# Kraken 2

Improved metagenomic analysis with Kraken 2
Derrick E. Wood, Jennifer Lu & Ben Langmead
Genome Biol 20, 257 (2019)
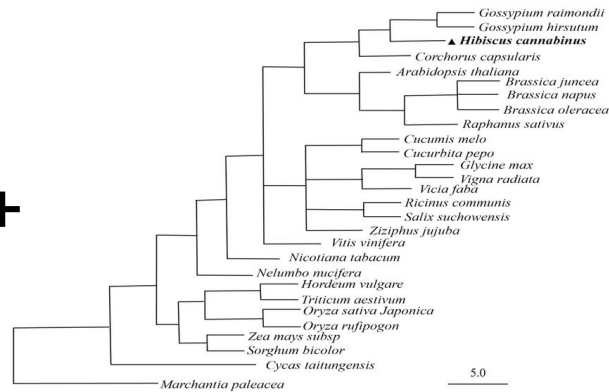https://doi.org/10.1186/s13059-019-1891-0

Konstantin Schütze
Advanced Bioinformatics 1
2021-09-15

# Metagenomic Sequence Classification
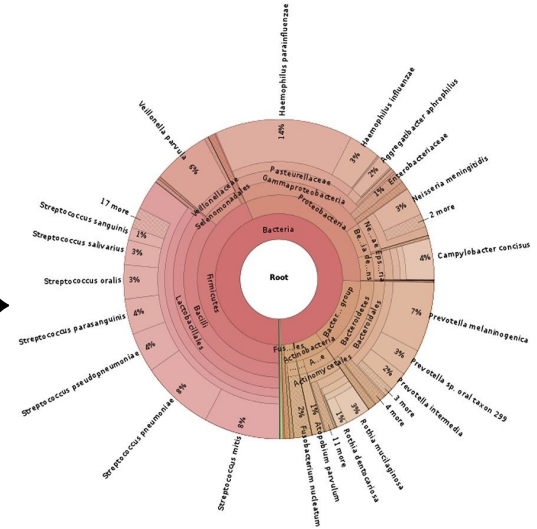
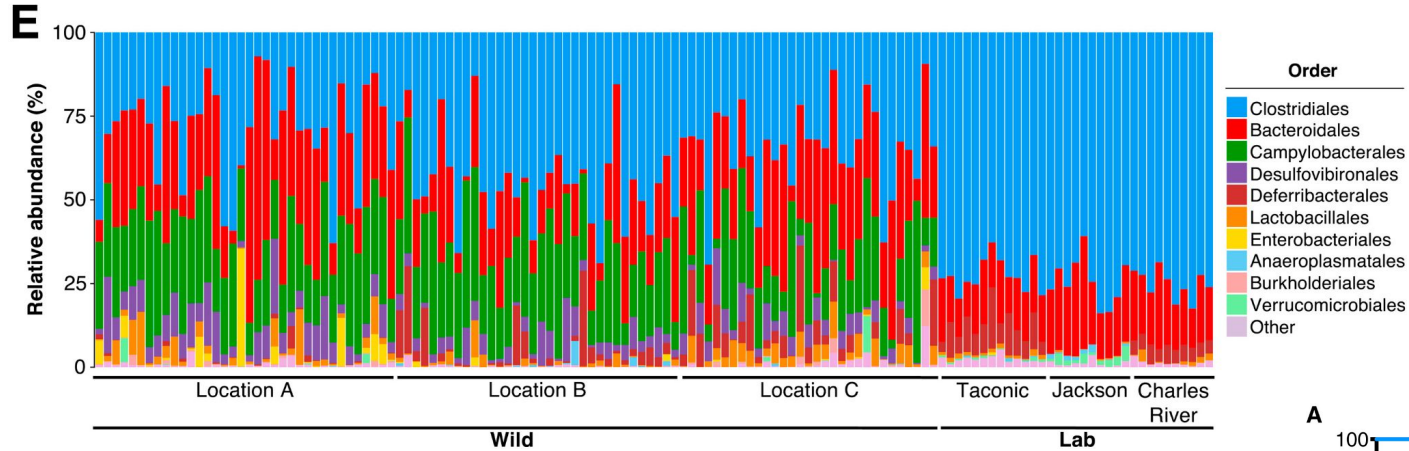# Reads + Reference Genomes → Taxonomic assignment



