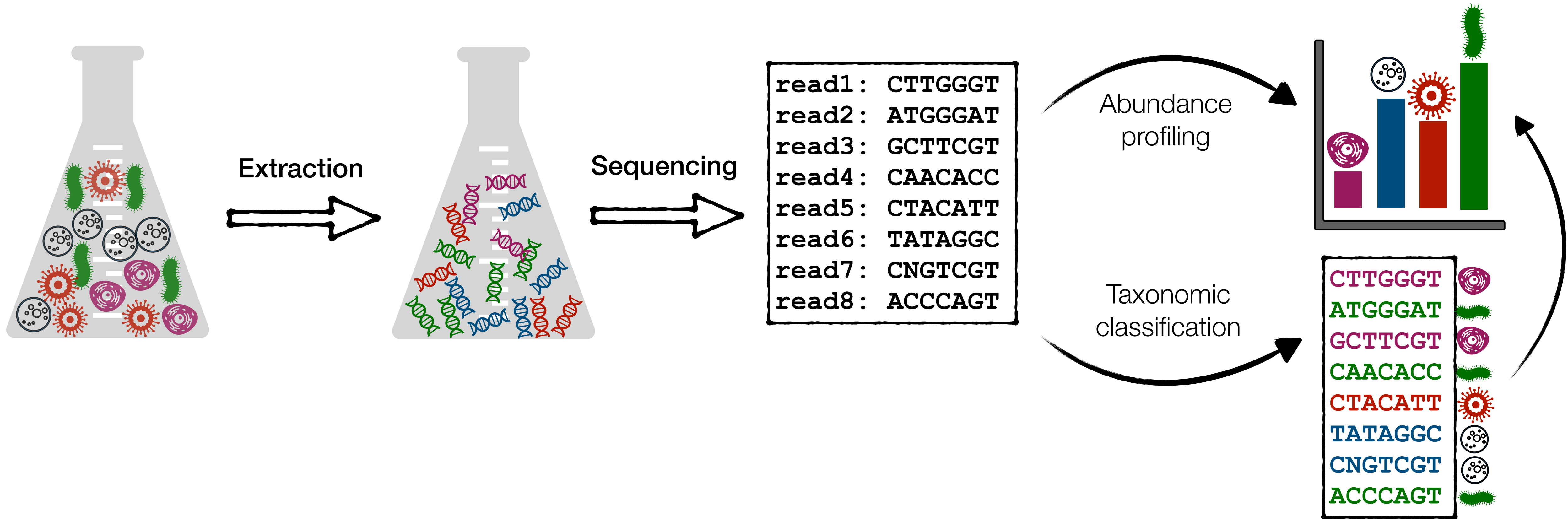
Ali Osman Berk Şapcı & Siavash Mirarab
UC San Diego
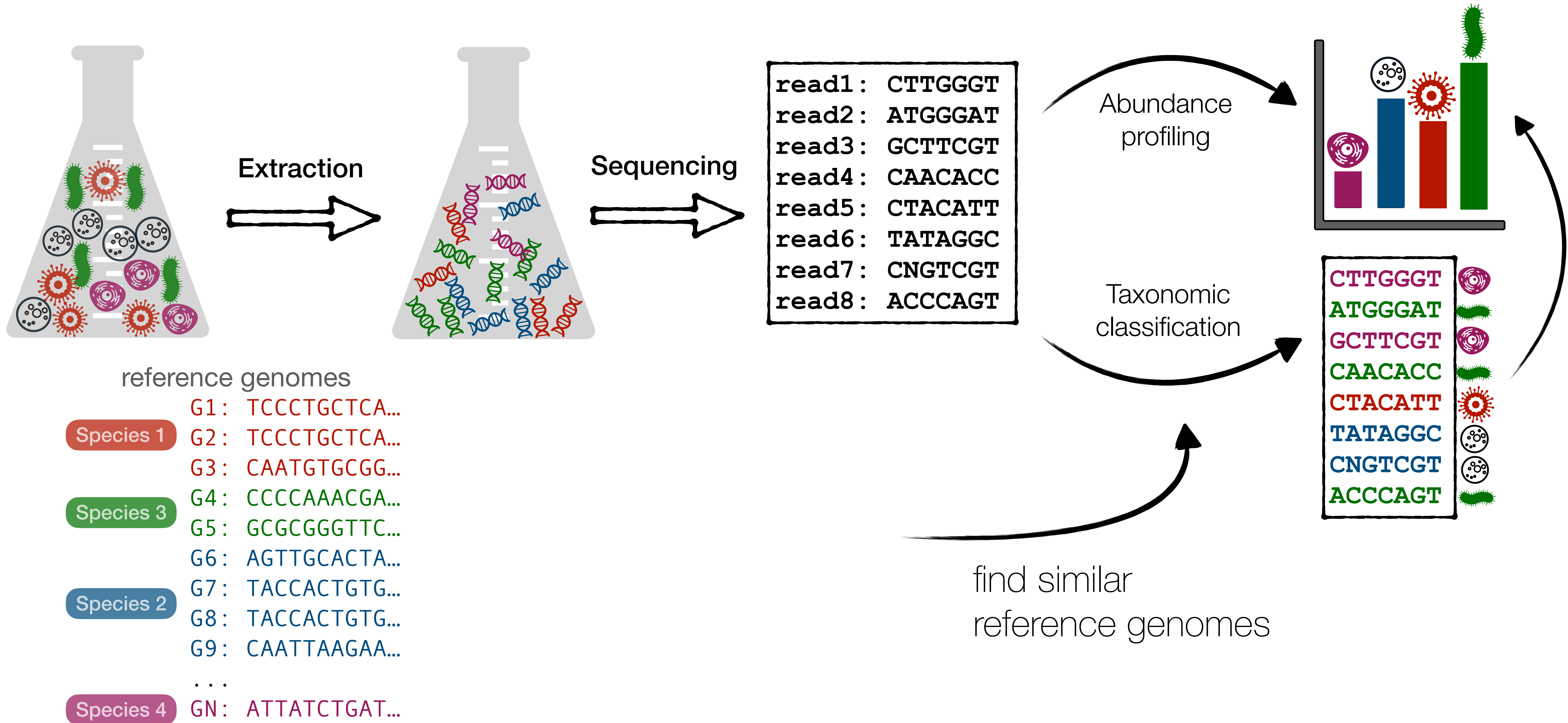
Bioinformatics & Systems Biology

UC San Diego
Electrical and Computer Engineering
JACOBS SCHOOL OF ENGINEERING
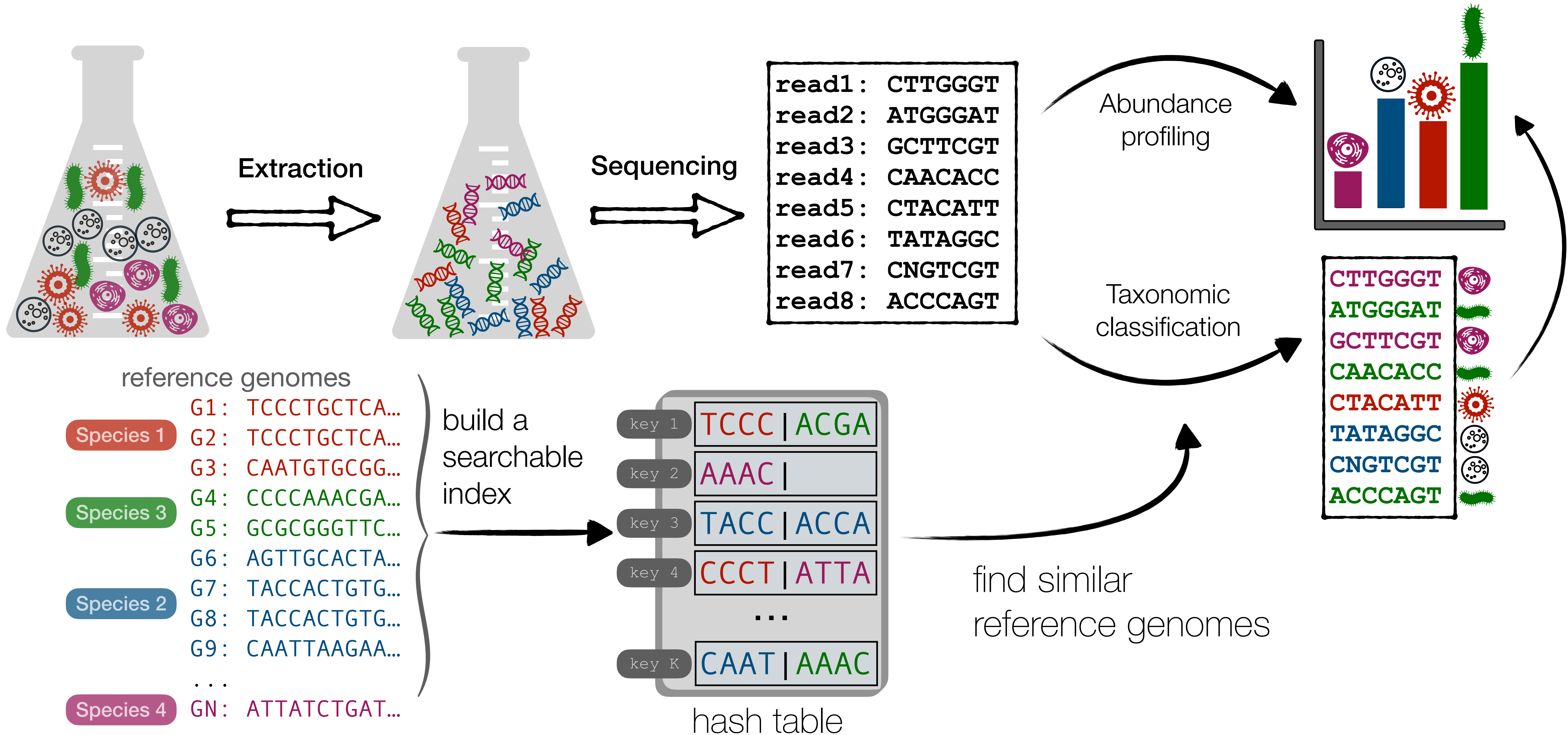
# Identifying metagenomic sequences

# Identifying metagenomic sequences

# Identifying metagenomic sequences

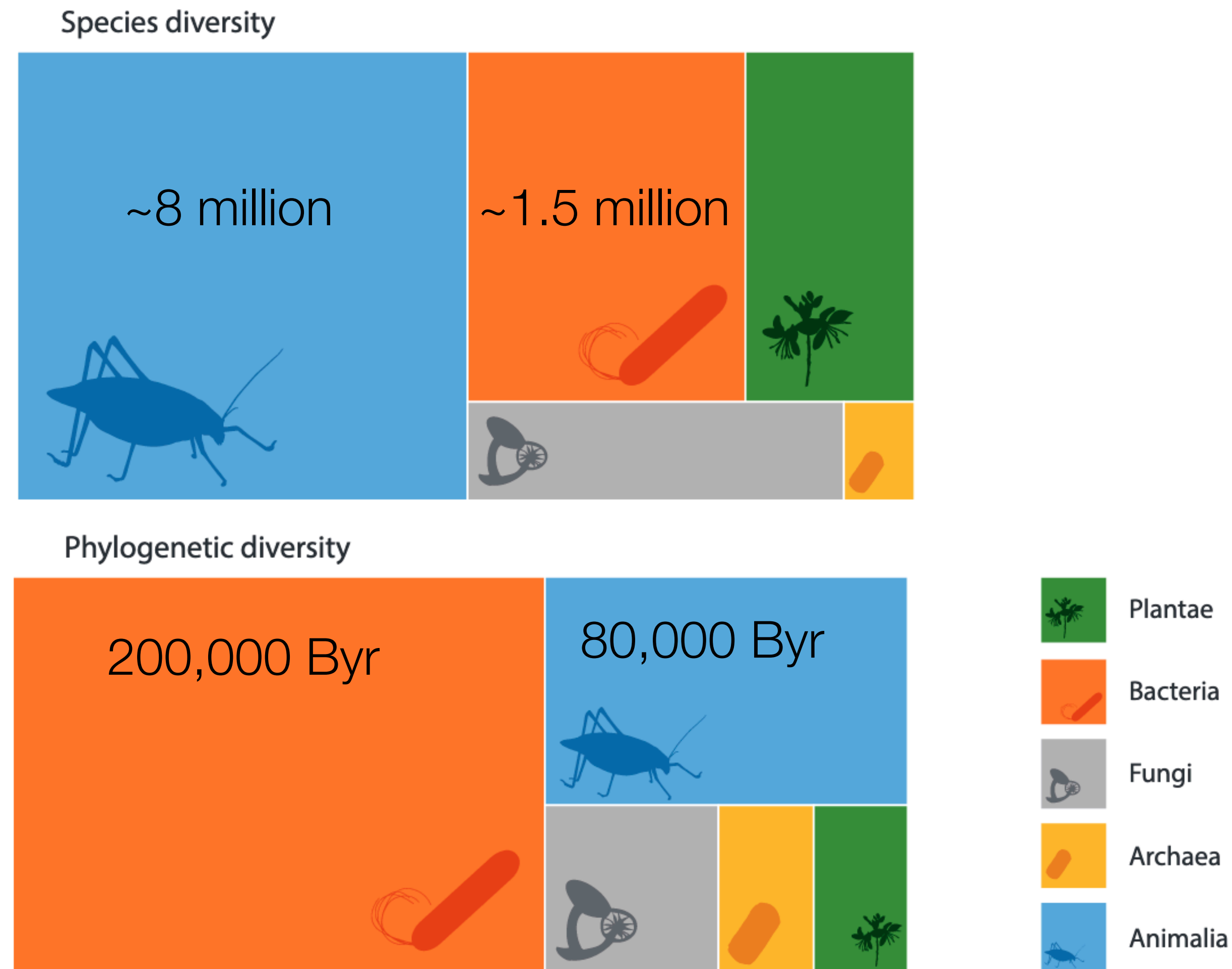# Identifying metagenomic sequences

# Novel sequences challenge popular tools

- Reference databases (and indexes) remain incomplete compared to all species…

  and there is a rich diversity within species!



Species diversity

~8 million   ~1.5 million

Phylogenetic diversity

200,000 Byr   80,000 Byr

Plantae

Bacteria

Fungi

Archaea

Animalia

[Díaz et al., 2022]

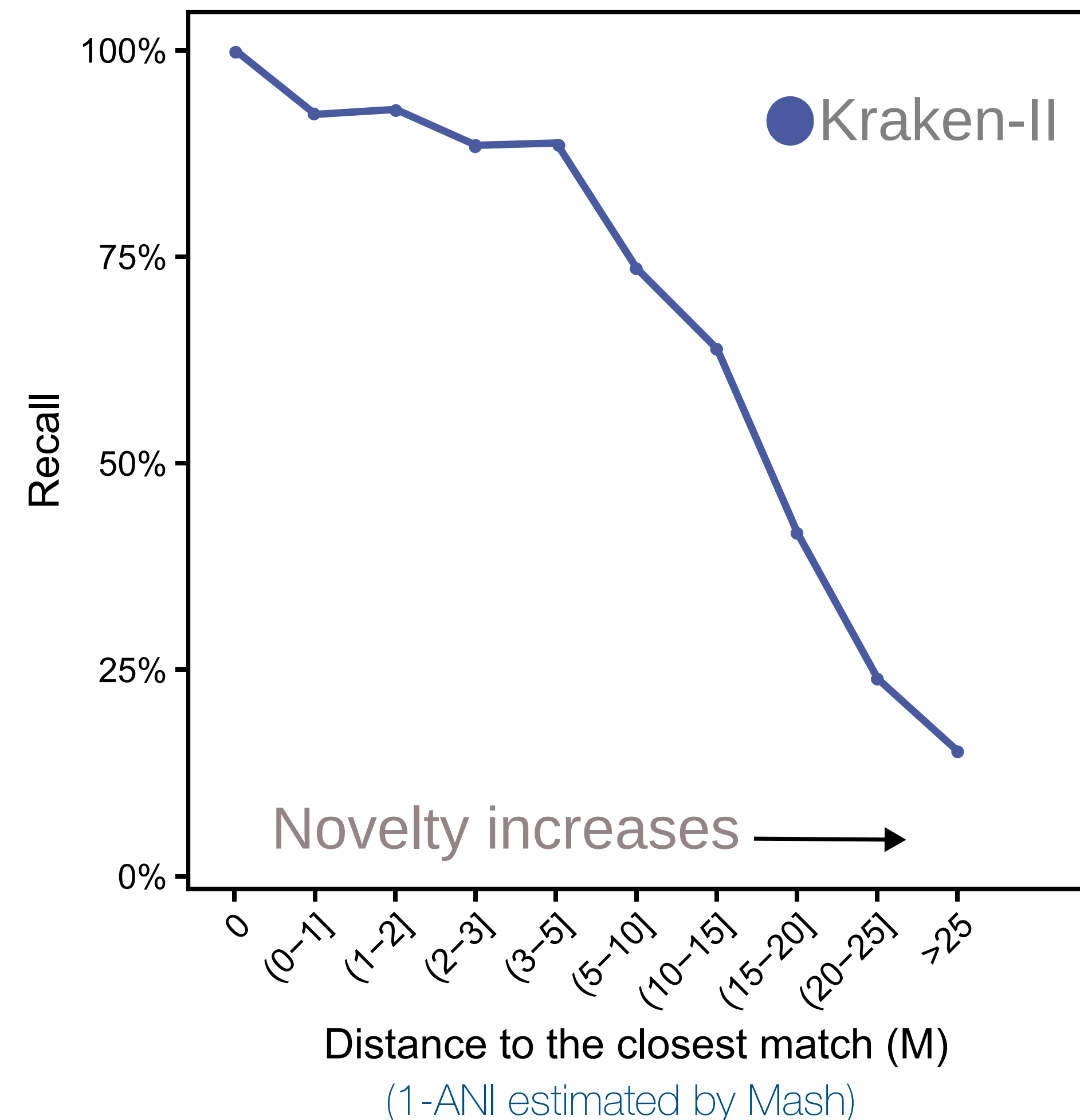# Novel sequences challenge popular tools

- Reference databases (and indexes) remain incomplete compared to all species…

  and there is a rich diversity within species!

