# Compressibility of the union of sketches
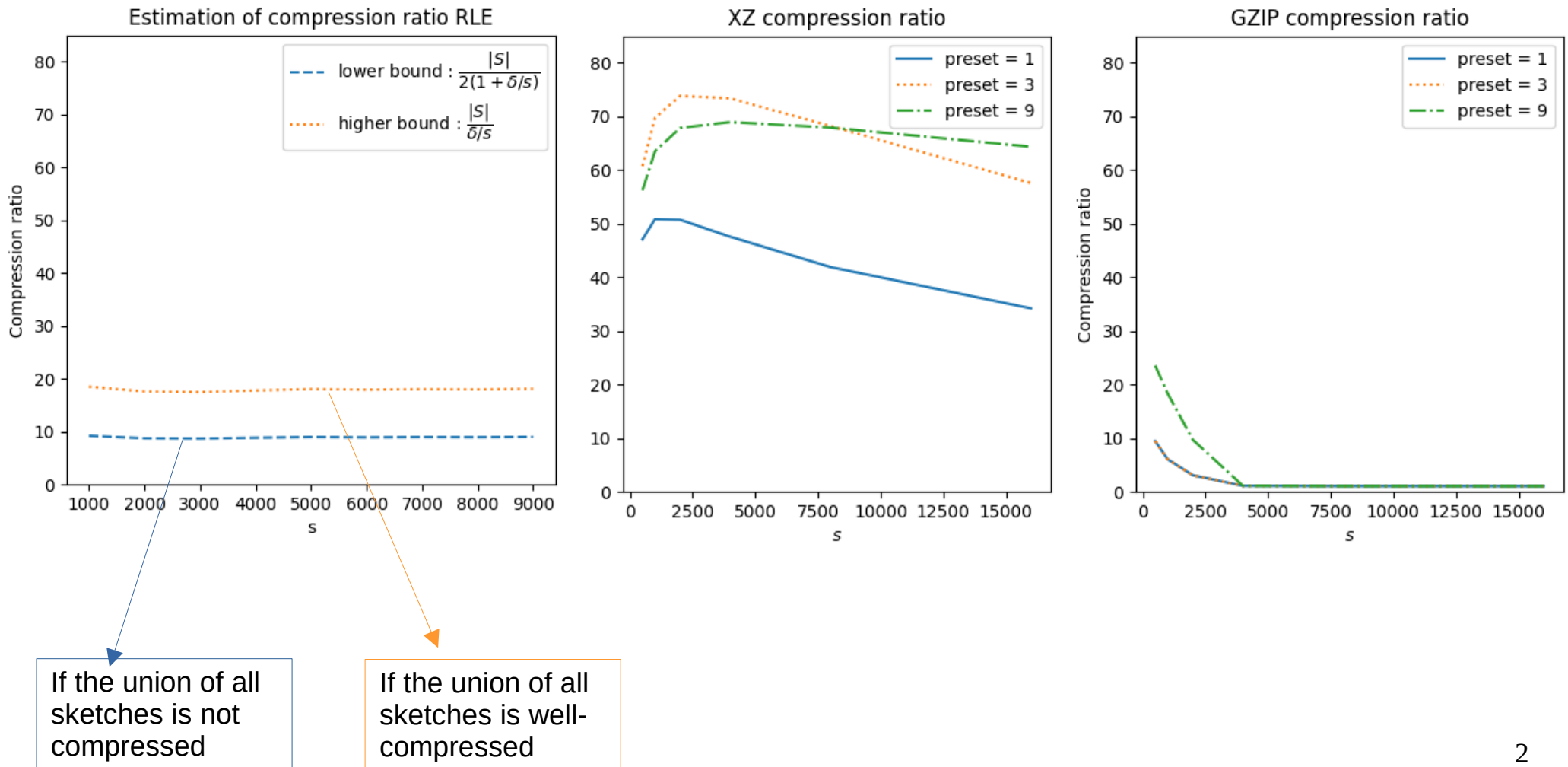## Elias-Fano and RLE

Marie Picard

–