Illumina HiSeq 4000

A taxomy
https://doi.org/10.1038/s41598-018-30297-w

Taxonomic distribution of saliva microbiome (kraken 1)

3

# Why do we even need that?

- Wild rat microbiome improves lab rat health
  https://doi.org/10.1016/j.cell.2017.09.016
- Zika virus evolution
  https://doi.org/10.1038/nature22402
- Finding the source of Pseudomonas aeruginosa infections in a hospital
  http://dx.doi.org/10.1136/bmjopen-2014-006278

# Rat Microbiome Diversity



High diversity in wild rats

Low diversity in lab rats

# Methods

1. **BLAST**: Search entire read
   - "Classic"
   - accurate
   - slow
2. **k-mers**: Split read into k-mers
   - faster
   - near-BLAST accuracy
3. **marker genes**: DB with only few genes (16S)
   - very fast
   - less accurate

# Kraken 1

1. Build fast mapping reference genomes k-mers -> taxonomy id
   - Minimizer
   - Offset table
2. Query all k-mers from read
3. Assign read by majority assignment of k-mers

# Minimizer (l-mer)

Minimizer: Lexicographically smallest l-mer inside the k-mer

```
ATGATCGTGCATCCATCAGTCATCAGTA <- initial k-mer

ATGATCGTGCATCCATCAGTCA
 TGATCGTGCATCCATCAGTCAT
  TATCGTGCATCCATCAGTCATC
   ATCGTGCATCCATCAGTCATCA <- minimizer¹ (l-mer)
    TCGTGCATCCATCAGTCATCAG
     CGTGCATCCATCAGTCATCAGT
      GTGCATCCATCAGTCATCAGTA
```
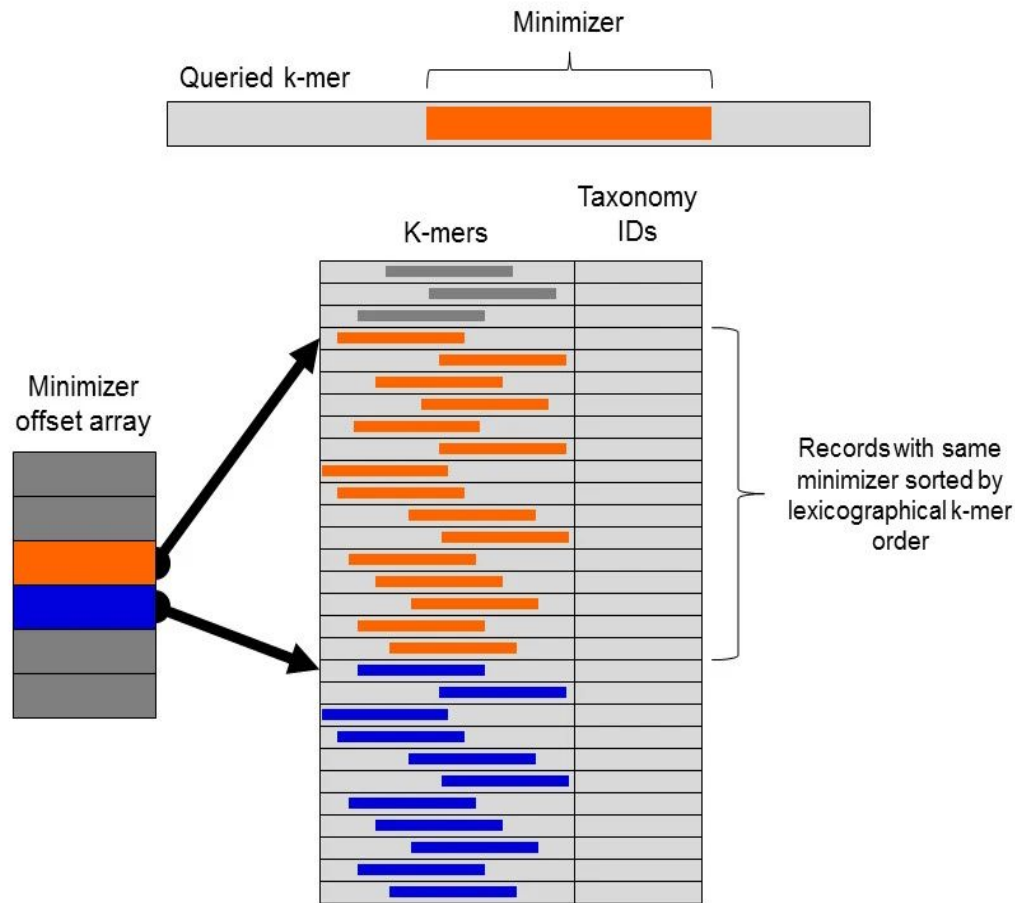
¹ignoring XOR low complexity masking