- **Novel sequences:** sequences which lack a close matching reference genome



Distance to the closest match (M)

(1-ANI estimated by Mash)

[Rachtman et al., 2019]

2

# Solutions for identifying novel queries w/ limited resources

▸ find distant matches → increase sensitivity of the search

▸ enhance the reference set → utilize more genomes & larger databases

# Solutions for identifying novel queries w/ limited resources
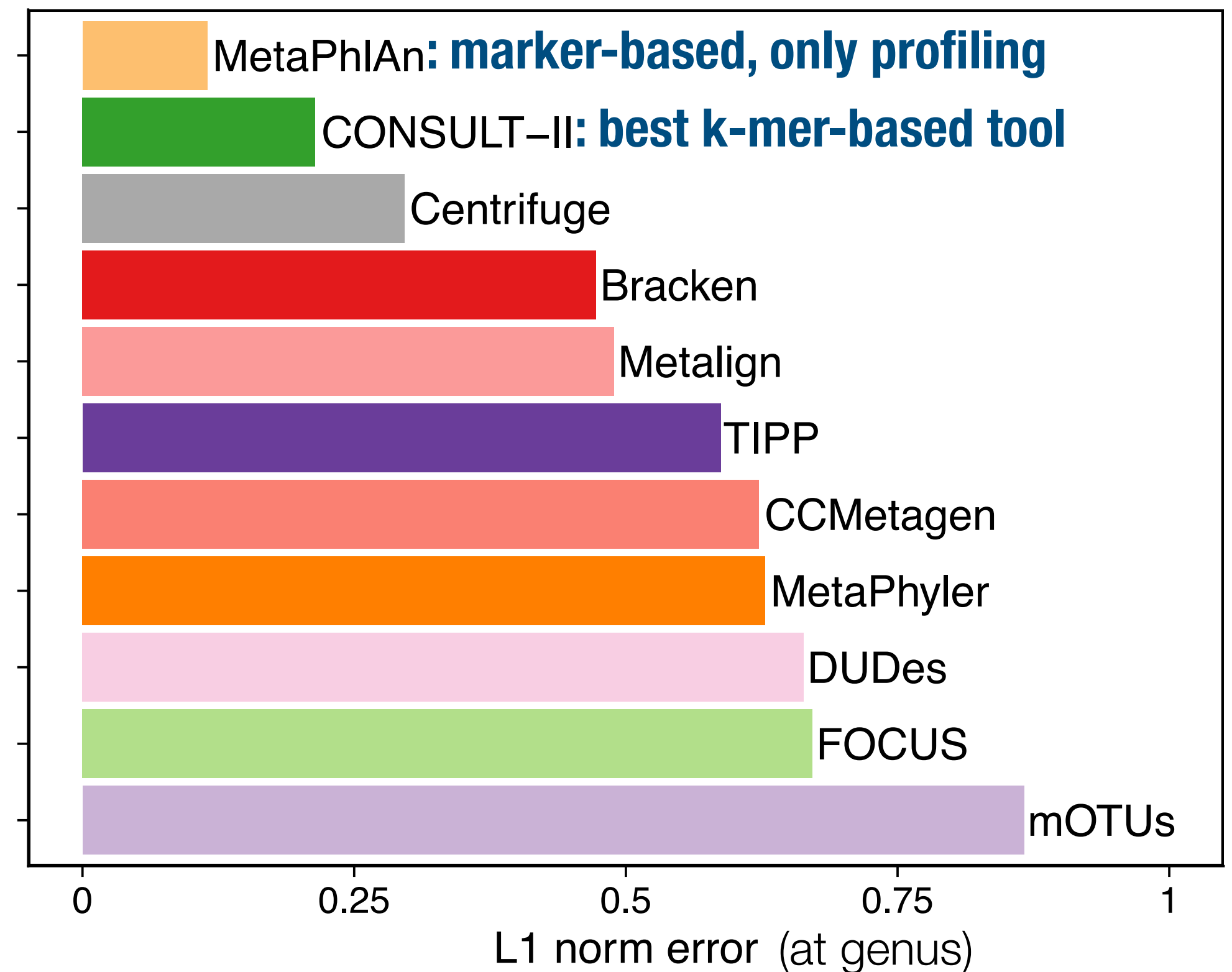
▸ find distant matches → increase sensitivity of the search

▸ enhance the reference set → utilize more genomes & larger databases

Computing the Hamming distances of inexact matches

**CONSULT-II: accurate taxonomic identification and profiling using locality-sensitive hashing**

Ali Osman Berk Şapcı [ID] [1], Eleonora Rachtman [ID] [1], Siavash Mirarab [ID] [1,2,*]

# Solutions for identifying novel queries w/ limited resources

▸ find distant matches → increase sensitivity of the search

▸ enhance the reference set → utilize more genomes & larger databases

**Strain–madness dataset** [CAMI-II]

Computing the Hamming distances of inexact matches

**CONSULT-II: accurate taxonomic identification and profiling using locality-sensitive hashing**

Ali Osman Berk Şapcı [1], Eleonora Rachtman [1], Siavash Mirarab [1,2,*]



(using a RefSeq snapshot from 2019 with ~130k genomes)

# Can we use more reference genomes?

▸ find distant matches → increase sensitivity of the search

▸ enhance the reference set → utilize more genomes & larger databases

# Can we use more reference genomes?

▸ find distant matches → increase sensitivity of the search

▸ enhance the reference set → utilize more genomes & larger databases

- **Challenge:** very large & diverse databases
  have too many $k$-mers to fit in the memory

  ▸ Limited to a selected subset

- **This talk:** — **KRANK**

  ▸ Selecting a representative subset of $k$-mers
    + classification/profiling using CONSULT-II

4

# Can we use more reference genomes?

▸ find distant matches → increase sensitivity of the search

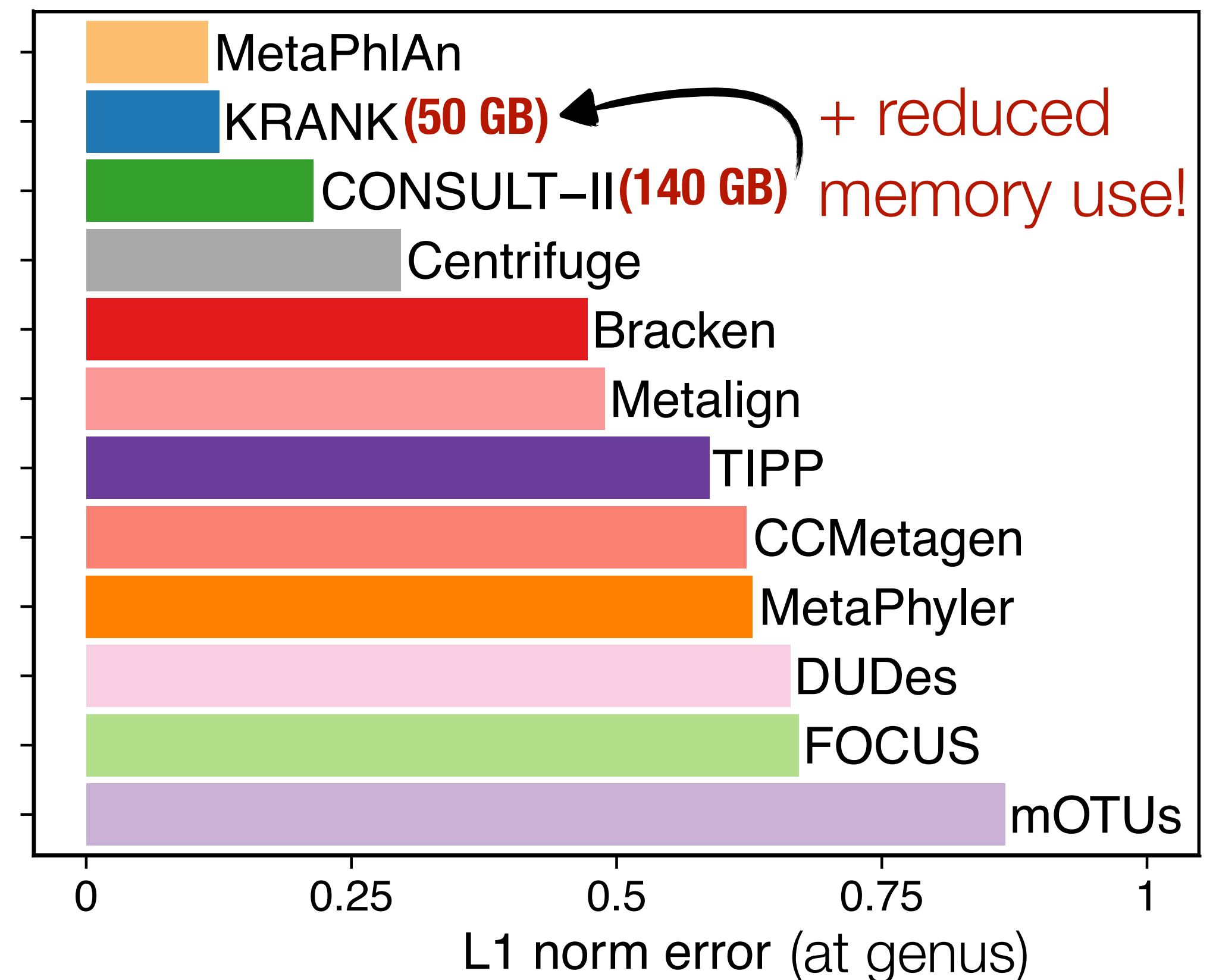▸ enhance the reference set → utilize more genomes & larger databases

**Strain–madness dataset** [CAMI-II]

- **Challenge:** very large & diverse databases have too many *k*-mers to fit in the memory

  ▸ Limited to a selected subset

- **This talk:** — **KRANK**

  ▸ Selecting a representative subset of *k*-mers + classification/profiling using CONSULT-II



+ reduced memory use!

L1 norm error (at genus)

(using a RefSeq snapshot from 2019 with ~130k genomes)

# Problem statement

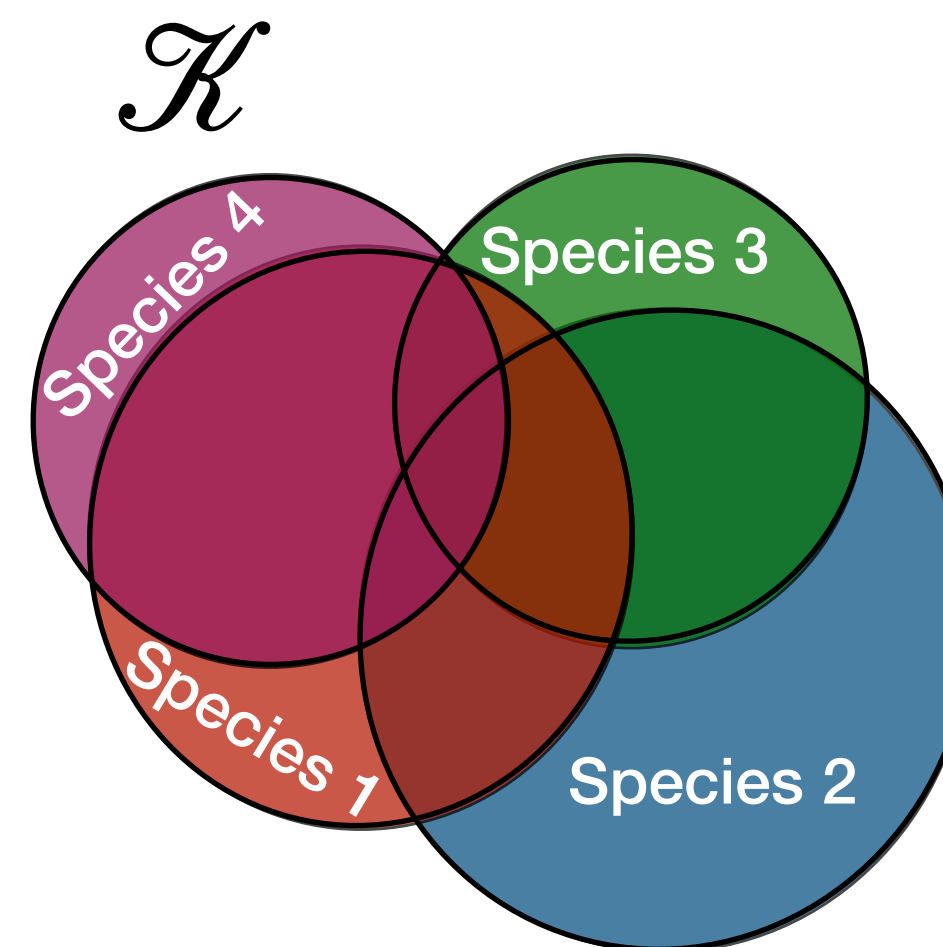- Given:

  1. *k*-mer set $\mathcal{K}$ of a large collection of genomes

  2. limited budget $M < |\mathcal{K}|$

  3. taxonomy

```
         G1:  TCCCTGCTCAGTGGTATATGGTTTTTGCTA…
Species 1 G2:  TCCCTGCTCAGCCCCATATGGTTTTTGCTA…
         G3:  CAATGTGCGGATGGCGTTACGACTTACTGG…
         G4:  CCCCAAACGATGCTGAAGGCTCAGGTTACA…
Species 3 G5:  GCGCGGGTTCCCGCCCTCAACCCGGGCCGA…
         G6:  AGTTGCACTACTTCTGCGACCCAAATGCAC…
         G7:  TACCACTGTGTTCGTGTCATCTAGGACGGG…
Species 2 G8:  TACCACTGTGTTCGTGTCATCTAGGACGGG…
         G9:  CAATTAAGAATACCTTATATTATTGTACAC…
            . . .
Species 4 GN:  ATTATCTGATTTTATATTATGATTTTAGTA…
```
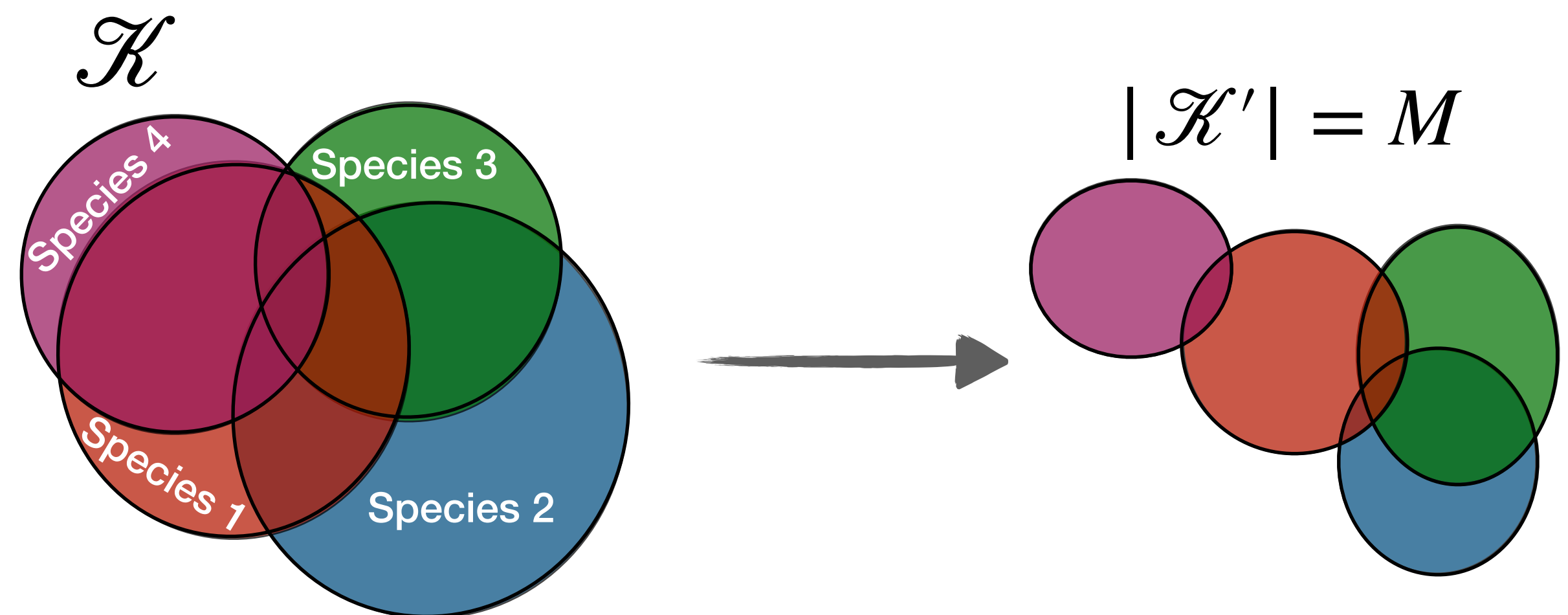
# Problem statement

- Given:
  1. *k*-mer set $\mathcal{K}$ of a large collection of genomes
  2. limited budget $M < |\mathcal{K}|$
  3. taxonomy

- Select a subset with size $M$ such that the collection is well represented



```
           G1: TCCCTGCTCAGTGGTATATGGTTTTTGCTA…
Species 1  G2: TCCCTGCTCAGCCCCATATGGTTTTTGCTA…
           G3: CAATGTGCGGATGGCGTTACGACTTACTGG…
           G4: CCCCAAACGATGCTGAAGGCTCAGGTTACA…
Species 3  G5: GCGCGGGTTCCCGCCCTCAACCCGGGCCGA…
           G6: AGTTGCACTACTTCTGCGACCCAAATGCAC…
           G7: TACCACTGTGTTCGTGTCATCTAGGACGGG…
Species 2  G8: TACCACTGTGTTCGTGTCATCTAGGACGGG…
           G9: CAATTAAGAATACCTTATATTATTGTACAC…
           . . .
Species 4  GN: ATTATCTGATTTTATATTATGATTTTAGTA…
```

$\mathcal{K}$

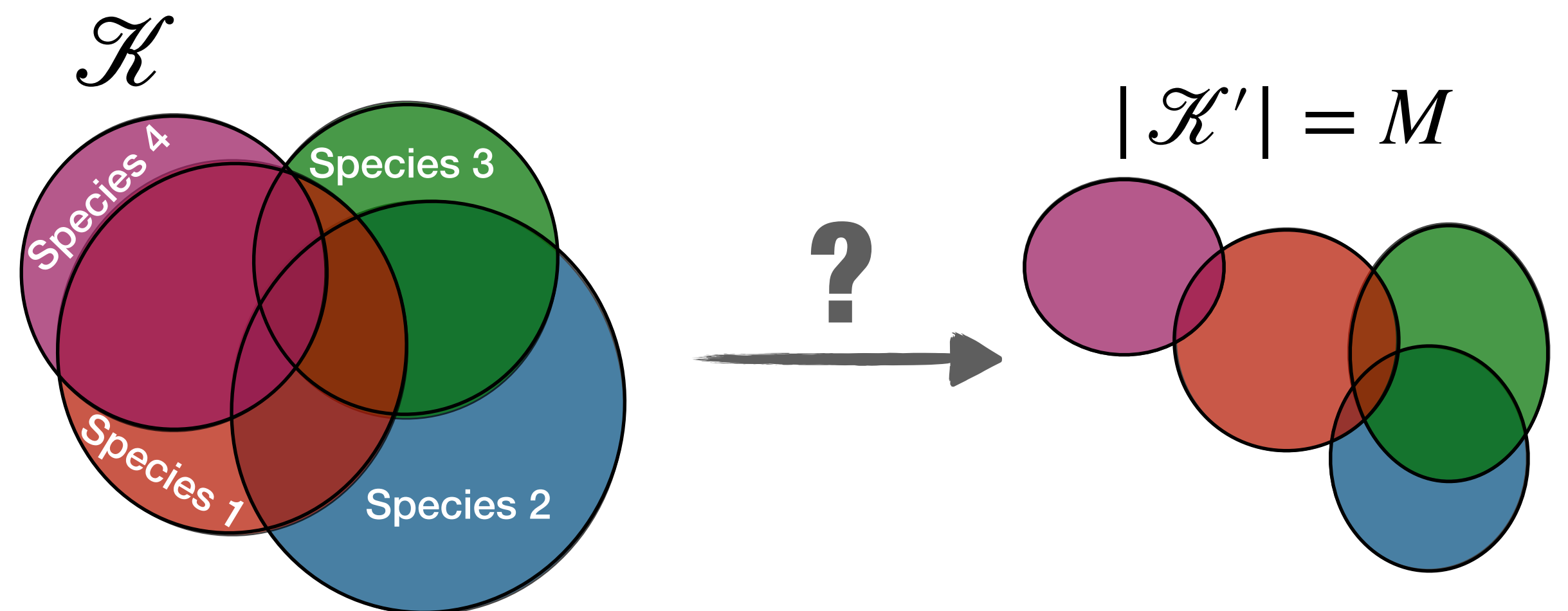$|\mathcal{K}'| = M$

# Problem statement

- Given:

  1. $k$-mer set $\mathcal{K}$ of a large collection of genomes

  2. limited budget $M < |\mathcal{K}|$

  3. taxonomy

  high accuracy in taxonomic identification

- Select a subset with size $M$ such that the collection is well represented



```
       G1:  TCCCTGCTCAGTGGTATATGGTTTTTGCTA…
Species 1  G2:  TCCCTGCTCAGCCCCATATGGTTTTTGCTA…
       G3:  CAATGTGCGGATGGCGTTACGACTTACTGG…
       G4:  CCCCAAACGATGCTGAAGGCTCAGGTTACA…
Species 3  G5:  GCGCGGGTTCCCGCCCTCAACCCGGGCCGA…
       G6:  AGTTGCACTACTTCTGCGACCCAAATGCAC…
       G7:  TACCACTGTGTTCGTGTCATCTAGGACGGG…
Species 2  G8:  TACCACTGTGTTCGTGTCATCTAGGACGGG…
       G9:  CAATTAAGAATACCTTATATTATTGTACAC…
       . . .
Species 4  GN:  ATTATCTGATTTTATATTATGATTTTAGTA…
```

$\mathcal{K}$

$|\mathcal{K}'| = M$

?