Supervisors : Karel Břinda, Leo Ackermann
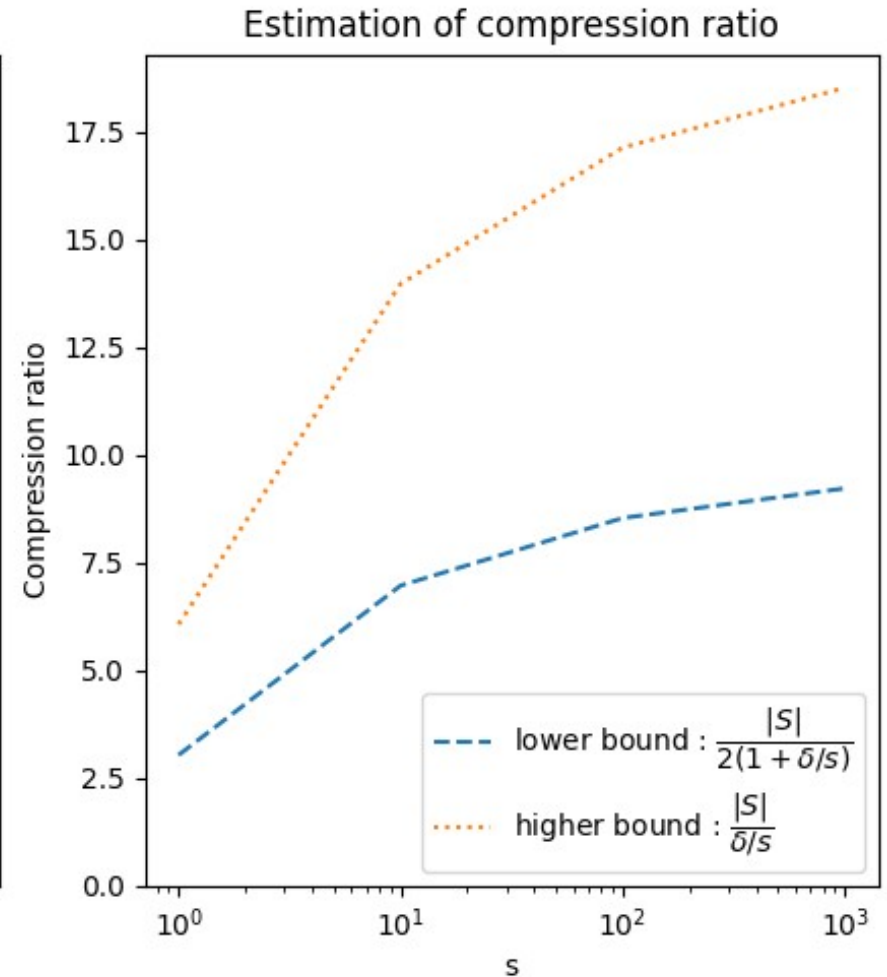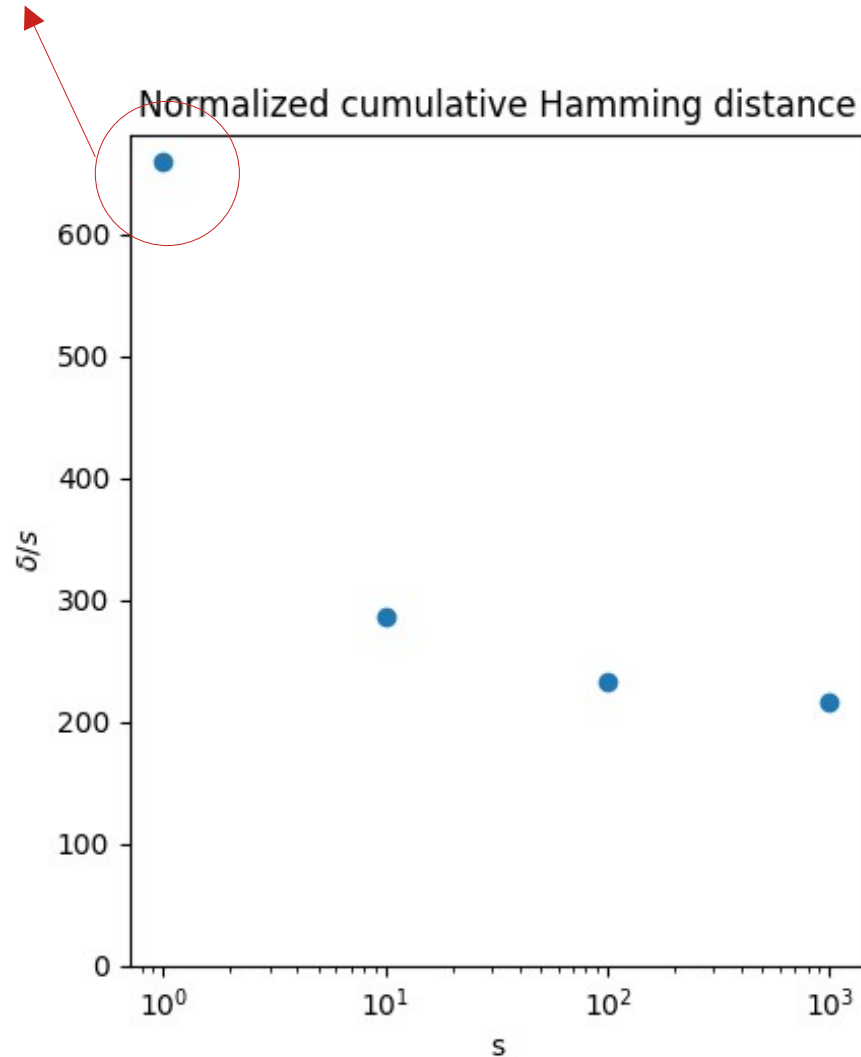
December 4, 2025