# Kraken 1

- Group k-mers by minimizers
- Minimizer offset array
  - $8 \times 4^l$ bytes for l-mers
- Good for CPU caching

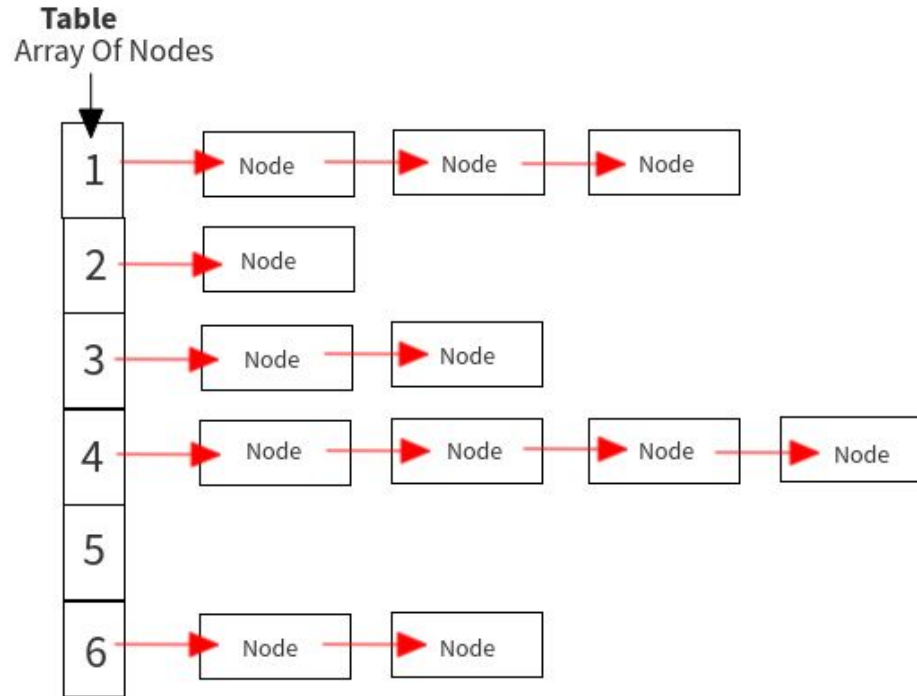# Can we compress k-mers even more?