5

# Reducing the reference set by selecting k-mers

G1: TCCCTGC
    CCCTGCT
     CCTGCTC
      CTGCTCA…
G2: TCGCTAC
    CGCTACG
     GCTACGC
      CTACGCG…
G3: CAATGTG
    AATGTGC
     ATGTGCG
      TGTGCGG…
G5: GCGCGGG
    CGCGGGT
     GCGGGTT
      CGGGTTC…
G4: CCCCAAA
    CCCAAAC
     CCAAACG
      CAAACGT…

6

# Reducing the reference set by selecting k-mers

- **Baseline:** random selection

G1: TCCCTGC
     CCCTGCT
      CCTGCTC
       CTGCTCA…
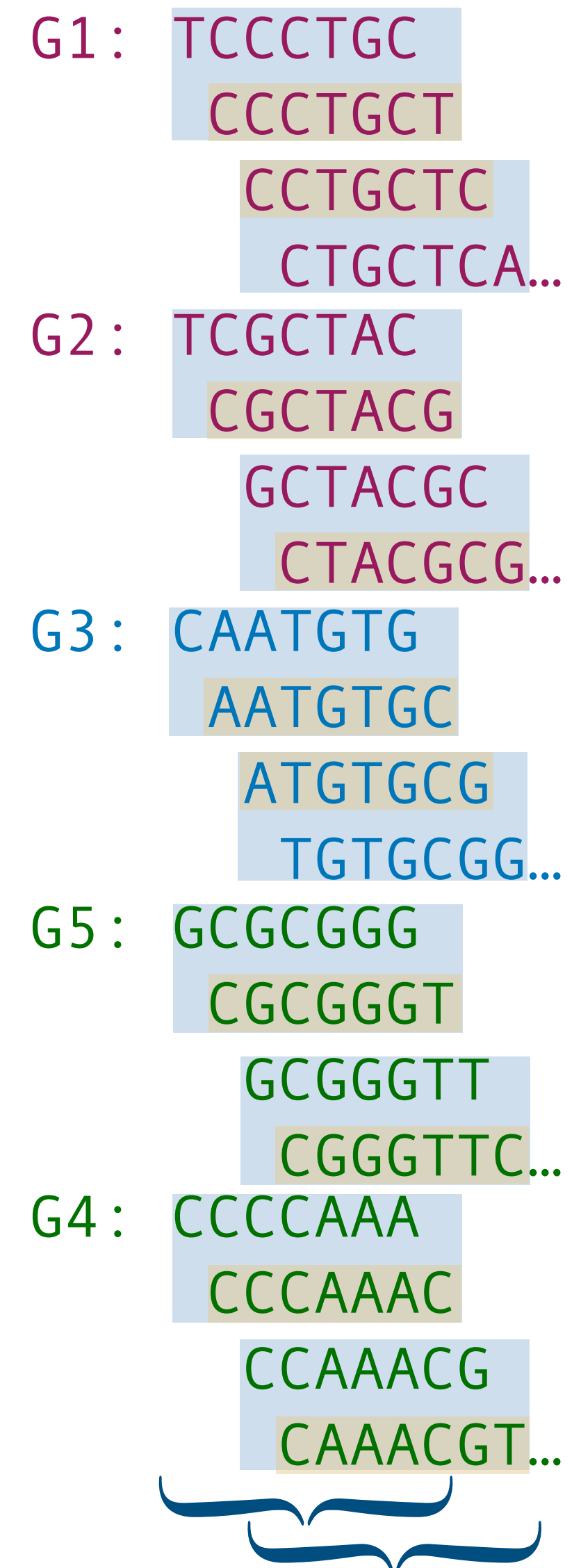G2: TCGCTAC
     CGCTACG
      GCTACGC
       CTACGCG…
G3: CAATGTG
     AATGTGC
      ATGTGCG
       TGTGCGG…
G5: GCGCGGG
     CGCGGGT
      GCGGGTT
       CGGGTTC…
G4: CCCCAAA
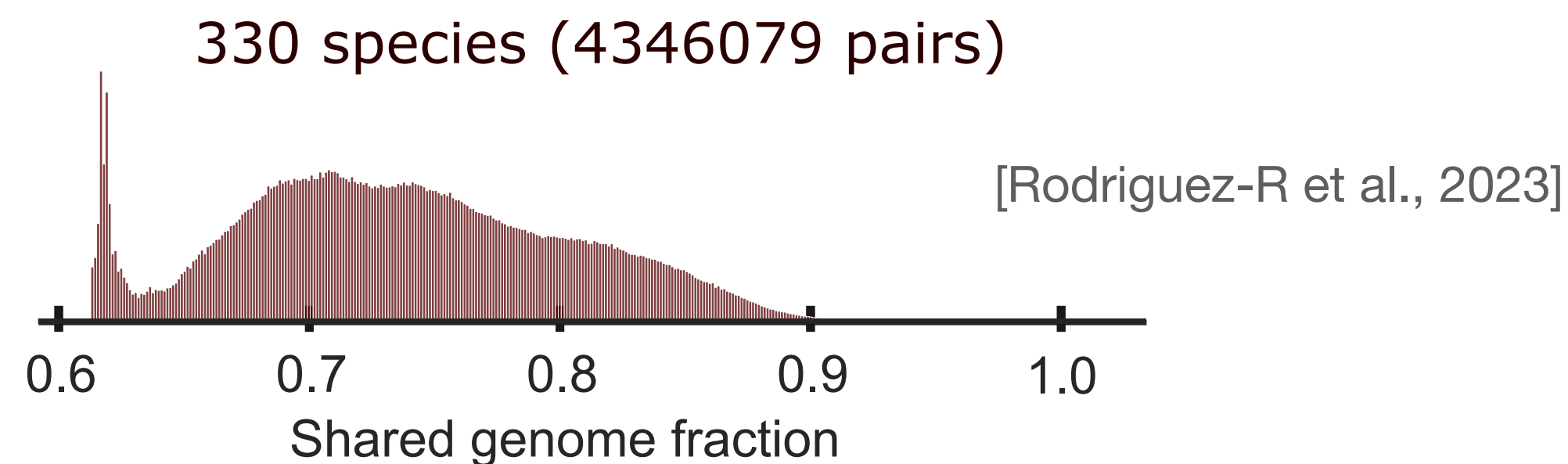     CCCAAAC
      CCAAACG
       CAAACGT…

# Reducing the reference set by selecting k-mers

- **Baseline:** random selection

- **Minimizers:** selecting one among overlapping *k*-mers with a sliding window

G1: TCCCTGC
      CCCTGCT
        CCTGCTC
          CTGCTCA...

G2: TCGCTAC
      CGCTACG
        GCTACGC
          CTACGCG...

G3: CAATGTG
      AATGTGC
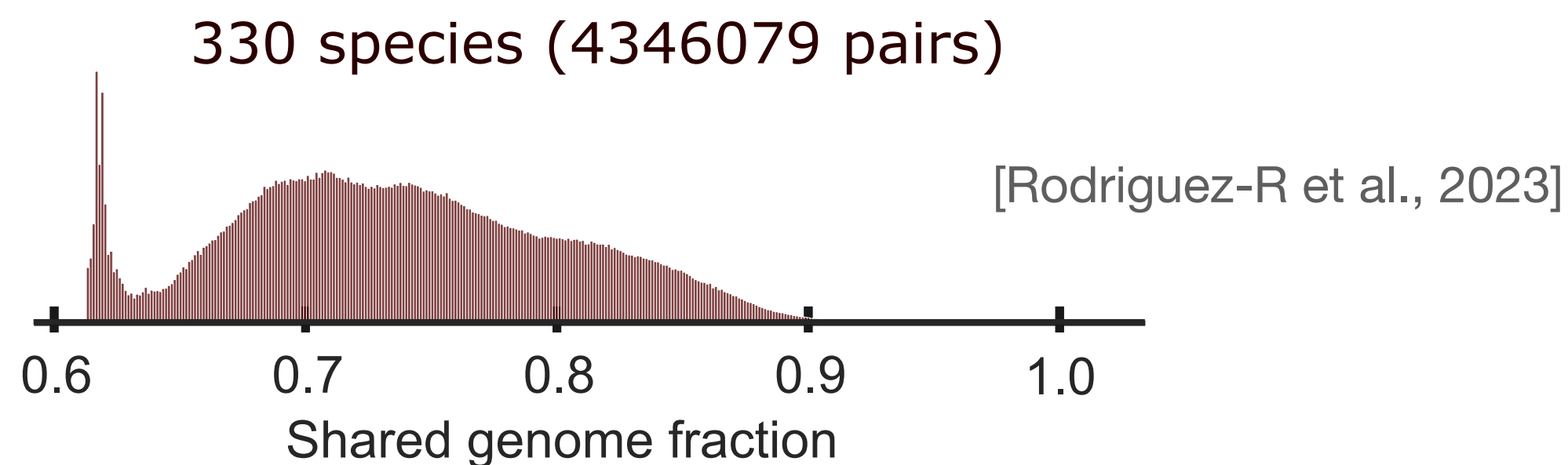        ATGTGCG
          TGTGCGG...

G5: GCGCGGG
      CGCGGGT
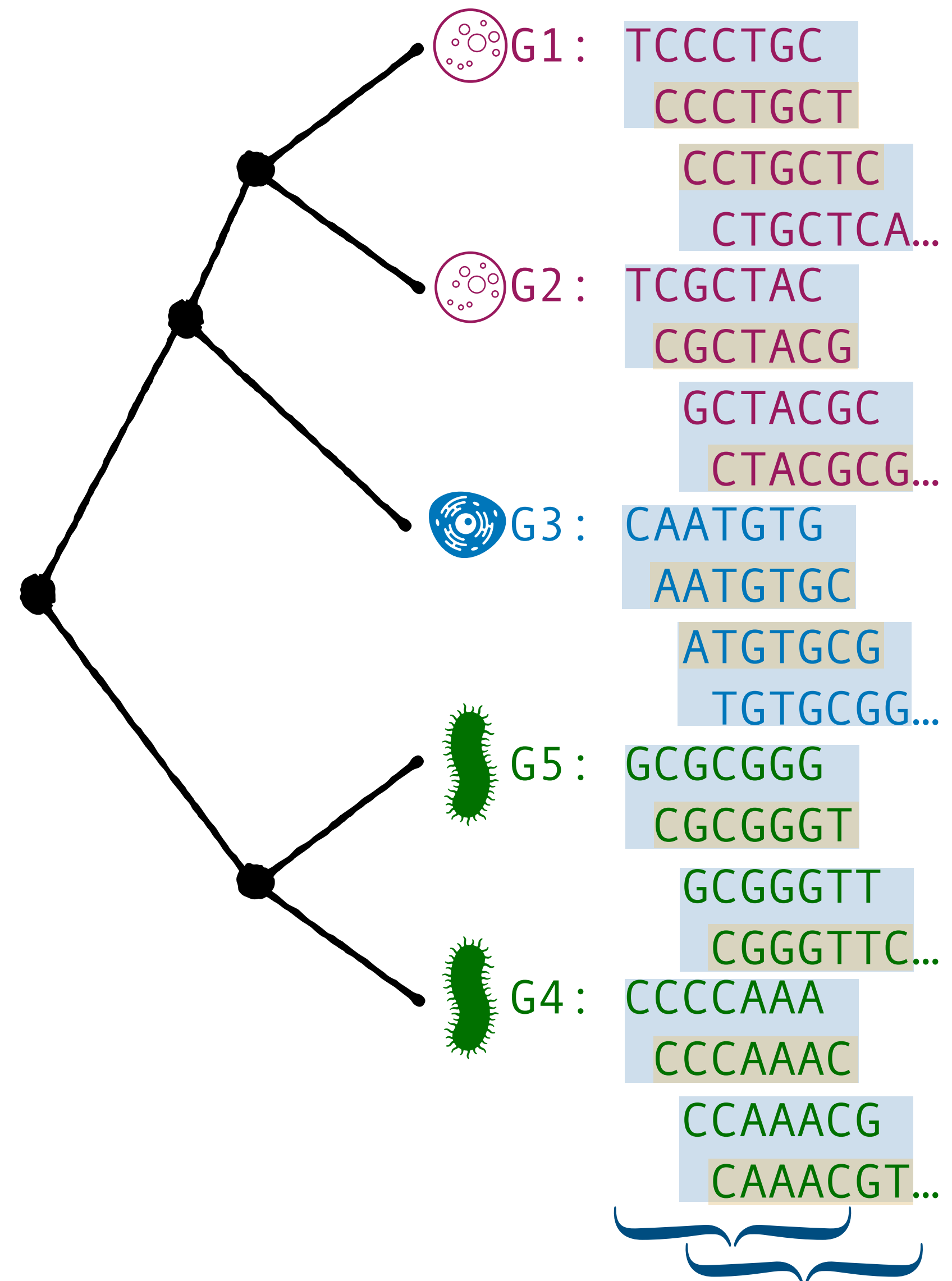        GCGGGTT
          CGGGTTC...

G4: CCCCAAA
      CCCAAAC
        CCAAACG
          CAAACGT...

# Reducing the reference set by selecting k-mers

- **Baseline:** random selection

- **Minimizers:** selecting one among overlapping *k*-mers with a sliding window

- Even with minimizers, number of distinct *k*-mers grows fast with the number of genomes

330 species (4346079 pairs)



[Rodriguez-R et al., 2023]

Shared genome fraction

G1: TCCCTGC
    CCCTGCT
    CCTGCTC
    CTGCTCA...

G2: TCGCTAC
    CGCTACG
    GCTACGC
    CTACGCG...

G3: CAATGTG
    AATGTGC
    ATGTGCG
    TGTGCGG...

G5: GCGCGGG
    CGCGGGT
    GCGGGTT
    CGGGTTC...

G4: CCCCAAA
    CCCAAAC
    CCAAACG
    CAAACGT...

# Reducing the reference set by selecting k-mers

- **Baseline:** random selection

- **Minimizers:** selecting one among overlapping *k*-mers with a sliding window

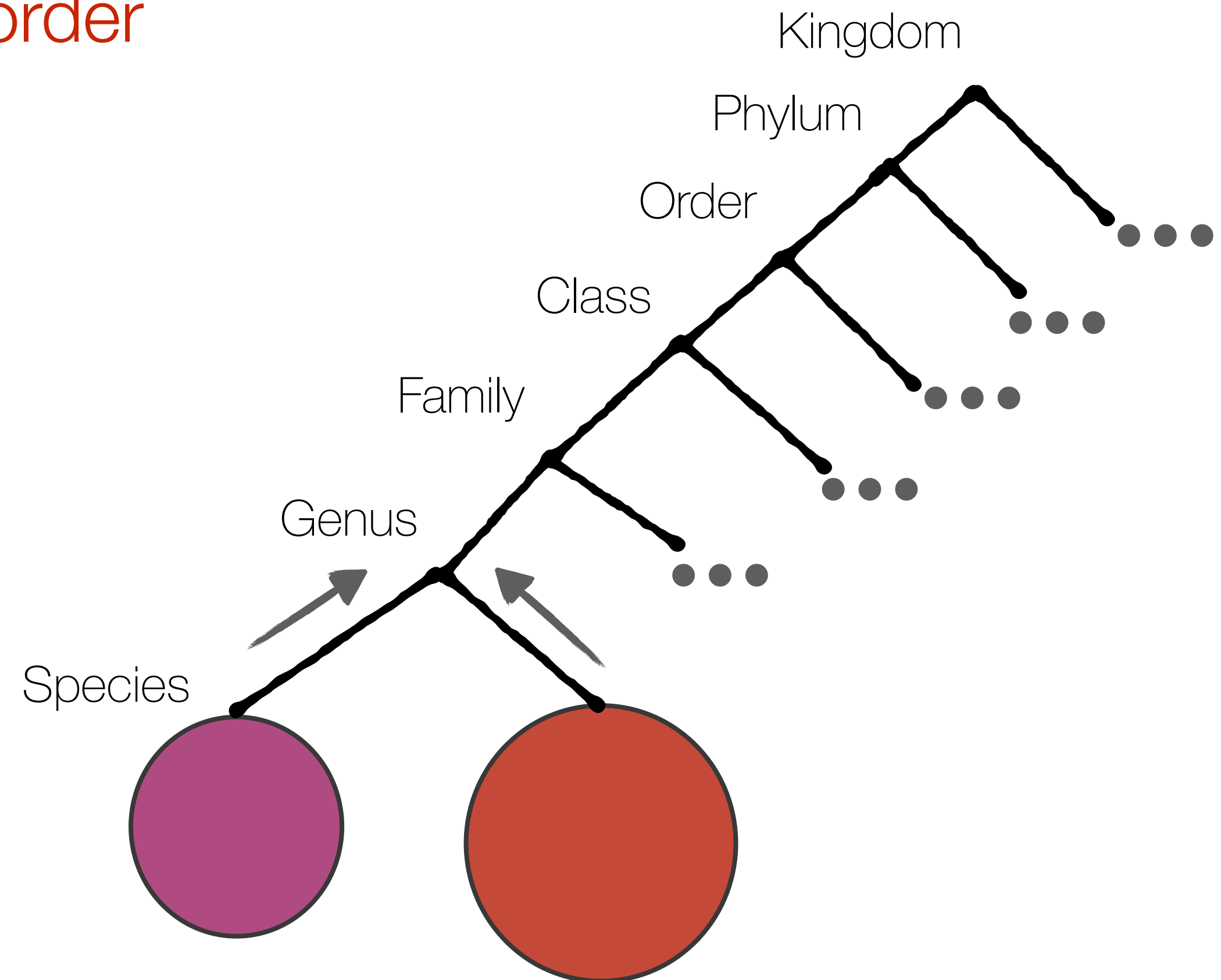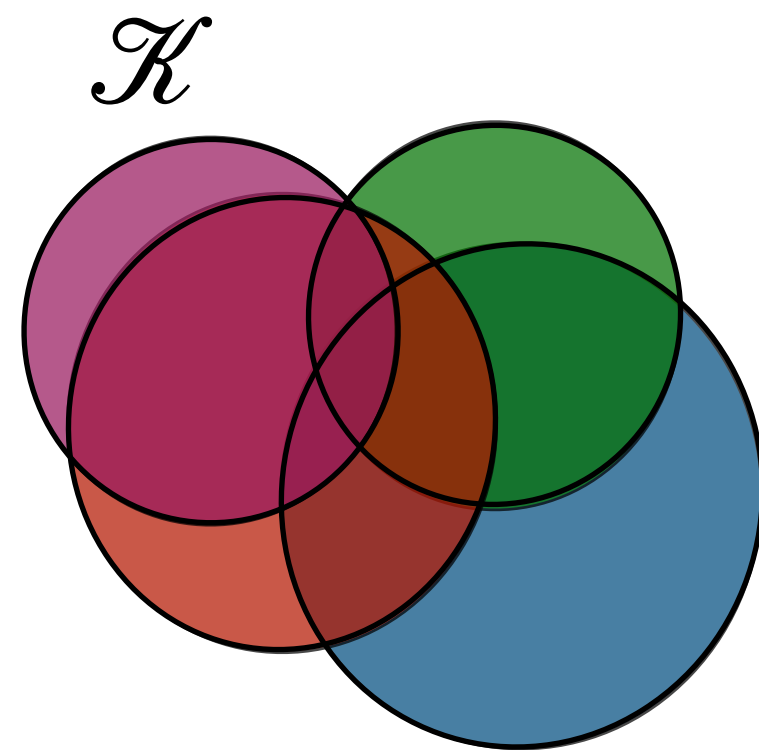- Even with minimizers, number of distinct *k*-mers grows fast with the number of genomes

330 species (4346079 pairs)

[Rodriguez-R et al., 2023]

Shared genome fraction

0.6    0.7    0.8    0.9    1.0

- Additionally, exploit the evolutionary dimension

G1 : TCCCTGC
       CCCTGCT
         CCTGCTC
           CTGCTCA...

G2 : TCGCTAC
       CGCTACG
         GCTACGC
           CTACGCG...

G3 : CAATGTG
       AATGTGC
         ATGTGCG
           TGTGCGG...

G5 : GCGCGGG
       CGCGGGT
         GCGGGTT
           CGGGTTC...