# Reminder of last time's results



*Neisseria gonorrhoeae 01, 4000 genomes

2

# Reminder of last time's results

4 different minimal values, one change every 6 genomes



Normalized cumulative Hamming distance

Estimation of compression ratio

lower bound : $\dfrac{|S|}{2(1 + \delta/s)}$

higher bound : $\dfrac{|S|}{\delta/s}$

*Neisseria gonorrhoeae 01, 4000 genomes

3

# What can we take away ?

- Factor 2 depending on how well the union of sketches is compressed

- For the binary matrix :
  - Better compression scheme than RLE
  - Is the order optimal ?

# What can we take away ?
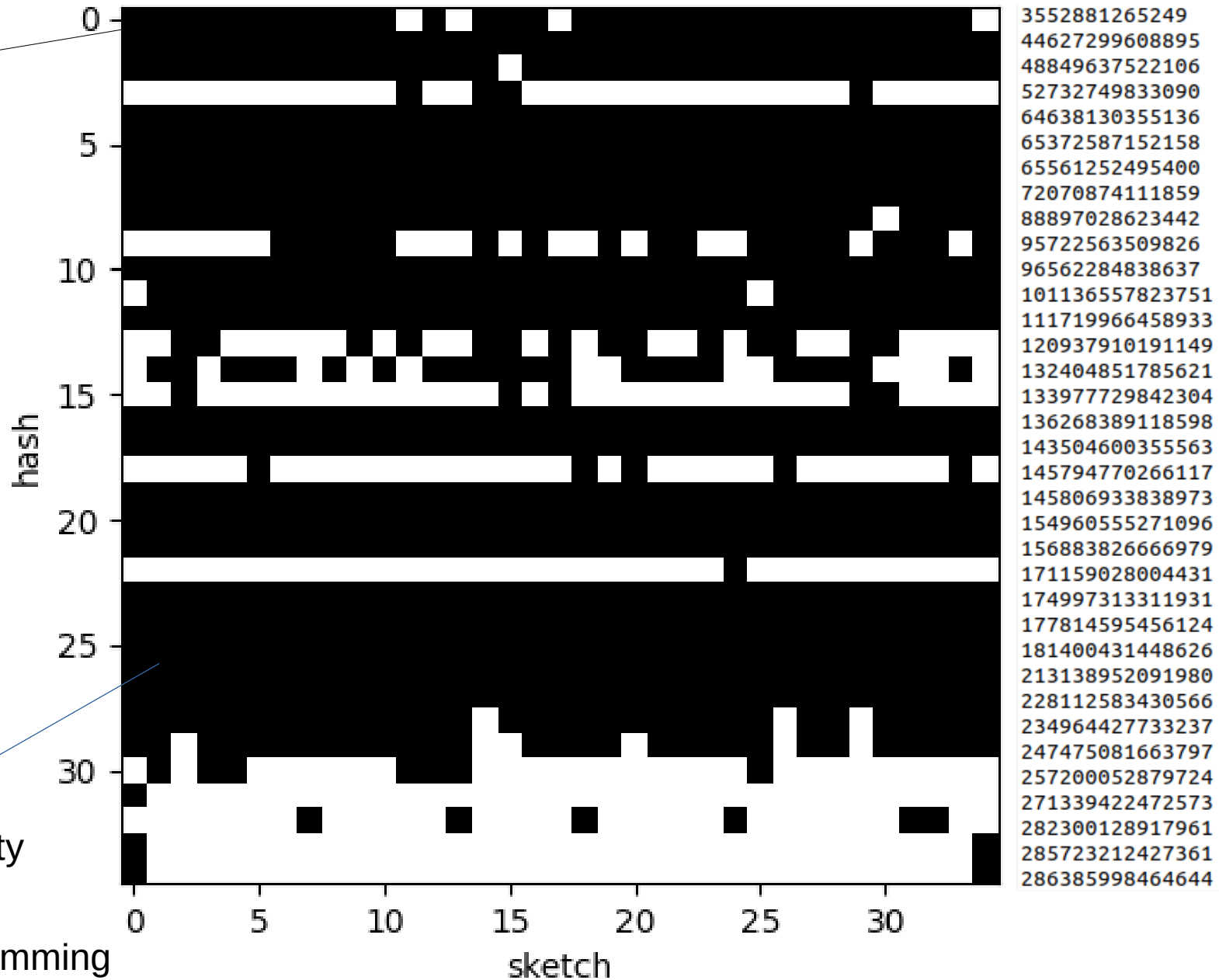
- Factor 2 depending on how well the union of sketches is compressed