# Spaced k-mers

Mask *s* positions in l-mers (default: *s*=7)

```
AGTATCGTGCATCGATCAGTCATCAGTA  <- original k-mer

    TCGTGCATCGATCAGTCATCA       <- initial minimizer

    TCGTGCATCGATCAG-C-T-A       <- spaced minimizer
```

# Kraken 2

- Spaced k-mers
- k-mers share minimizers -> store only distinct minimizers
  - all k-mers -> few l-mers
- Only store hash of minimizer
- Defaults: k=35 l=31

# Textbook Hash Maps

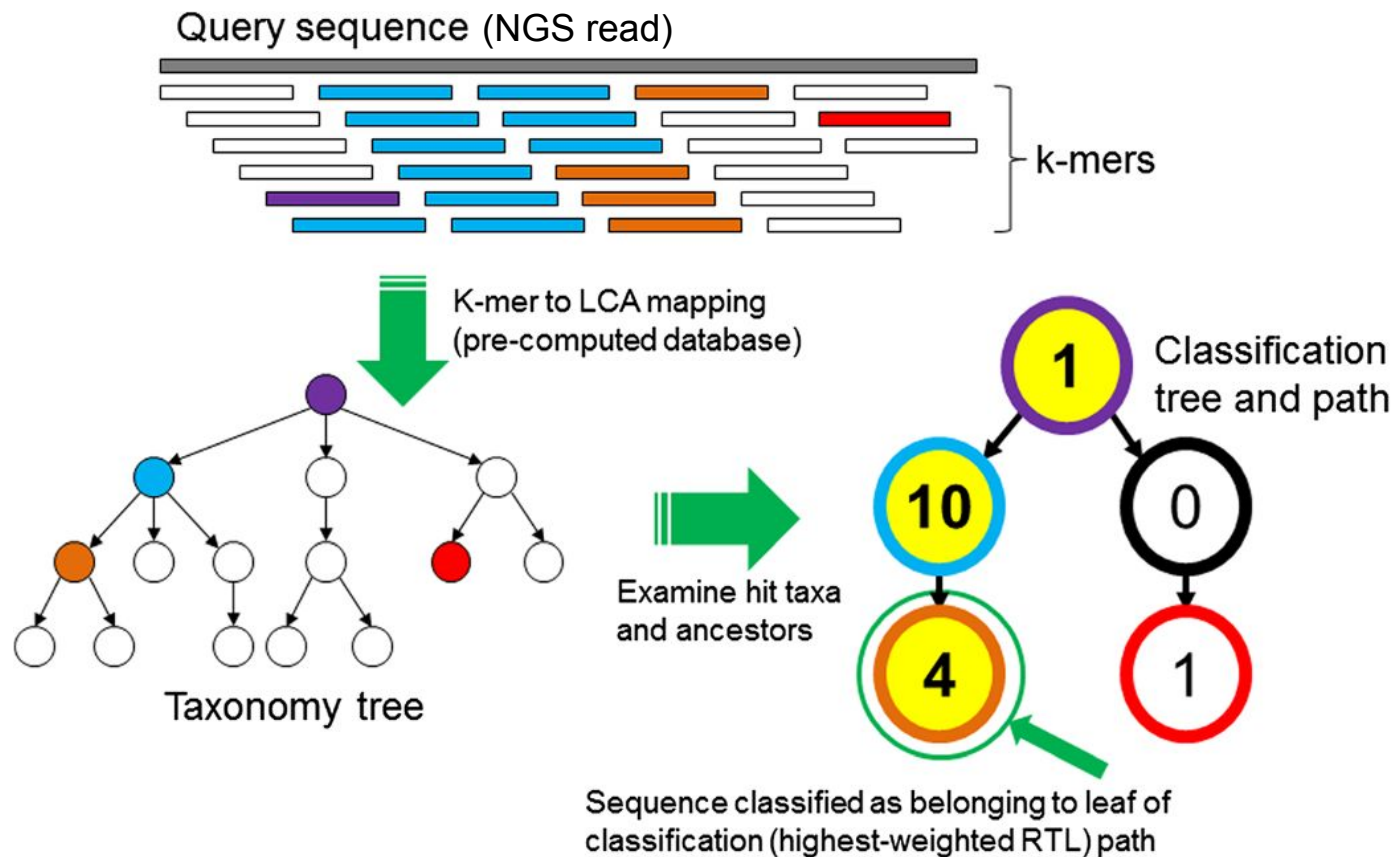# Compact Hash Table and Linear Probing



Load: 45%

Load: 73%

# Kraken 2: Memory Usage

- 4 bytes per entry: 15 bits hash, 17 bytes taxonomy id
- 70% load factor and D distinct minifier: **4D/0.7** bytes
- 85% less memory and 5x faster than Kraken 1