G4 : CCCCAAA
       CCCAAAC
         CCAAACG
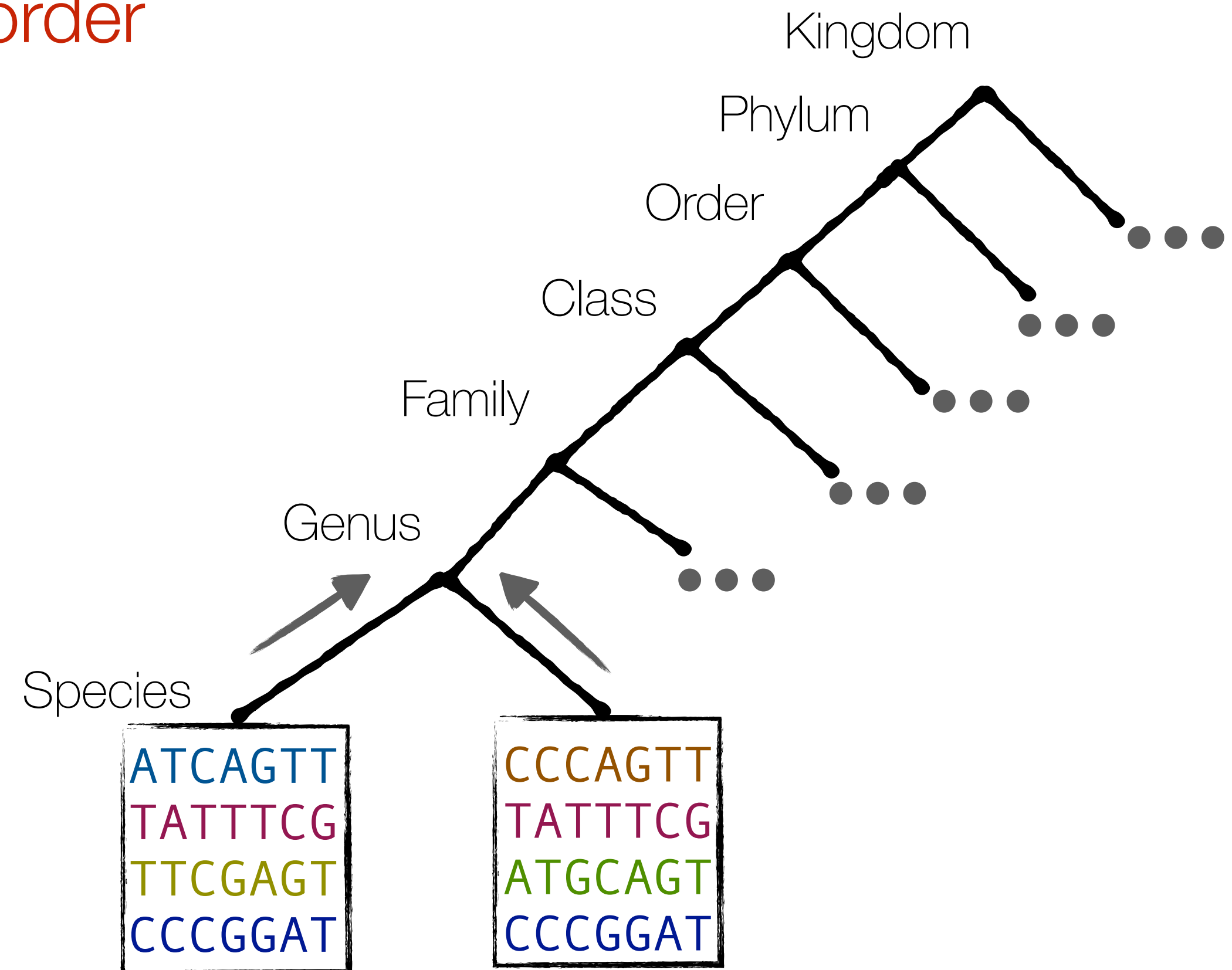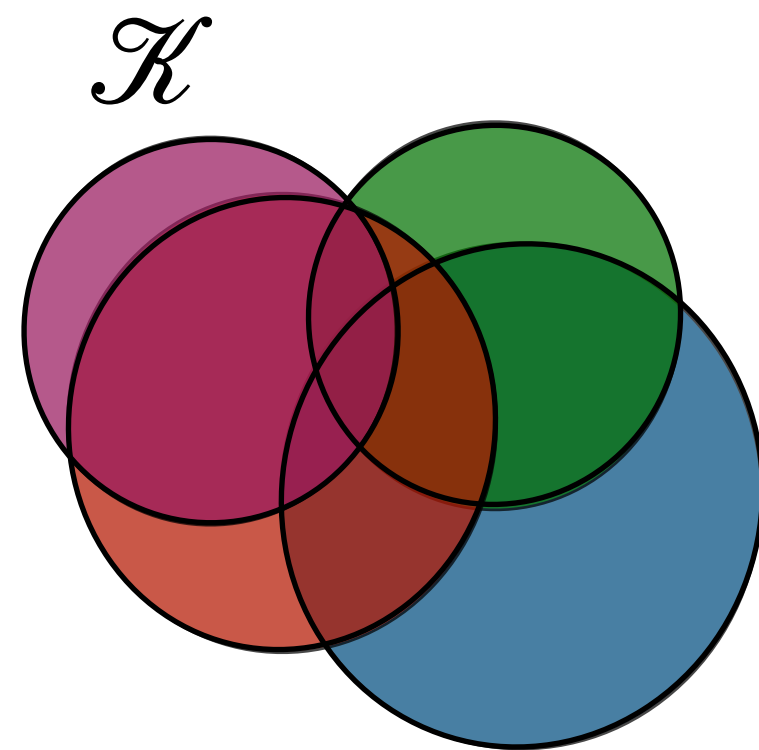           CAAACGT...

# **KRANK** selects a representative *k*-mer subset in a memory-bound manner!

**Core idea:** instead of computing all intersections; hierarchical subsampling through a post order traversal of the taxonomic tree
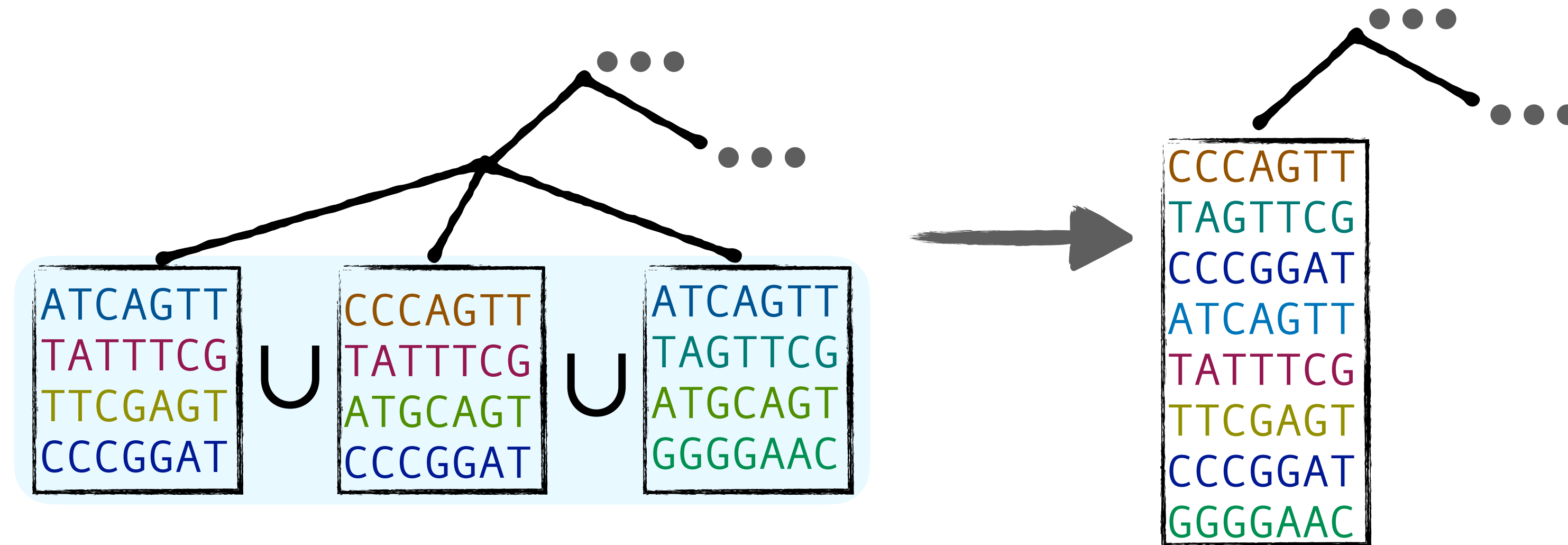
# KRANK selects a representative *k*-mer subset in a memory-bound manner!

**Core idea:** instead of computing all intersections; hierarchical subsampling through a post order traversal of the taxonomic tree
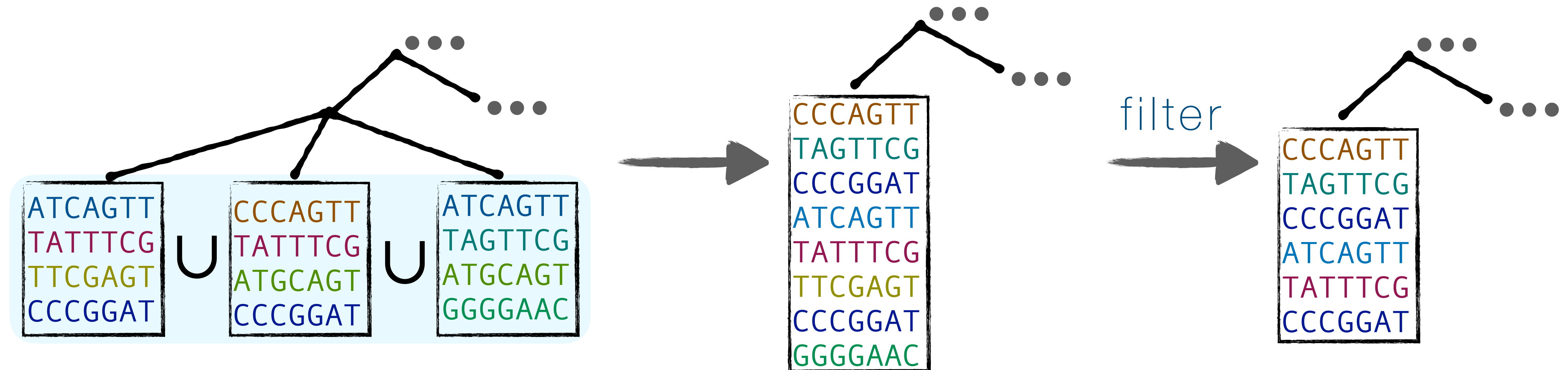
# Gradual filtering of k-mers at internal nodes
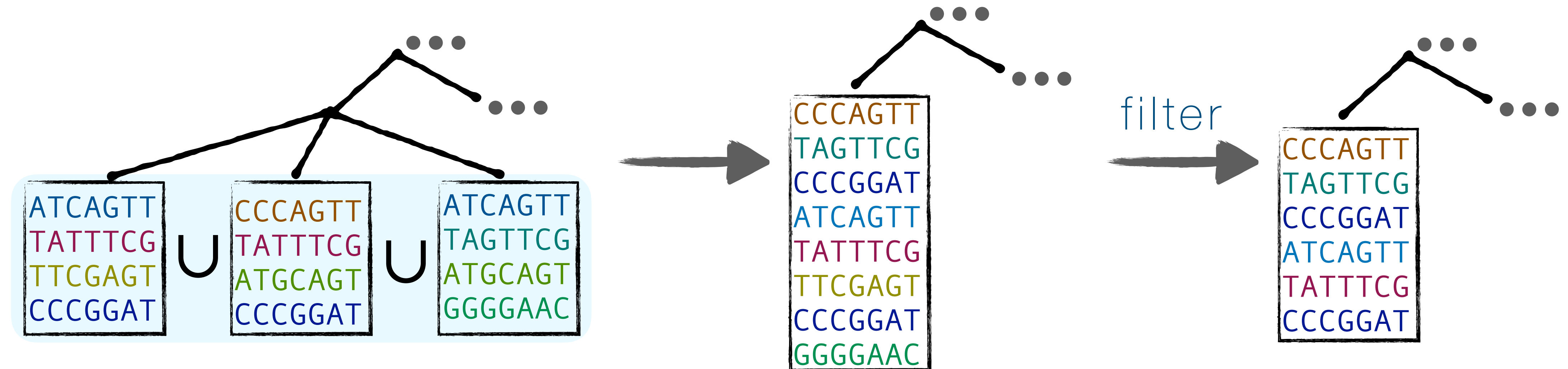
- Recursively take the union of sibling taxa

# Gradual filtering of k-mers at internal nodes

- Recursively take the union of sibling taxa

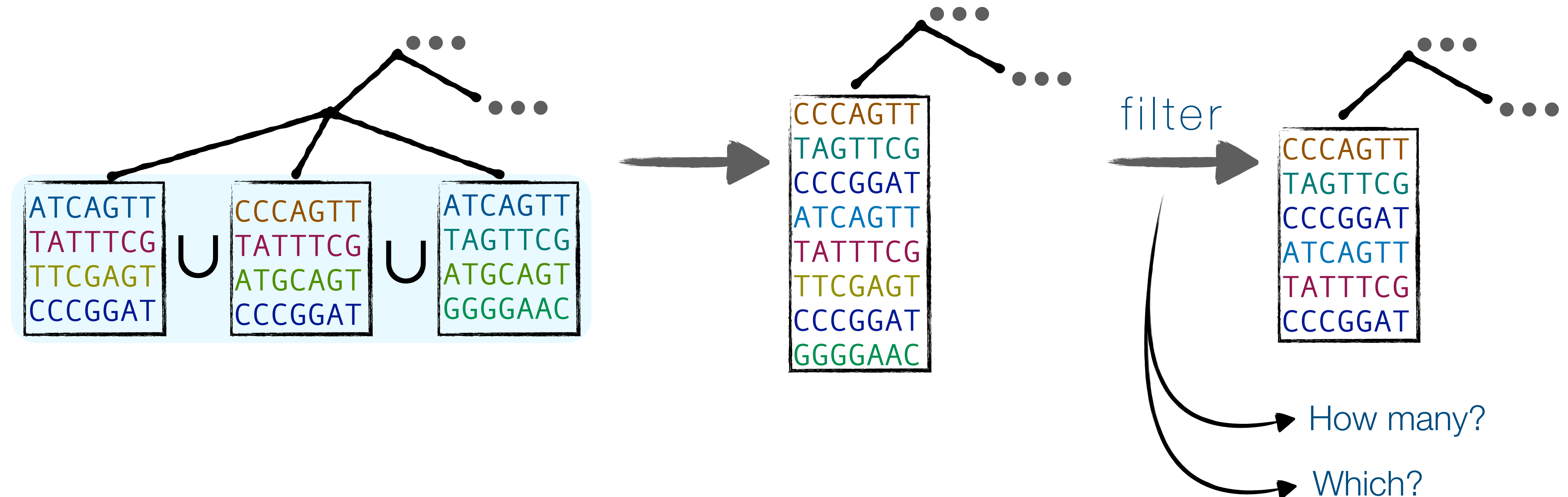- Filter *some number* of *k*-mers based on *a ranking*

# Gradual filtering of k-mers at internal nodes

- Recursively take the union of sibling taxa

- Filter *some number* of *k*-mers based on *a ranking*

- At the root, we obtain the final library with size $M$

# Gradual filtering of k-mers at internal nodes

- Recursively take the union of sibling taxa

- Filter *some number* of *k*-mers based on *a ranking*

- At the root, we obtain the final library with size $M$

**Q1:** How many *k*-mers should we remove from each node/taxon?

**Q2:** How do we rank *k*-mers to assess which one(s) should be kept?

- **Baseline:** no gradual filtering — wait & select $M$ randomly at the root

- **Baseline:** no gradual filtering — wait & select $M$ randomly at the root

Given total budget $M$,

$\mathbb{E}[\text{\# of selected } k\text{-mers for a taxon } t]$ is

$$M \frac{|\mathcal{K}_t|}{|\mathcal{K}|}$$

set of $k$-mers
under the taxon $t$

set of all
reference $k$-mers

- **Baseline:** no gradual filtering — wait & select $M$ randomly at the root

Given total budget $M$,

$\mathbb{E}[\text{\# of selected } k\text{-mers for a taxon } t]$ is
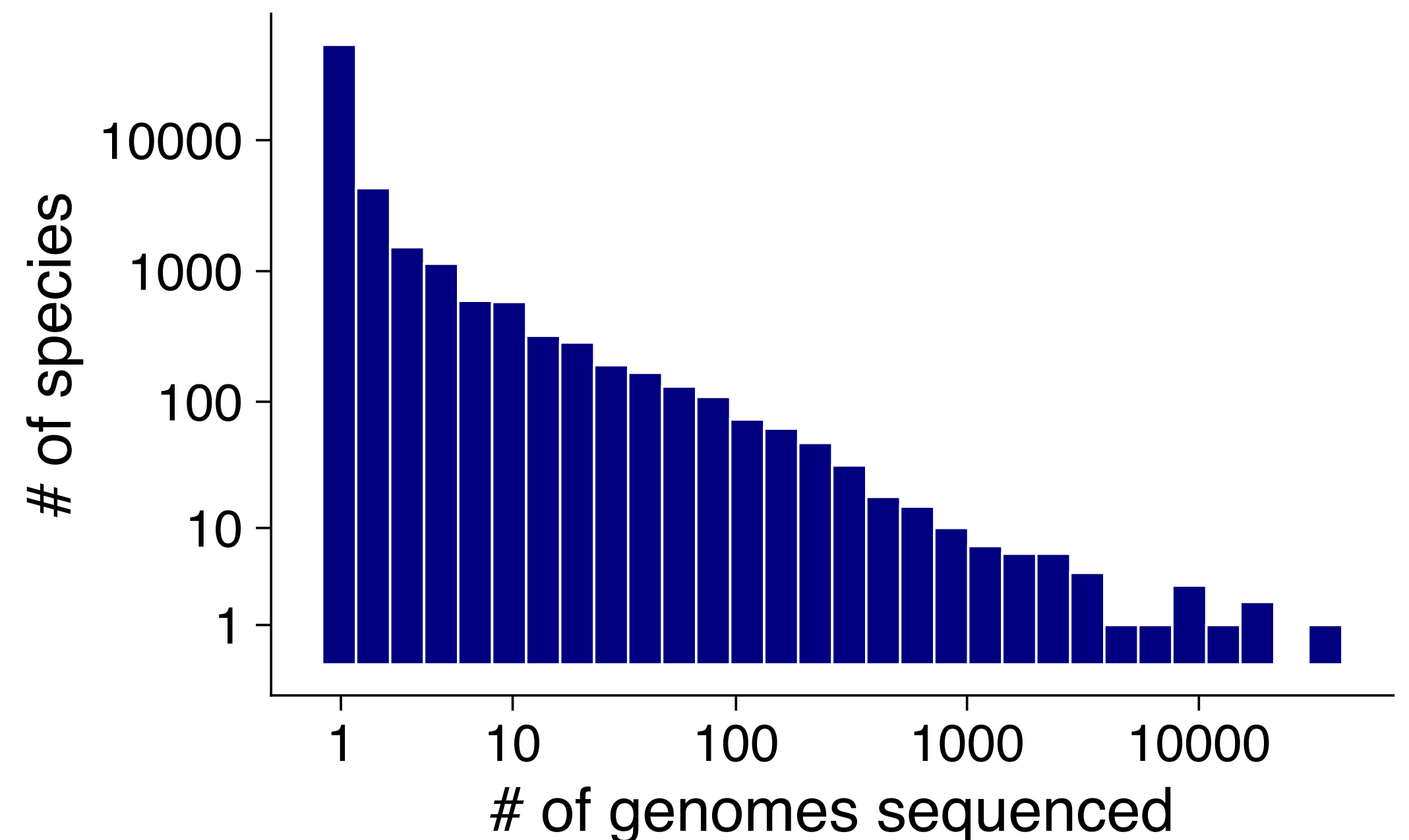
- Proportional contribution $\rightarrow$
  - ▸ taxa with low sampling get little representation
  - ▸ highly-sampled groups dominates (e.g., *E. coli*)

$$M \frac{|\mathscr{K}_t|}{|\mathscr{K}|}$$

set of *k*-mers under the taxon *t*

set of all reference *k*-mers



# of species (y-axis) vs # of genomes sequenced (x-axis)

# Gradual filtering is making some decisions earlier
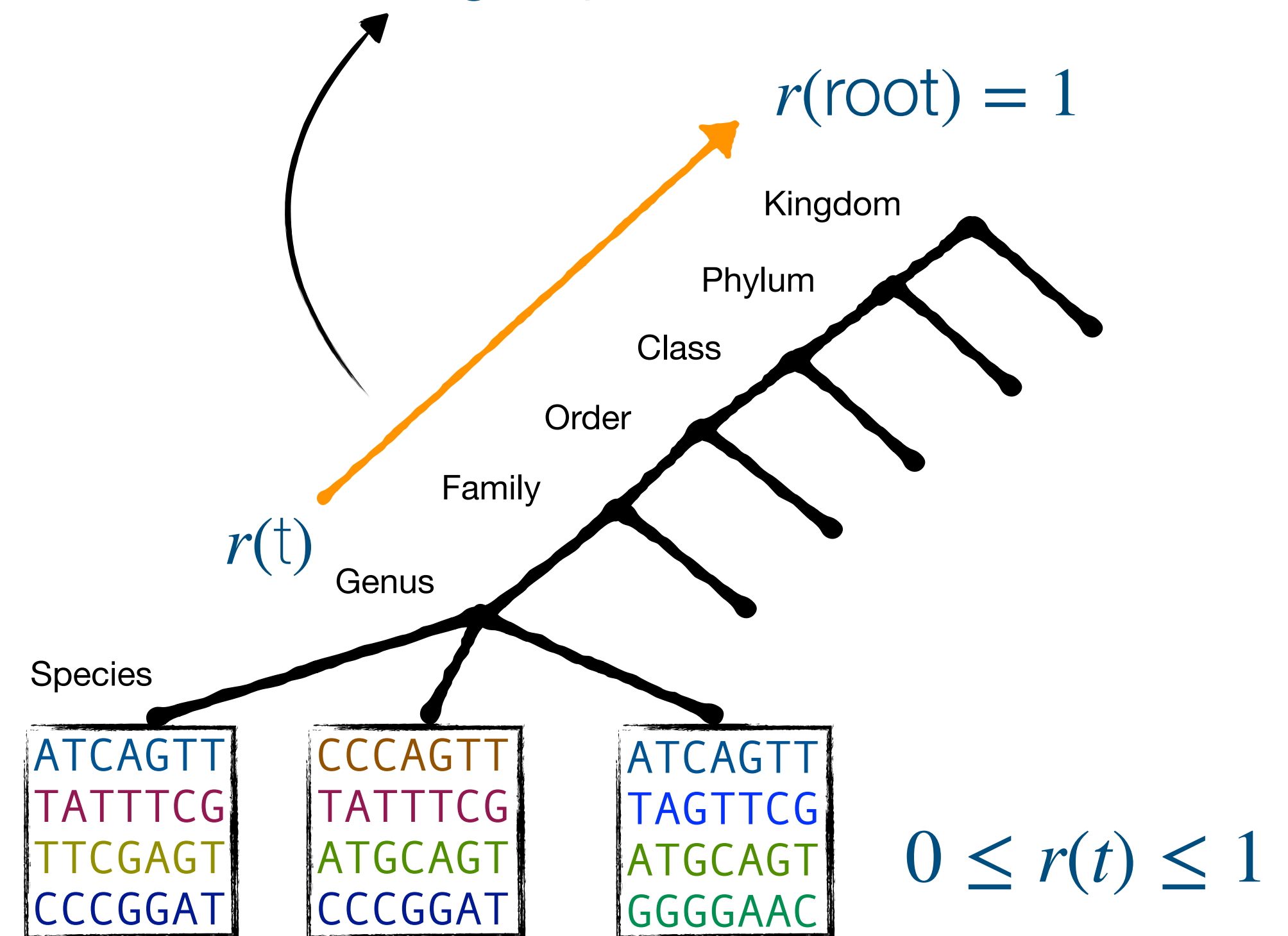
**Goal:** remove *k*-mers from bloated taxa earlier & delay decisions for smaller taxa

# Gradual filtering is making some decisions earlier

**Goal:** remove *k*-mers from bloated taxa earlier & delay decisions for smaller taxa

- Adaptive size constraint, $r(t)M$, on internal nodes

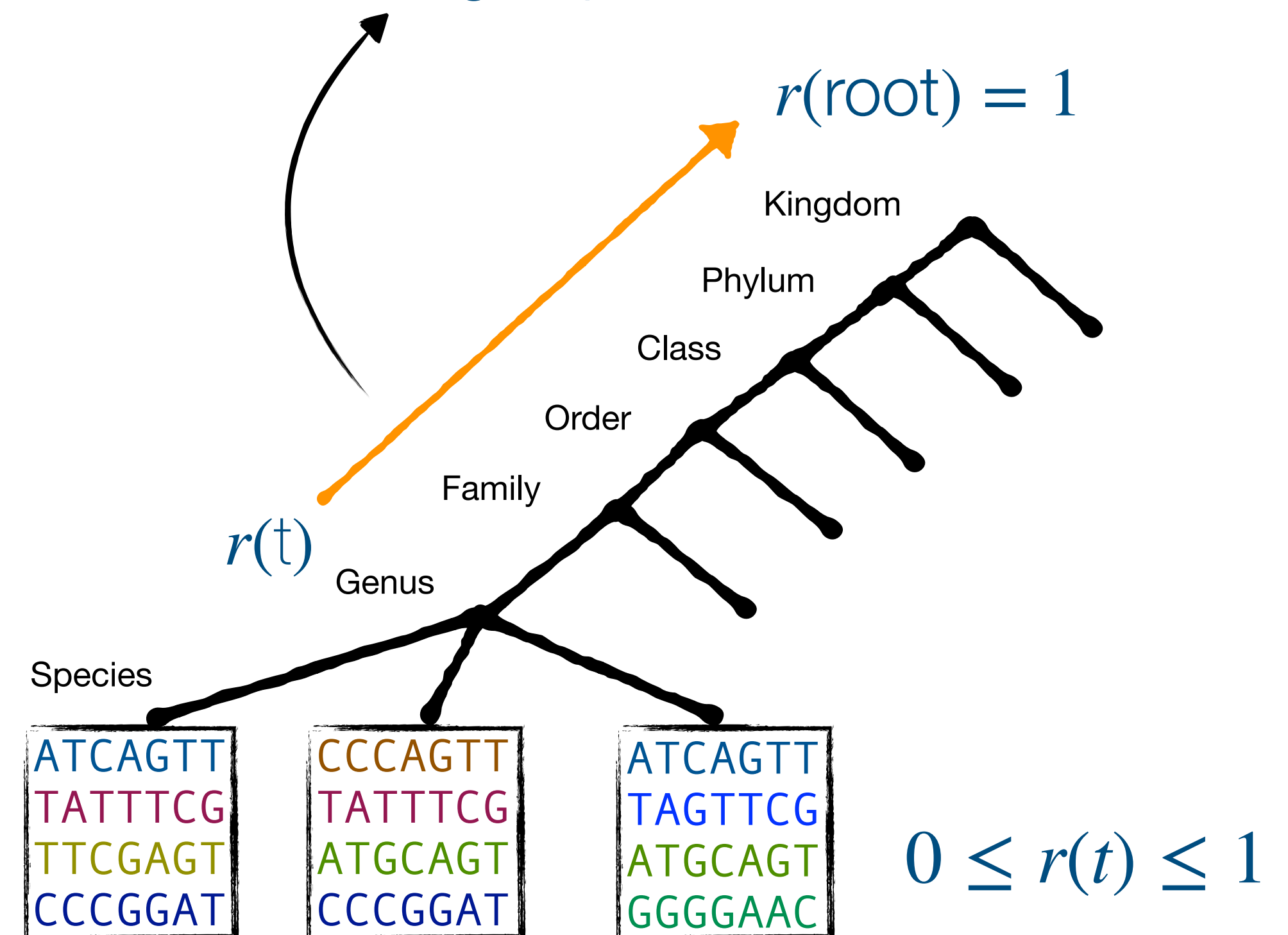increases as we go up in the tree!

$r(\text{root}) = 1$

Kingdom

Phylum

Class

Order

Family

$r(\text{t})$

Genus

Species

| ATCAGTT | CCCAGTT | ATCAGTT |
| TATTTCG | TATTTCG | TAGTTCG |
| TTCGAGT | ATGCAGT | ATGCAGT |
| CCCGGAT | CCCGGAT | GGGGAAC |

$0 \leq r(t) \leq 1$

# Gradual filtering is making some decisions earlier

**Goal:** remove *k*-mers from bloated taxa earlier & delay decisions for smaller taxa

increases as we go up in the tree!

$r(\text{root}) = 1$

- Adaptive size constraint, $r(t)M$, on internal nodes

- $r(t)$ is a heuristic: square root of ratio of *k*-mers under $t$

Kingdom

Phylum

Class

Order

Family

$r(\text{t})$

Genus

Species

| ATCAGTT | CCCAGTT | ATCAGTT |
| TATTTCG | TATTTCG | TAGTTCG |
| TTCGAGT | ATGCAGT | ATGCAGT |
| CCCGGAT | CCCGGAT | GGGGAAC |

$0 \leq r(t) \leq 1$

# Gradual filtering is making some decisions earlier

**Goal:** remove *k*-mers from bloated taxa earlier & delay decisions for smaller taxa

- Adaptive size constraint, $r(t)M$, on internal nodes

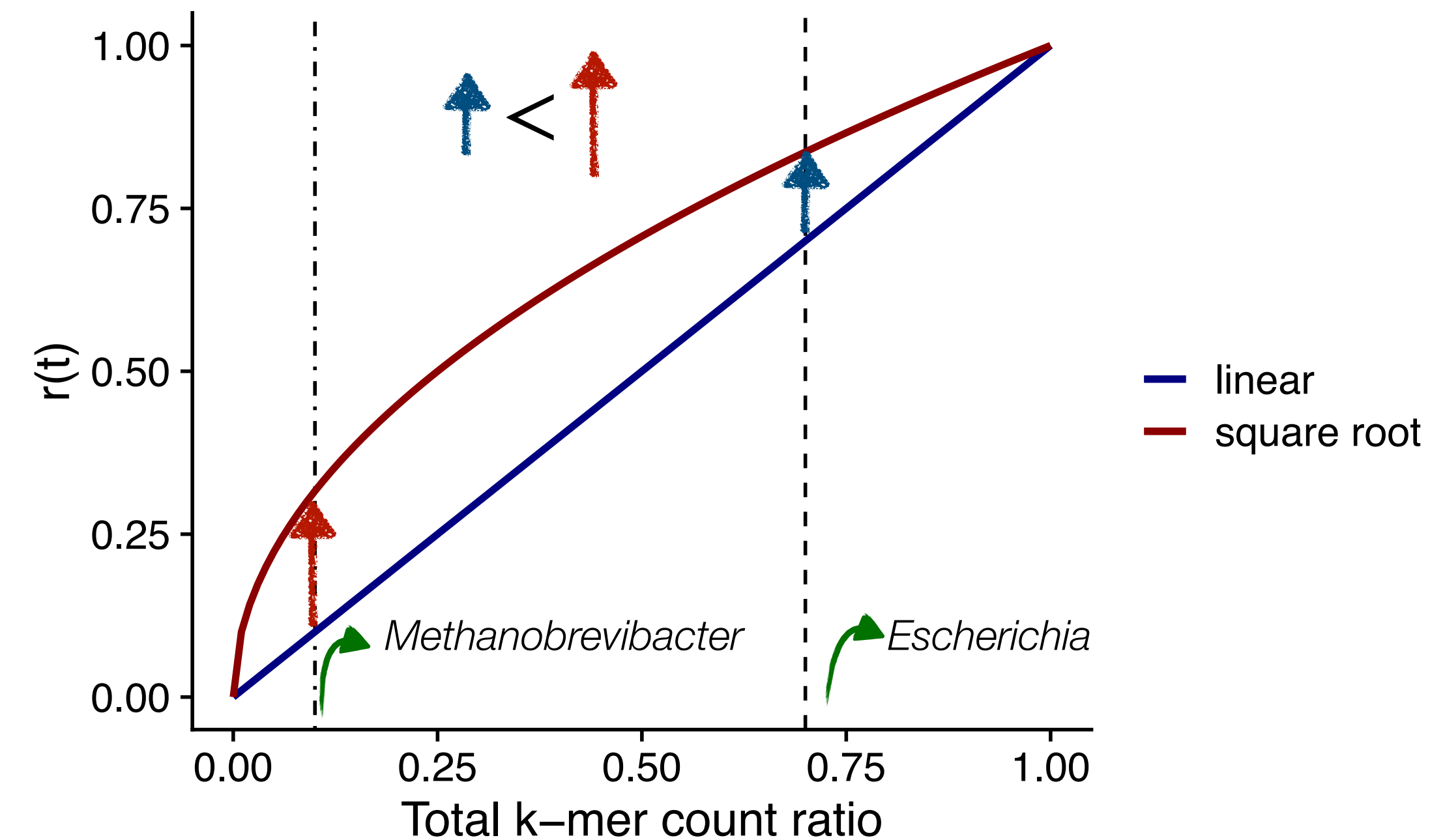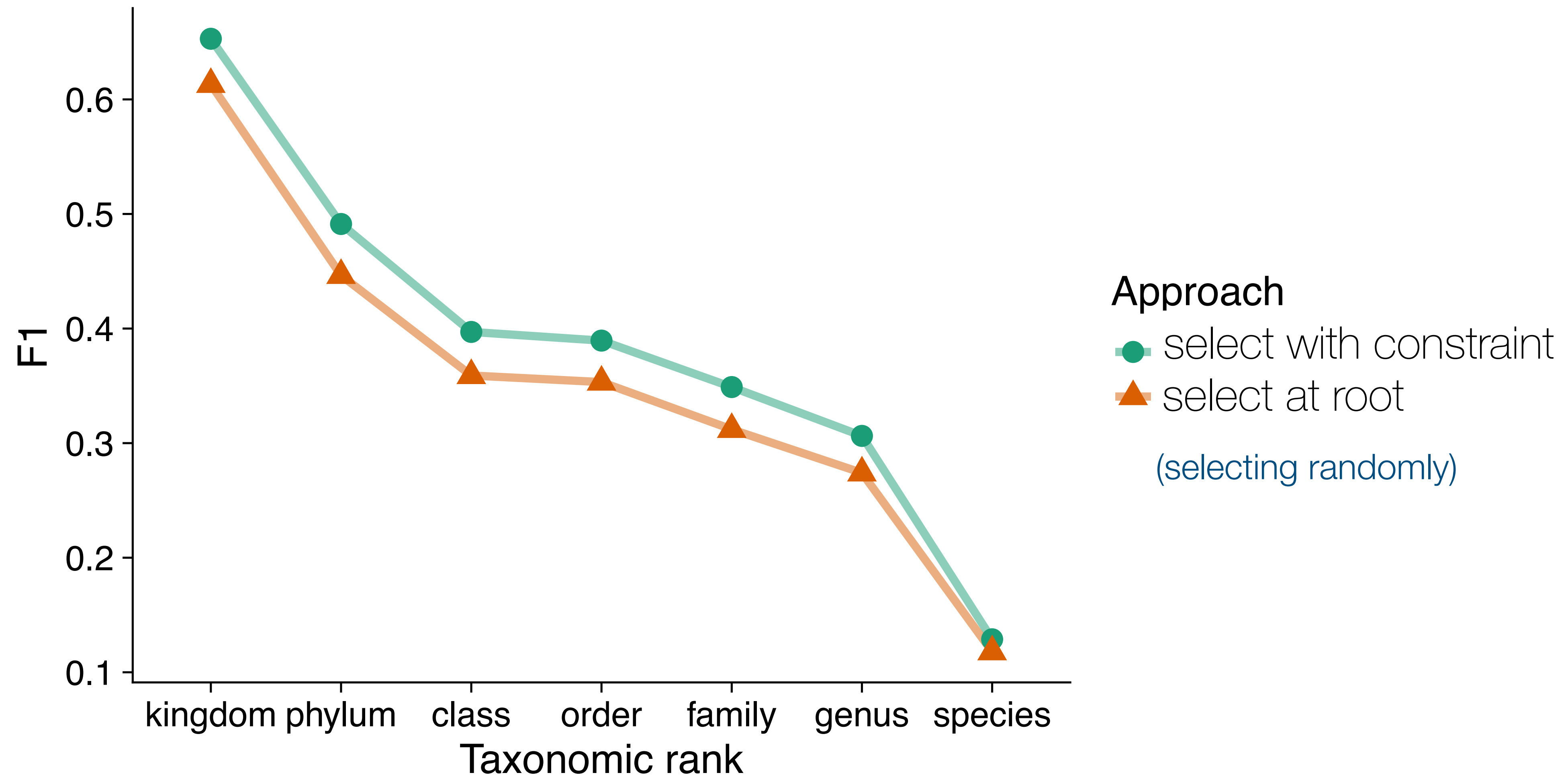- $r(t)$ is a heuristic: square root of ratio of *k*-mers under $t$

- Concavity of $r(t)$ favors taxa with fewer k-mers (less diversity or sparsely sampled)

# Adaptive size constraint improves classification



(empirical analysis using 3.2Gb, in WoL-v1 with 9k species, 10k genomes)

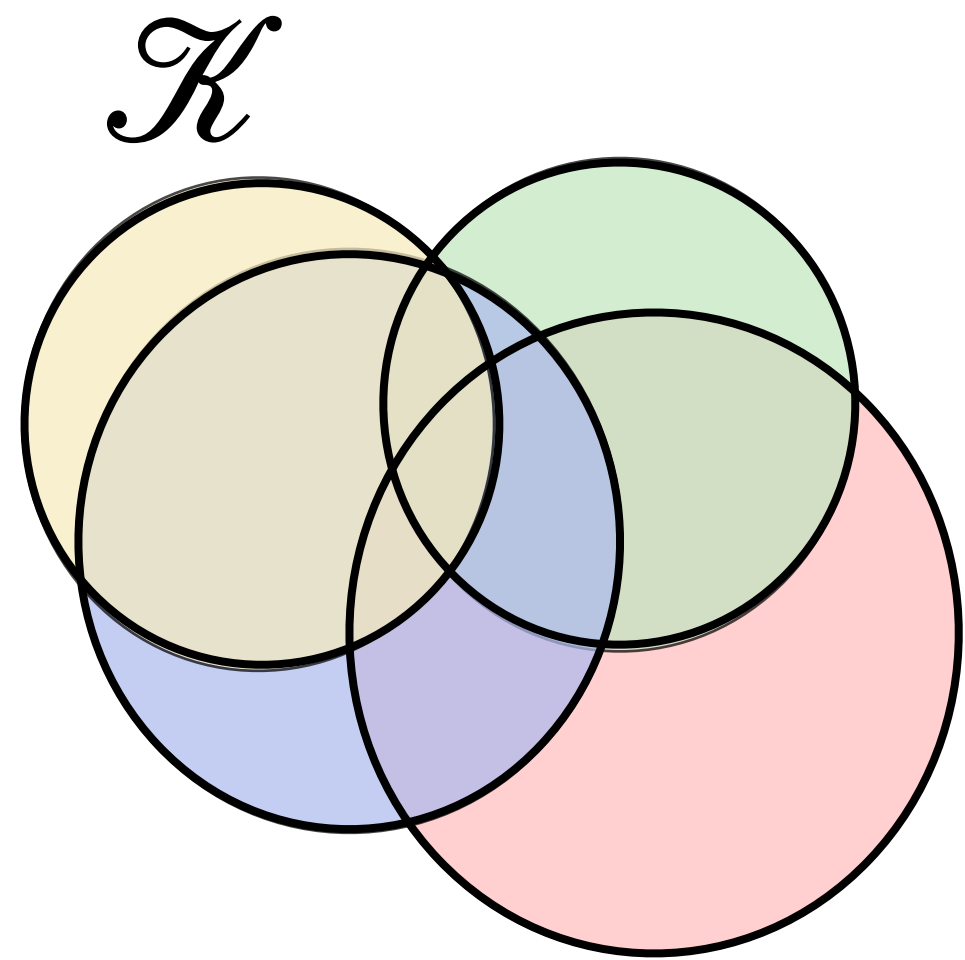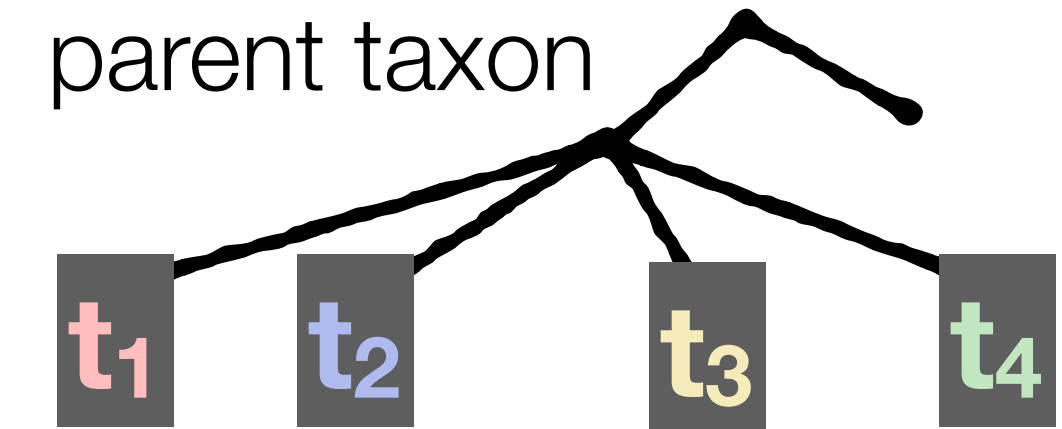**Q1:** How many *k*-mers should we remove from each node/taxon?