# Building the Database

- When two minimizers collide we set them to the *Lowest Common Ancestor* (LCA) of both source species
- Example: Insulin in human and chimp -> LCA is Primate

```
NM_001008996.    1 AGCCCTCCAGGACAGGCTGCATCAGAAGAGGCCATCAAGCAGATCACTGT    50
                   ||||||||||||||||||||||||||||||||||||||||||||||||||
NM_000207.3      1 AGCCCTCCAGGACAGGCTGCATCAGAAGAGGCCATCAAGCAGATCACTGT    50

NM_001008996.   51 CCTTCTGCCATGGCCCTGTGGATGCGCCTCCTGCCCCTGCTGGTGCTGCT   100
                   |||||||||||||||||||||||||||||||||||||||||·||||||
NM_000207.3     51 CCTTCTGCCATGGCCCTGTGGATGCGCCTCCTGCCCCTGCTGGCGCTGCT   100

NM_001008996.  101 GGCCCTCTGGGGACCTGACCCAGCCTCGGCCTTTGTGAACCAACACCTGT   150
                   |||||||||||||||||||||||||||·|·||||||||||||||||||||
NM_000207.3    101 GGCCCTCTGGGGACCTGACCCAGCCGCAGCCTTTGTGAACCAACACCTGT   150

NM_001008996.  151 GCGGCTCCCACCTGGTGGAAGCTCTCTACCTAGTGTGCGGGGAACGAGGC   200
                   ||||||||·|||||||||||||||||||||||||||||||||||||||||
NM_000207.3    151 GCGGCTCACACCTGGTGGAAGCTCTCTACCTAGTGTGCGGGGAACGAGGC   200
```

# k-mer classification -> read classification

Query sequence (NGS read)

k-mers

K-mer to LCA mapping
(pre-computed database)

Taxonomy tree

Examine hit taxa
and ancestors

Classification
tree and path

Sequence classified as belonging to leaf of
classification (highest-weighted RTL) path

# Evaluation

| Type | Classifier | Memory Required | Time Required | Reference |
|---|---|---|---|---|
| DNA | CLARK-S | 170 Gb | 40 min | Ounit and Lonardi, 2016 |
| | Kraken | 190 Gb | 1 min | Wood and Salzberg, 2014 |
| | Kraken2 | 36 Gb | 1 min | Wood and Salzberg, 2014 |
| | KrakenUniq | 200 Gb | 1 min | Breitwieser et al., 2018 |
| Protein | DIAMOND | 110 Gb (varies) | 10 min | Buchfink et al., 2015 |
| | MMseqs2 | 85 Gb (varies) | 9 h | Steinegger and Söding, 2017 |
| Markers | MetaPhlAn2 | 2 Gb | 1 min | Truong et al., 2015 |

Benchmarking Metagenomics Tools for Taxonomic Classification, Ye et al. 2019
https://doi.org/10.1016/j.cell.2019.07.010

**A** — Default database species L2 distance

High similarity between Kraken 1 and Kraken 2

# System Requirements and Colab

- 100 GB of disk space
- 29 GB of RAM
- "strongly suggest against using NFS storage"
- Use viral genomes instead

https://colab.research.google.com/drive/1DHLIjov7hSTg1fliDOCErUvl2dW-EM6X?usp=sharing