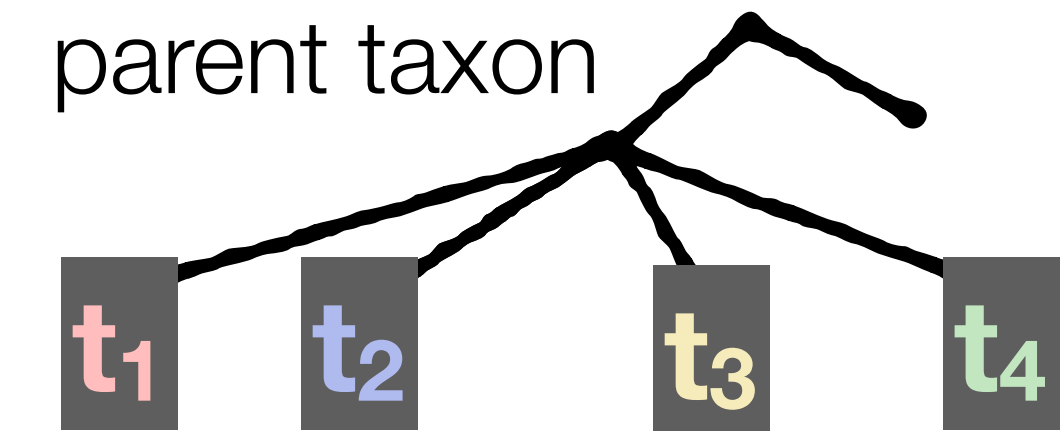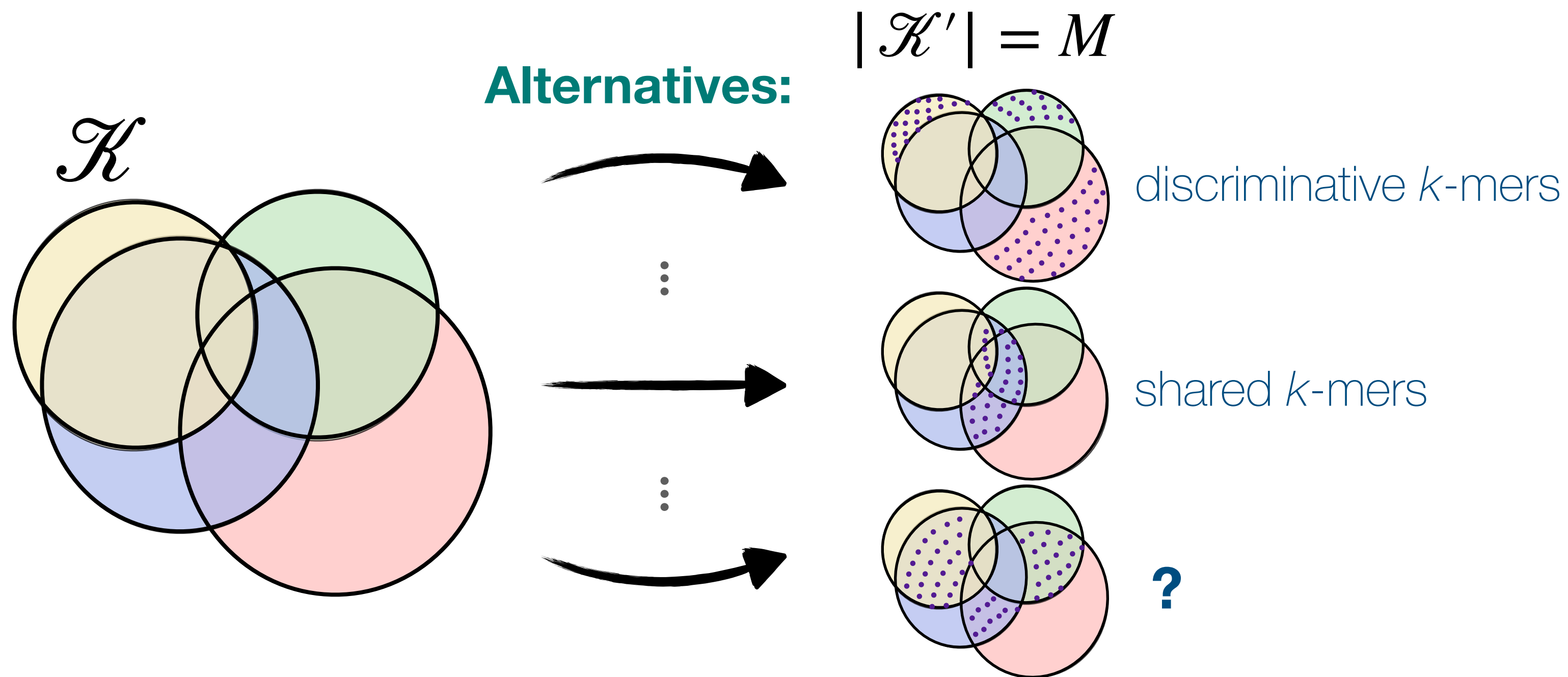**Q2:** How do we rank *k*-mers to assess which one(s) should be kept?

# Which k-mers would provide better representation?

**Baseline:** selecting randomly until the constraint is satisfied
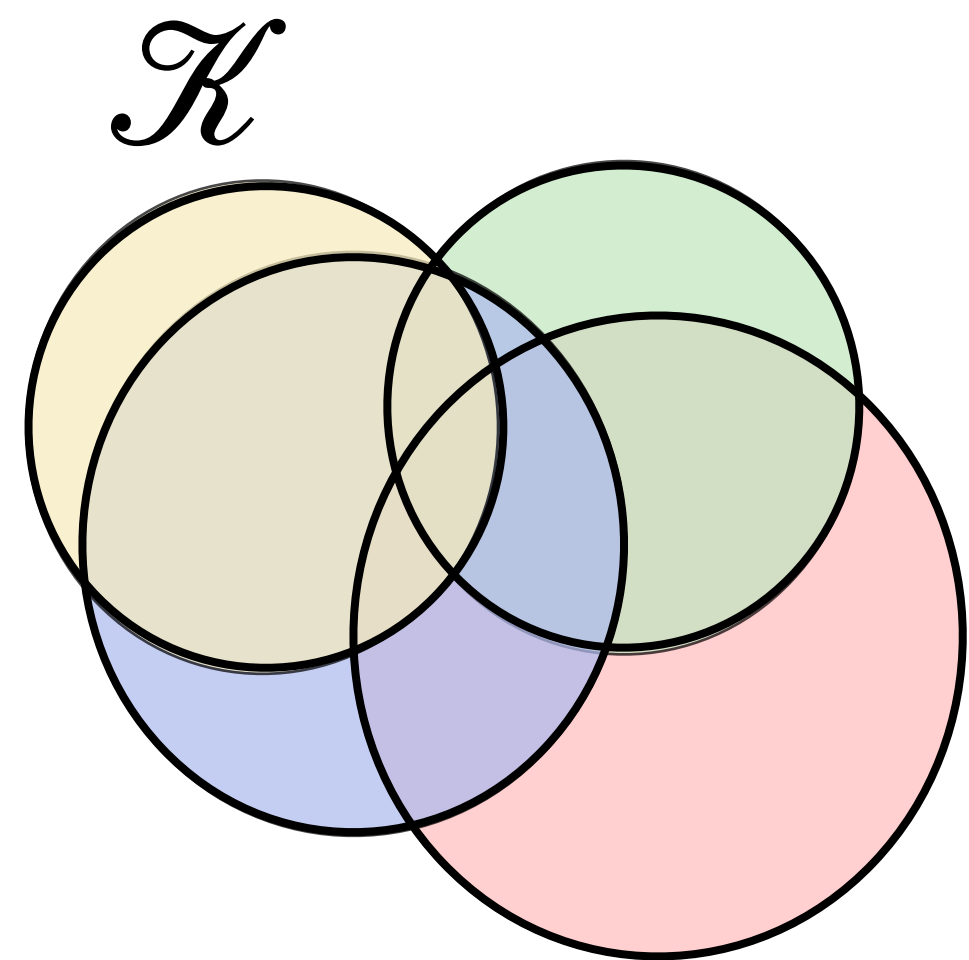
# Which k-mers would provide better representation?

**Baseline:** selecting randomly until the constraint is satisfied

parent taxon

$t_1$  $t_2$  $t_3$  $t_4$

$\mathcal{K}$

# Which k-mers would provide better representation?

**Baseline:** selecting randomly until the constraint is satisfied



parent taxon

$t_1$   $t_2$   $t_3$   $t_4$

$\mathscr{K}$

**Alternatives:**

$|\mathscr{K}'| = M$

discriminative *k*-mers

shared *k*-mers

**?**

# Which k-mers would provide better representation?

**Baseline:** selecting randomly until the constraint is satisfied

parent taxon

$t_1$   $t_2$   $t_3$   $t_4$

$\mathscr{K}$

**Alternatives:**

$|\mathscr{K}'| = M$

discriminative *k*-mers

shared *k*-mers

**?**

# of species under t with k-mer x

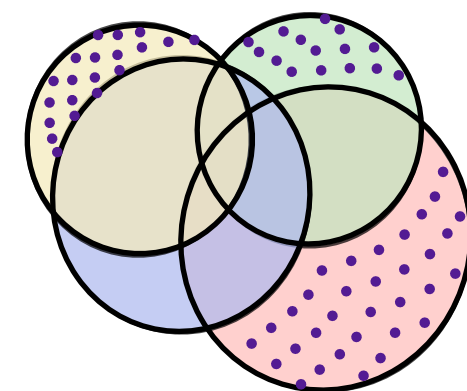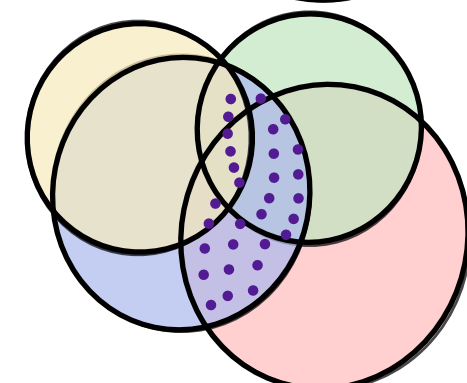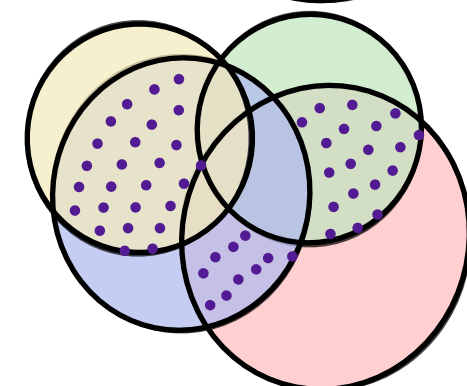|       | $x_1$ | $x_2$ | $x_3$ | ... | $x_{|\mathscr{K}'|}$ |
|-------|-------|-------|-------|-----|------|
| $t_1$ | 4     | 7     | 0     | ... | 3    |
| $t_2$ | 0     | 0     | 2     | ... | 0    |
| $t_3$ | 0     | 0     | 1     | ... | 1    |
| $t_4$ | 2     | 2     | 1     | ... | 0    |
| **Score:** | **6** | **9** | **4** | **...** | **4** |

# Which k-mers would provide better representation?

**Baseline:** selecting randomly until the constraint is satisfied
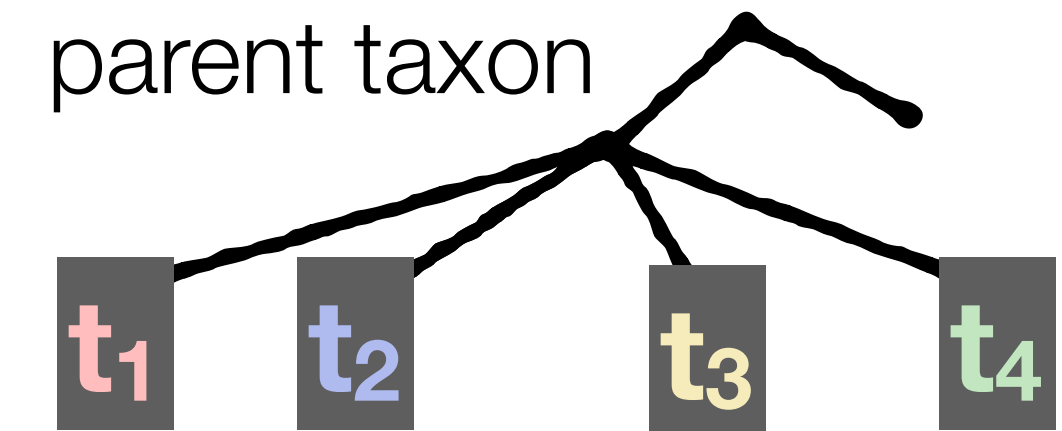
parent taxon



$$|\mathcal{K}'| = M$$

**Alternatives:**

$\mathcal{K}$

discriminative *k*-mers
low scores

shared *k*-mers
high scores

**?**

# of species under t with k-mer x

| | $x_1$ | $x_2$ | $x_3$ | ... | $x_{|\mathcal{K}'|}$ |
|---|---|---|---|---|---|
| $t_1$ | 4 | 7 | 0 | ... | 3 |
| $t_2$ | 0 | 0 | 2 | ... | 0 |
| $t_3$ | 0 | 0 | 1 | ... | 1 |
| $t_4$ | 2 | 2 | 1 | ... | 0 |
| **Score:** | **6** | **9** | **4** | **...** | **4** |

# Neither discriminative nor shared k-mers improve the baseline



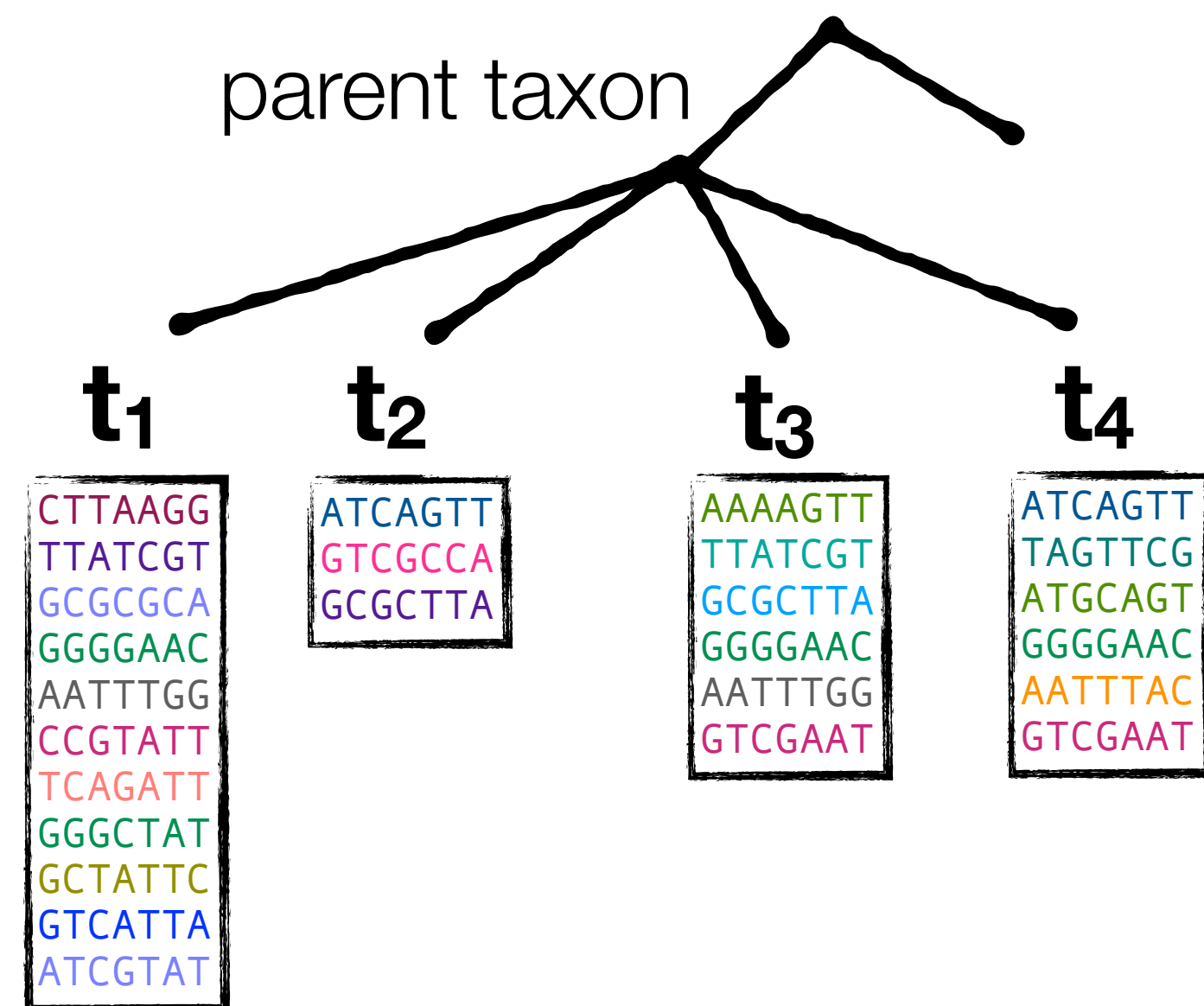(empirical analysis using 3.2Gb, in WoL-v1 with 9k species, 10k genomes)

# Incorporating taxon coverage in ranking

**Intuition:** keep shared *k*-mers but ensure no group is left uncovered



parent taxon

**t₁** **t₂** **t₃** **t₄**

**t₁:** Afford to remove more!

**t₂:** Needs to be prioritized!

# Incorporating taxon coverage in ranking

**Intuition:** keep shared *k*-mers but ensure no group is left uncovered

**Scalable heuristic:** down-weight the impact of taxa that are highly covered among surviving k-mers
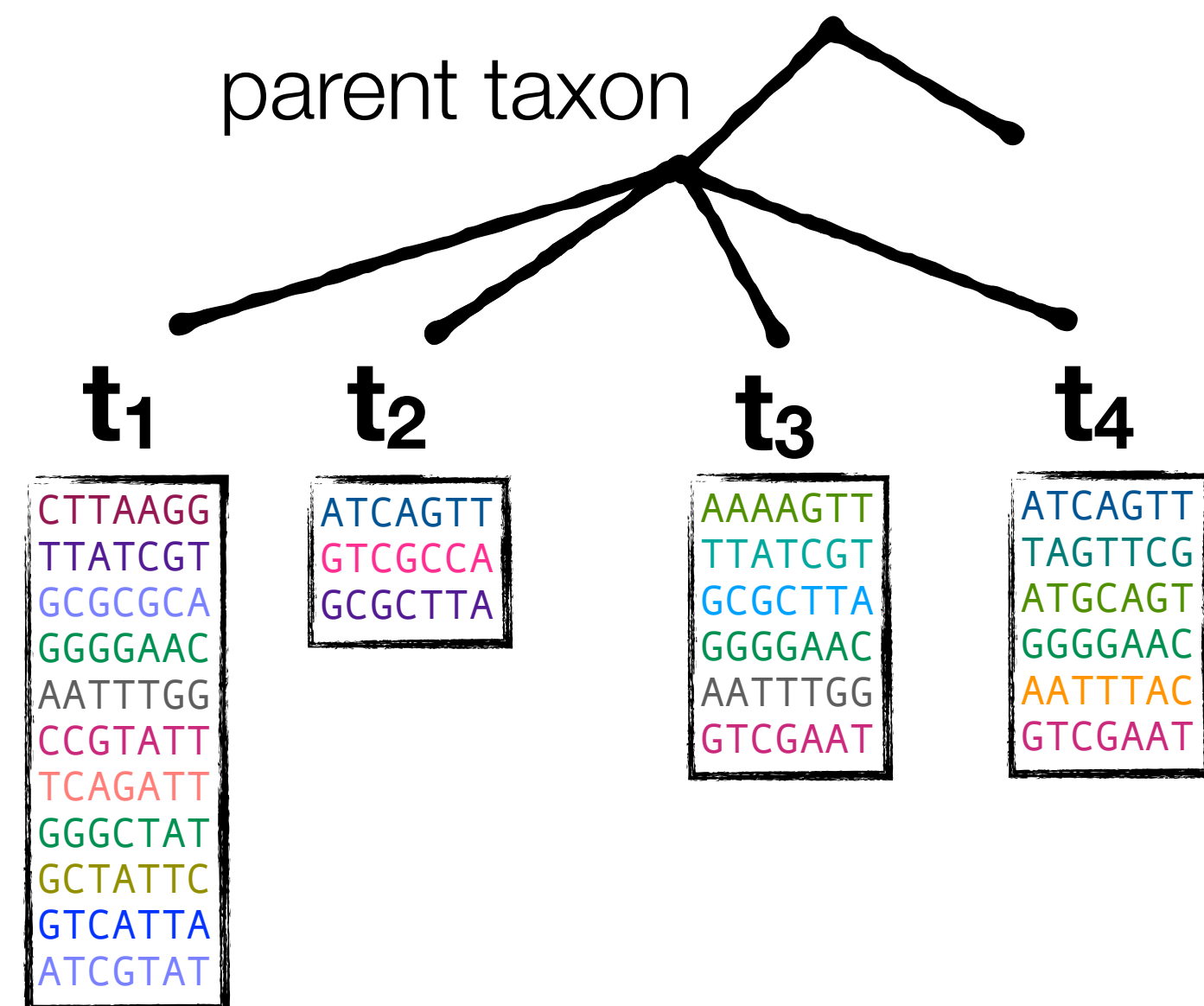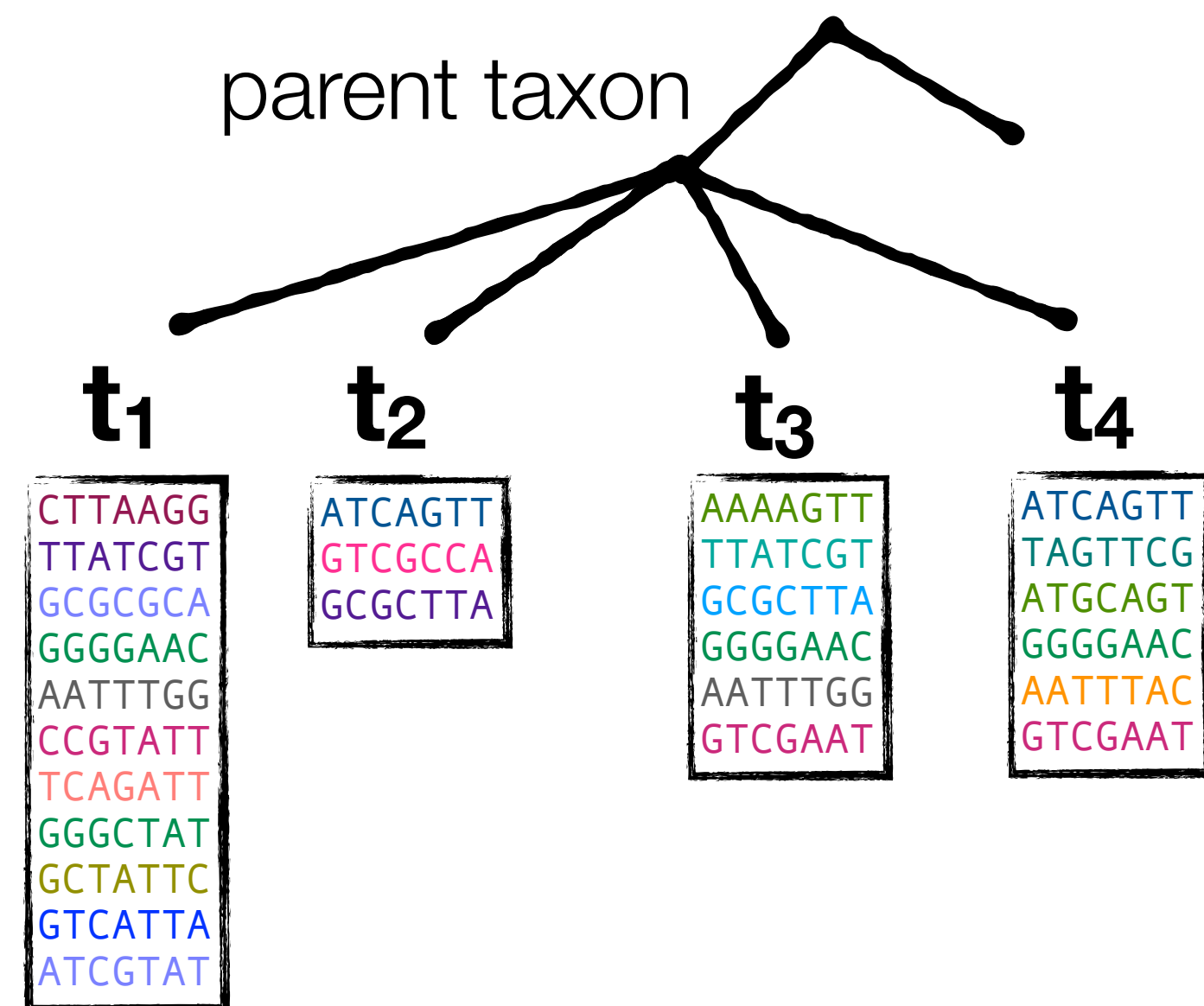


**t₁:** Afford to remove more!

**t₂:** Needs to be prioritized!

# Incorporating taxon coverage in ranking

**Intuition:** keep shared *k*-mers but ensure no group is left uncovered

**Scalable heuristic:** down-weight the impact of taxa that are highly covered among surviving k-mers
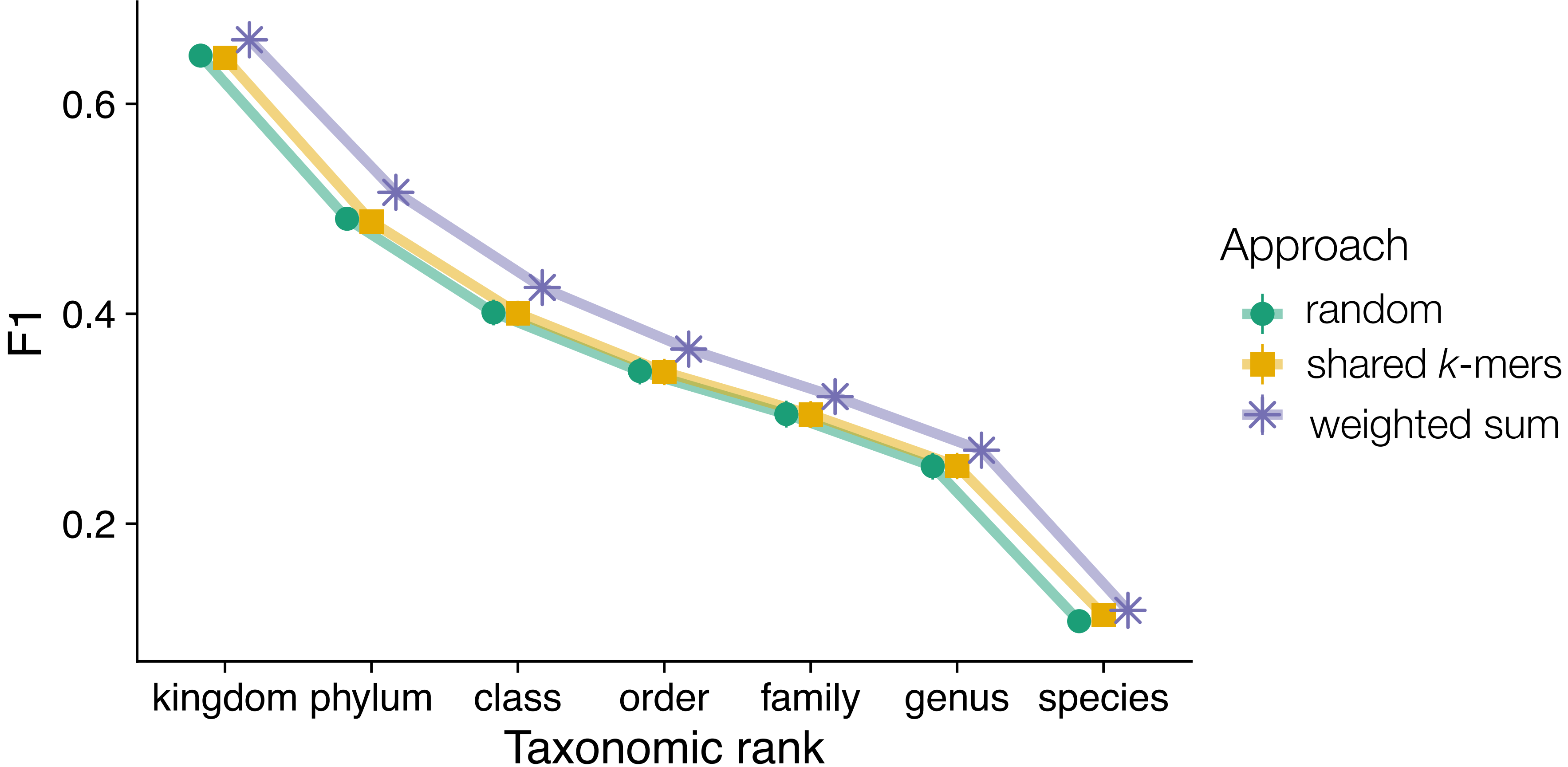


parent taxon

$t_1$ $t_2$ $t_3$ $t_4$

| $t_1$ | $t_2$ | $t_3$ | $t_4$ |
|---|---|---|---|
| CTTAAGG | ATCAGTT | AAAAGTT | ATCAGTT |
| TTATCGT | GTCGCCA | TTATCGT | TAGTTCG |
| GCGCGCA | GCGCTTA | GCGCTTA | ATGCAGT |
| GGGGAAC | | GGGGAAC | GGGGAAC |
| AATTTGG | | AATTTGG | AATTTAC |
| CCGTATT | | GTCGAAT | GTCGAAT |
| TCAGATT | | | |
| GGGCTAT | | | |
| GCTATTC | | | |
| GTCATTA | | | |
| ATCGTAT | | | |

# of species under t with k-mer x

weights of taxa

| | | $x_1$ | $x_2$ | $x_3$ | ... | $x_{|\mathcal{K}'|}$ |
|---|---|---|---|---|---|---|
| **0.09** | $t_1$ | 4 | 7 | 0 | ... | 3 |
| **0.33** | $t_2$ | 0 | 0 | 2 | ... | 0 |
| **0.17** | $t_3$ | 0 | 0 | 1 | ... | 1 |
| **0.17** | $t_4$ | 2 | 2 | 1 | ... | 0 |
| **Score:** | | 0.7 | 0.97 | 1 | ... | 0.44 |

•

# Neither discriminative nor shared k-mers improve the baseline
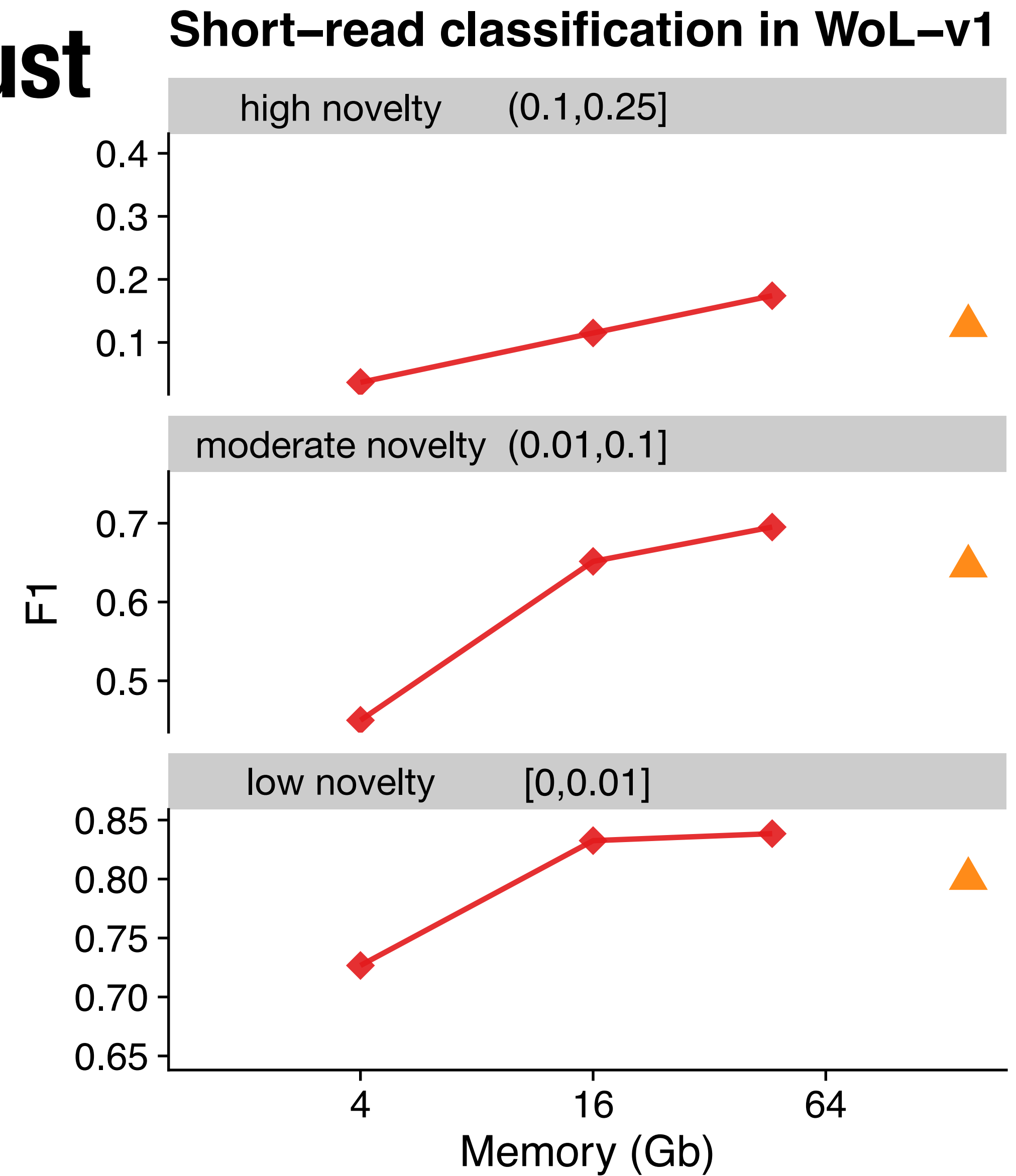


(empirical analysis using 3.2Gb, in WoL-v1 with 9k species, 10k genomes)

- **<u>KRANK</u>** puts all these heuristics together:
  - ‣ weighted-sum ranking + adaptive size constraint
  - ‣ other minor tricks
  - ‣ highly optimized and scalable implementation

# KRANK builds lightweight and robust reference libraries
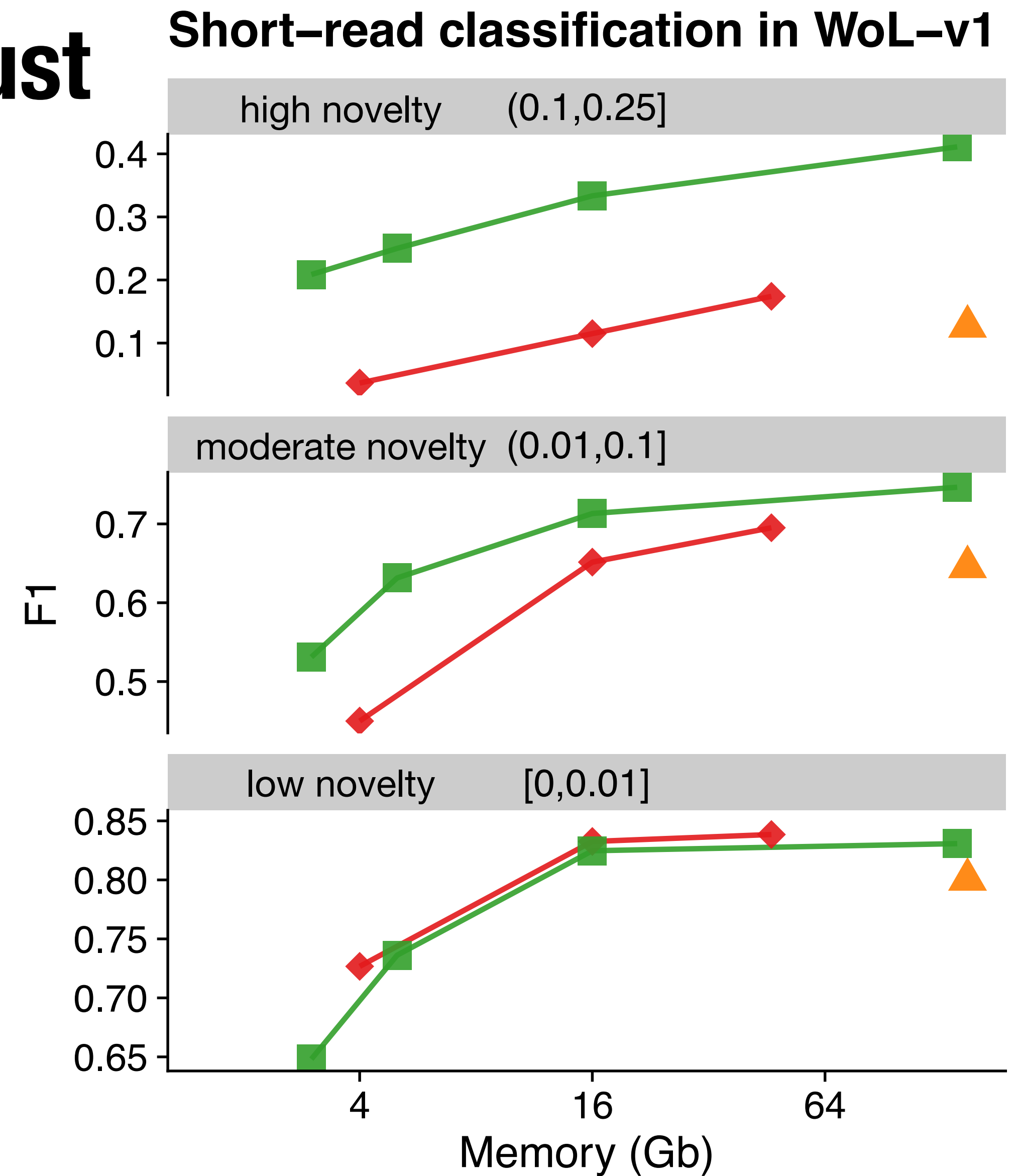
- Simulated reads across different novelty levels



**Short–read classification in WoL–v1**

10k genomes
9k species

Kraken–II v2.1.3    CLARK v1.2.6.1
CONSULT–II v0.4.0    KRANK v0.3.2

# KRANK builds lightweight and robust reference libraries

- Simulated reads across different novelty levels

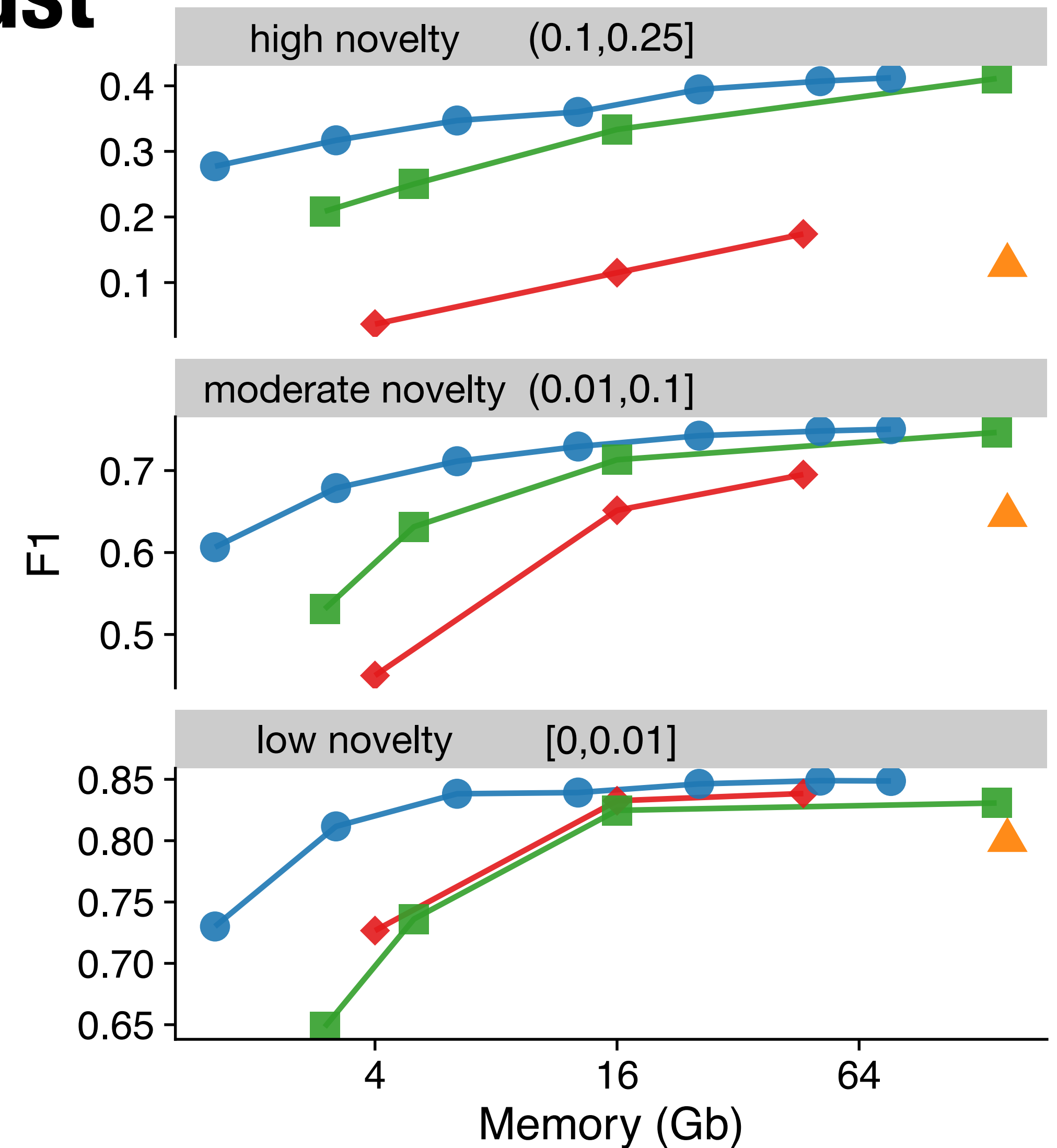- Adjusting the memory usage and observing the impact on the performance



Short−read classification in WoL−v1

10k genomes
9k species

◆ Kraken−II v2.1.3      ▲ CLARK v1.2.6.1

■ CONSULT−II v0.4.0    ● KRANK v0.3.2

# KRANK builds lightweight and robust reference libraries



**Short−read classification in WoL−v1**

- Simulated reads across different novelty levels

- Adjusting the memory usage and observing the impact on the performance

- KRANK preserves the same level of robust performance with much smaller *k*-mer subsets