~End~

**Kraken 2**
$k = 35, \ell = 31$

**Kraken 1**
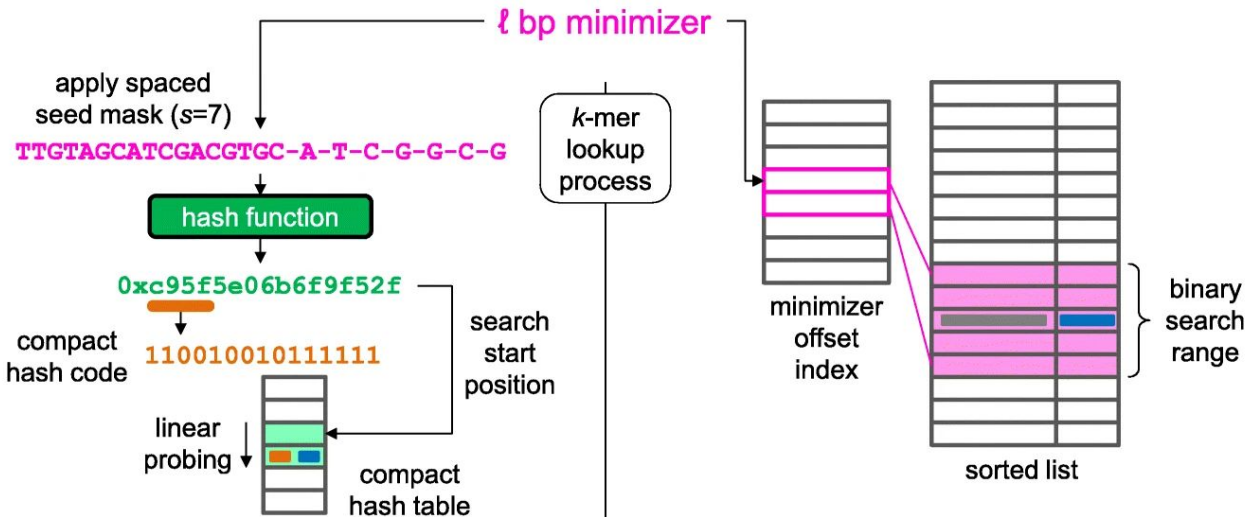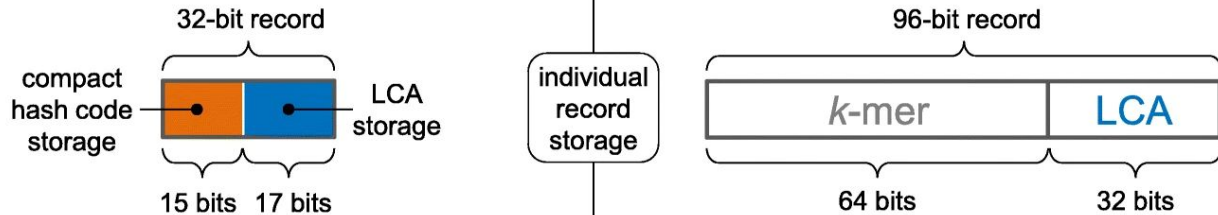$k = 31, \ell = 15$
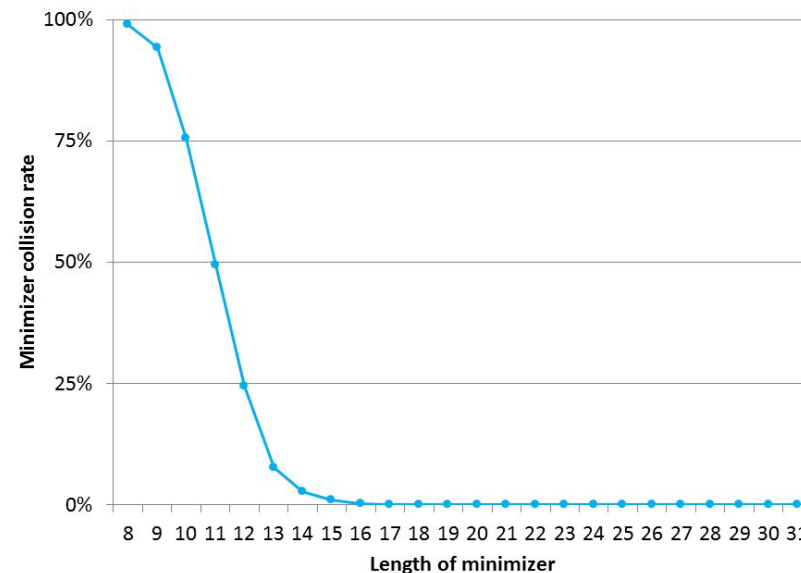
**a**

GCA**TTGTAGCATCGACGTGCCACTGCCGAGTCAG**A $k$-mer

$\ell$ bp minimizer

**b**

apply spaced seed mask ($s$=7)

TTGTAGCATCGACGTGC-A-T-C-G-G-C-G

hash function

0xc95f5e06b6f9f52f

compact hash code

110010010111111

search start position

linear probing

compact hash table

$k$-mer lookup process

minimizer offset index

sorted list

binary search range

**c**

32-bit record

compact hash code storage

LCA storage

15 bits   17 bits

individual record storage

96-bit record

$k$-mer   LCA

64 bits   32 bits

26

# Hash Collisions? Not a Problem

# Why do we even need that?

- Antibiotics can harm stem cell transplantation recipients
  https://doi.org/10.1126/scitranslmed.aaf2311

## Taxonomy



(a)

Simple taxonomy with tax Ids

## Insertions

(b) Insert 8-bit minimizer hash with 0011 as least-sig bits, 0101 as most-sig (01010011), and 4 as taxon of origin. Immediately find empty slot.

(c) Insert 11000100, with 5 as taxon of origin. Immediately find empty slot.

(d) Insert 01010011, with 6 as taxon of origin. Immediately find slot with matching most-sig bits, set value = LCA(4, 6) = 2.

(e) Insert 11000011, with 7 as taxon of origin. Scan forward 1 row before finding match for 1100. Note: **this is a hash table collision**. Set value = LCA(5, 7) = 1.

(f) Insert 10001100, with 8 as taxon of origin. Immediately find empty slot.

(g) Insert 11001101, with 9 as taxon of origin. Immediately find empty slot.

(h) most-sig(hash)  LCA

| | |
|---|---|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0101 | 2 |
| 1100 | 1 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1000 | 8 |
| 1100 | 9 |
| 0 | 0 |
| 0 | 0 |

1100

Query 8-bit minimizer hash 10001100. Find 8, the correct answer for that hash.

(i) most-sig(hash)  LCA

0011

| | |
|---|---|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0101 | 2 |
| 1100 | 1 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1000 | 8 |
| 1100 | 9 |
| 0 | 0 |
| 0 | 0 |

Query 01010011. Find 2, the LCA of the taxa for the two identical hashes we inserted in steps b and d.

(j) most-sig(hash)  LCA

0100

| | |
|---|---|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0101 | 2 |
| 1100 | 1 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1000 | 8 |
| 1100 | 9 |
| 0 | 0 |
| 0 | 0 |

Query 11000100. Find 1, the LCA of the taxa for two hashes, one matching the query and one (11000011) not. These were added in steps c and e.

(k) most-sig(hash)  LCA

0011

| | |
|---|---|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0101 | 2 |
| 1100 | 1 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1000 | 8 |
| 1100 | 9 |
| 0 | 0 |
| 0 | 0 |

Query 11000011. Scan forward 1 row before finding matching least-sig bits. Find 1, the LCA of the taxa for two hashes, one matching the query and one (11000100) not. These were added in steps c and e.

(l) most-sig(hash)  LCA

| | |
|---|---|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0101 | 2 |
| 1100 | 1 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1000 | 8 |
| 1100 | 9 |
| 0 | 0 |
| 0 | 0 |

1100

Query 11001100. Scan forward 1 row before finding matching least-sig bits. This yields the LCA 9, **which is spurious because this query hash was never inserted**.