10k genomes
9k species

Kraken−II v2.1.3   CLARK v1.2.6.1

CONSULT−II v0.4.0   KRANK v0.3.2

# Boosting the performance in CAMI-II with a smaller subset

- Library construction: 3-hours (36 nodes ✕ 14 cores) for RefSeq genomes (2019)

# Boosting the performance in CAMI-II with a smaller subset

- Library construction: 3-hours (36 nodes ✕ 14 cores) for RefSeq genomes (2019)

- Consistently improves CONSULT-II across all ranks



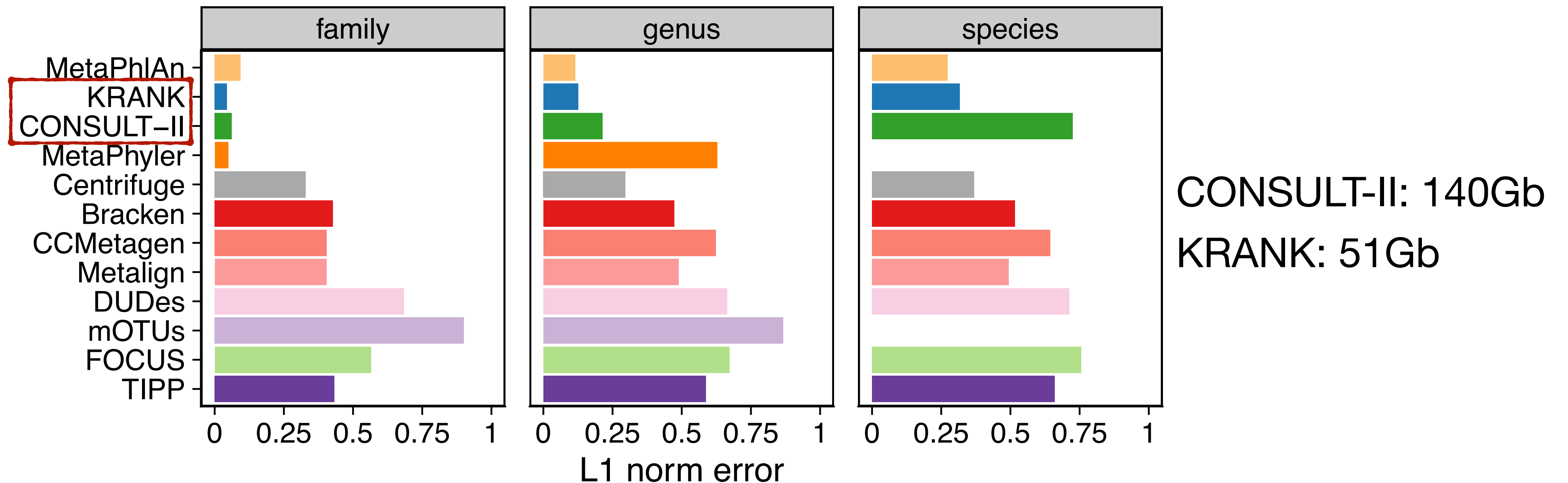**Strain–madness dataset of CAMI–II**

CONSULT-II: 140Gb

KRANK: 51Gb

# Boosting the performance in CAMI-II with a smaller subset

- Library construction: 3-hours (36 nodes ✕ 14 cores) for RefSeq genomes (2019)

- Consistently improves CONSULT-II across all ranks

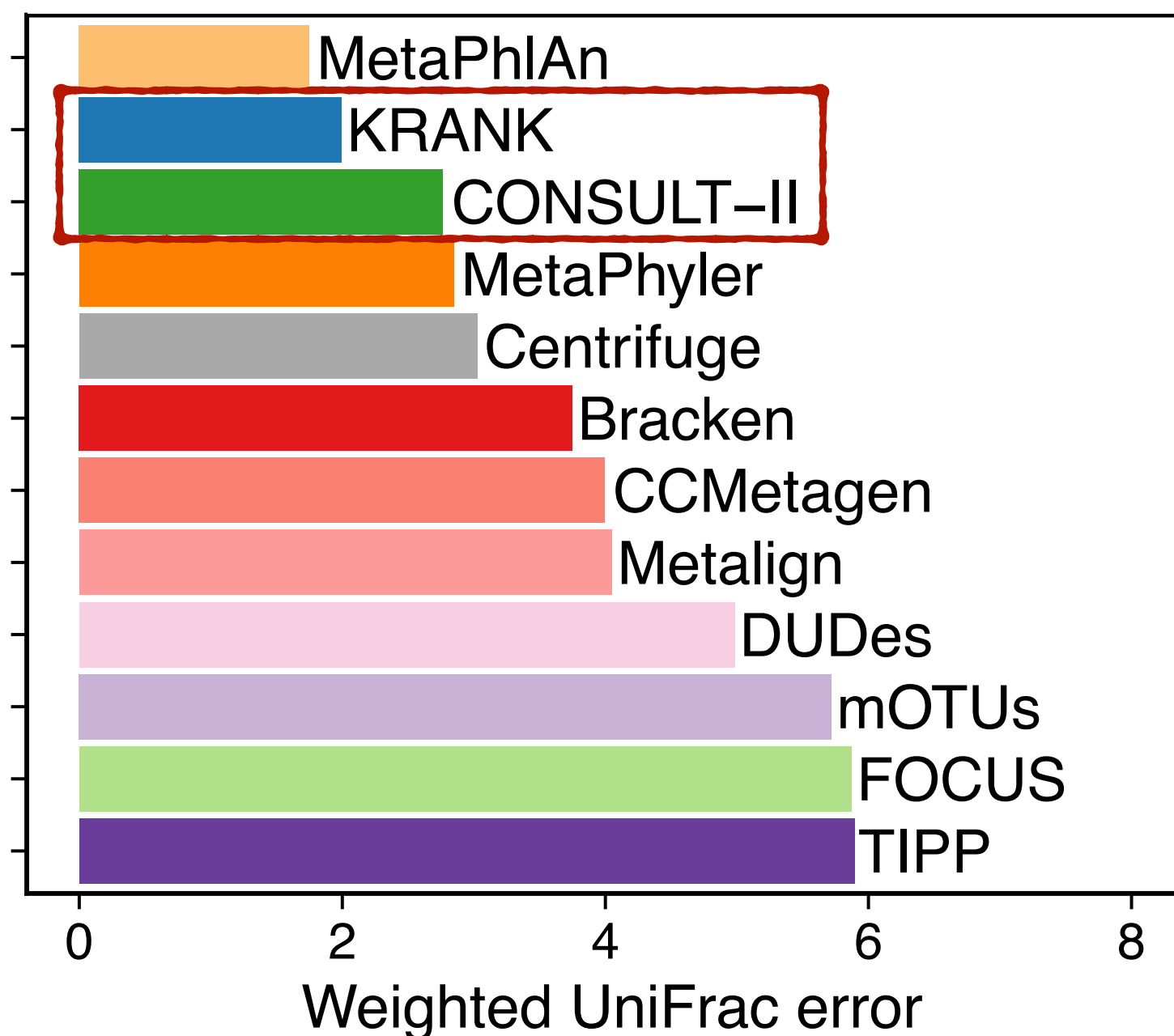- Second-best tool according to rank-invariant UniFrac error

**Strain–madness dataset of CAMI–II**



CONSULT-II: 140Gb

KRANK: 51Gb

- KRANK uses taxonomy to subsample large *k*-mer databases
  - ‣ based on <span style="color:darkred">carefully chosen heuristics</span>
  - ‣ used <span style="color:darkred">in combination with minimizers</span>

- Future work includes:
  - ‣ exploring <span style="color:darkred">alternatives strategy</span> a more <span style="color:darkred">principled approach</span>
    - better modeling of imbalance
    - using a phylogenetic tree
  - ‣ pairing KRANK with other classification methods
  - ‣ pairing with sketching algorithms

- KRANK uses taxonomy to subsample large *k*-mer databases
  - ‣ based on carefully chosen heuristics
  - ‣ used in combination with minimizers

- Future work includes:
  - ‣ exploring alternatives strategy a more principled approach
    - − better modeling of imbalance
    - − using a phylogenetic tree
  - ‣ pairing KRANK with other classification methods
  - ‣ pairing with sketching algorithms

bioRxiv

THE PREPRINT SERVER FOR BIOLOGY

# Extra Slides

# The case against discriminative k-mers

- **Problem:** considerably small portion of *k*-mers are shared within a group! (it gets worse for upper ranks)
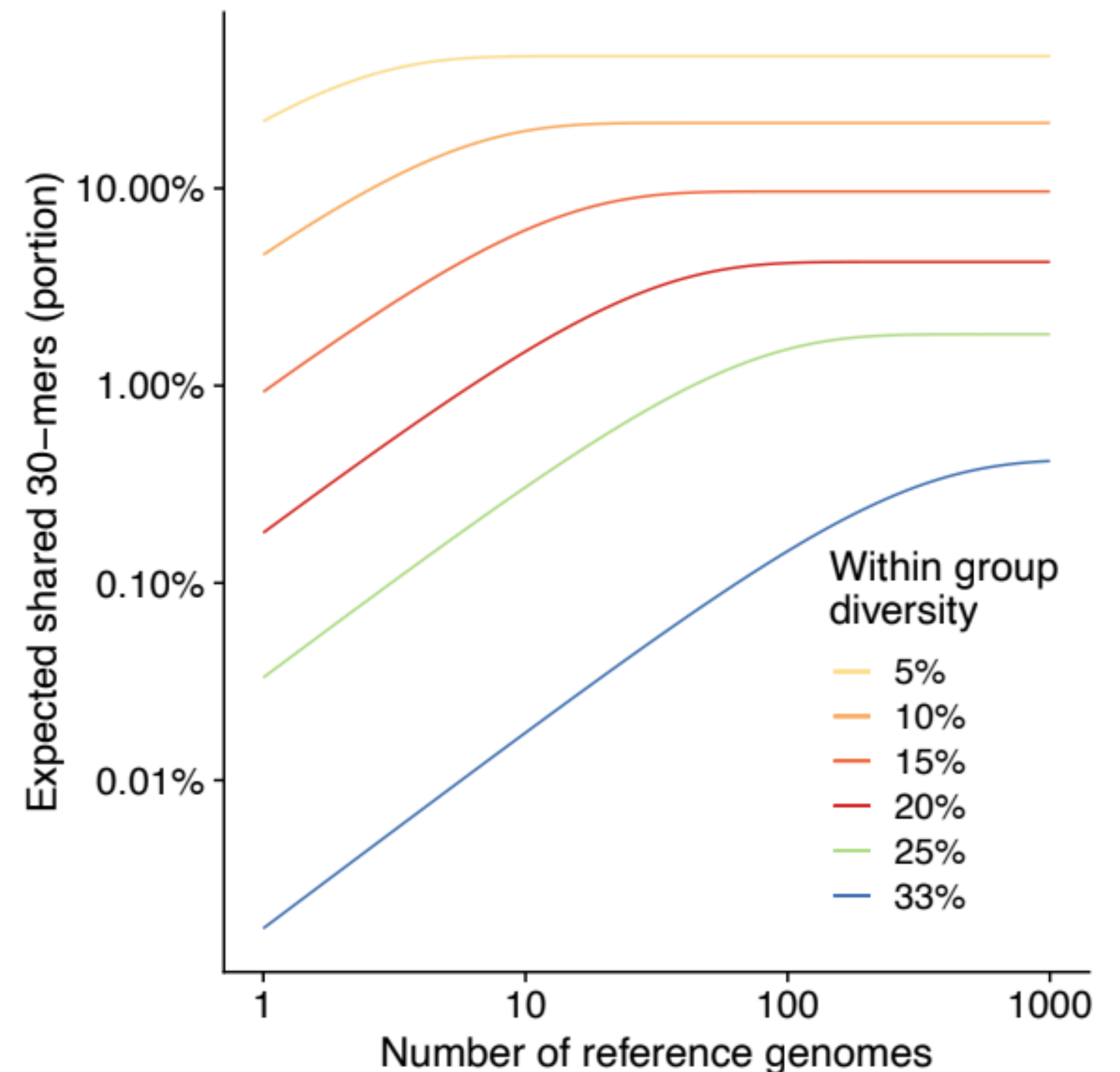
  - ◉ **Claim:** Removing common *k*-mers will make it difficult to find matches!

Given a query genome, what is the expected portion of shared *k*-mers in a reference set with $N$ genomes within $2d$ distance?

$$(1-d)^k\big(1 - (1 - (1-d)^k)^N\big)$$

*k*-mer from the ancestor stays same

*k*-mer from the ancestor changes in all $N$
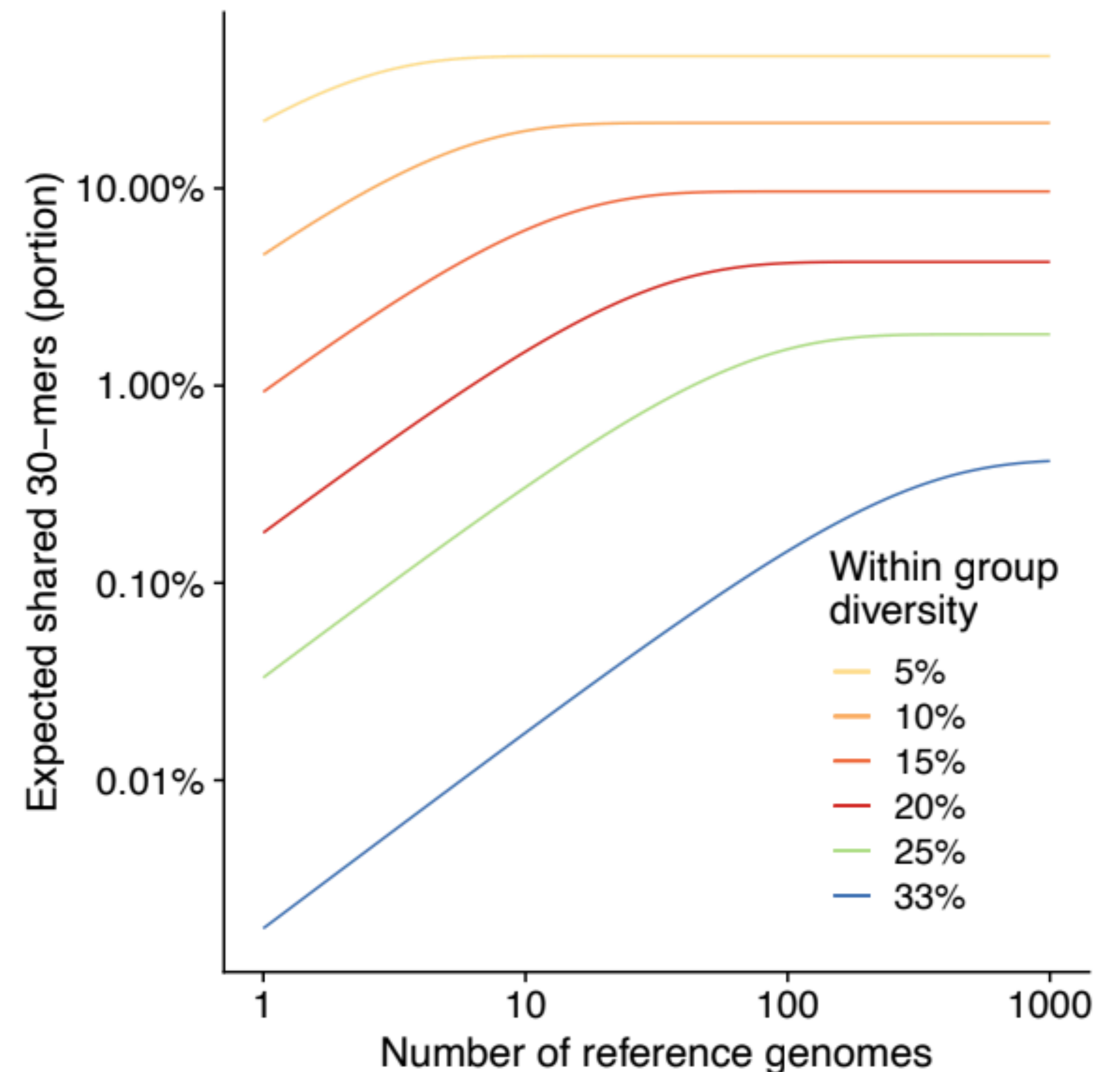
# The case against discriminative k-mers

- **Problem:** considerably small portion of *k*-mers are shared within a group! (it gets worse for upper ranks)

  ⊙ **Claim:** Removing common *k*-mers will make it difficult to find matches!

**Example:** within d = 20% diversity (~genus)

▸ $N = 5$: 0.7% of query 30-mers,

▸ $N \to \infty$: 4.2% of query 30-mers,

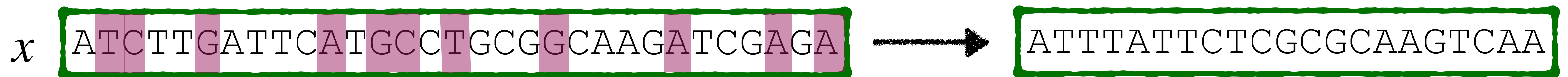will be found in at least one reference.

# Bonus: compact k-mer encodings

CONSULT-II used 2 bits per letter: 64bit for 32-mers.

We only compute HD between $k$-mers that have the same hash value!

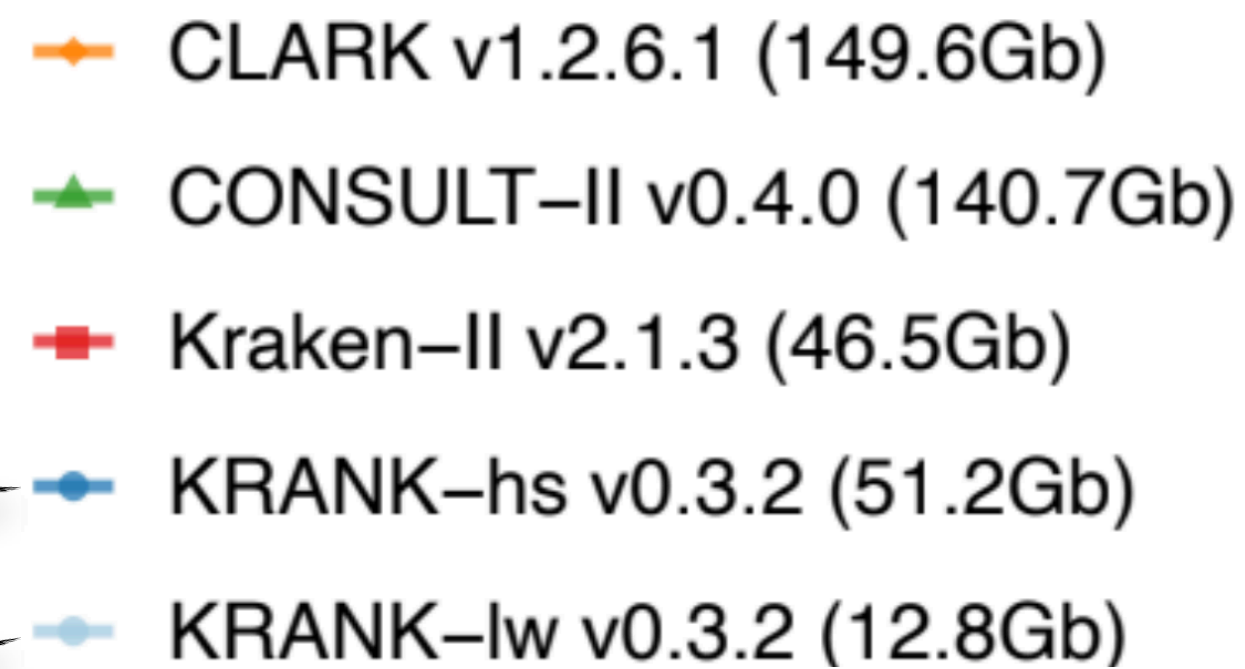We do not need $h$ positions used to compute LSH; they are already the same!

$x$ `ATCTTGATTCATGCCTGCGGCAAGATCGAGA` $\longrightarrow$ `ATTTATTCTCGCGCAAGTCAA`

Just drop LSH positions and store the rest: $k = 32, h = 16 \rightarrow 32$bit

# Improvements are pronounced at higher ranks

- KRANK 13Gb competes with CONSULT-II 144Gb.

- Novel queries were accurately classified at higher ranks.

- With little memory, KRANK+CONSULT-II is highly sensitive.

high-sensitivity memory level

lightweight memory level



Legend:
- CLARK v1.2.6.1 (149.6Gb)
- CONSULT-II v0.4.0 (140.7Gb)
- Kraken-II v2.1.3 (46.5Gb)
- KRANK-hs v0.3.2 (51.2Gb)
- KRANK-lw v0.3.2 (12.8Gb)

SR classification in WoL-v1