- For the binary matrix :
    - Better compression scheme than RLE
    - Is the order optimal ?

# Binary matrix representation

- $|S| = 35$ (number of genomes sketched)
- s = 25 (number of hashes per sketch)
- Types of genomes sketched :
  - Neisseria gonorrhoeae (part 54, n°01)
  - Dustbin (part 24, n°23)
- Phylogenetic order
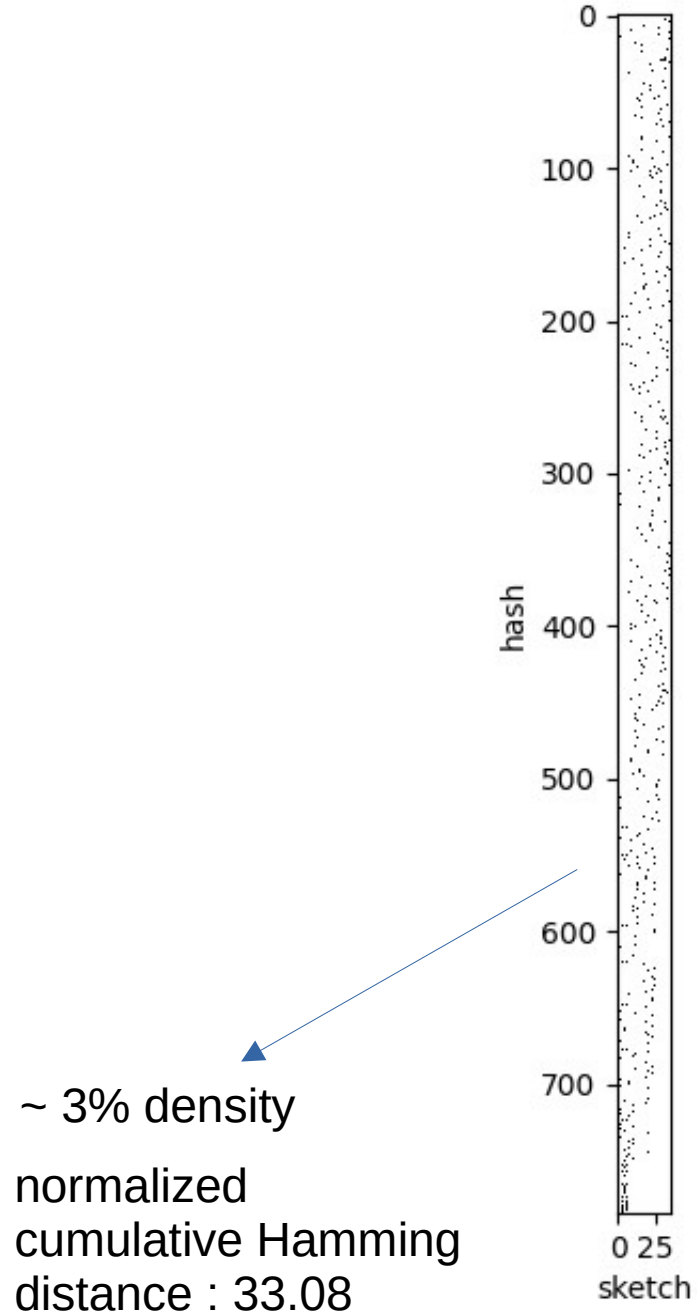
Presence-absence matrix for ngono - 35 sketches - s = 25

Black = presence
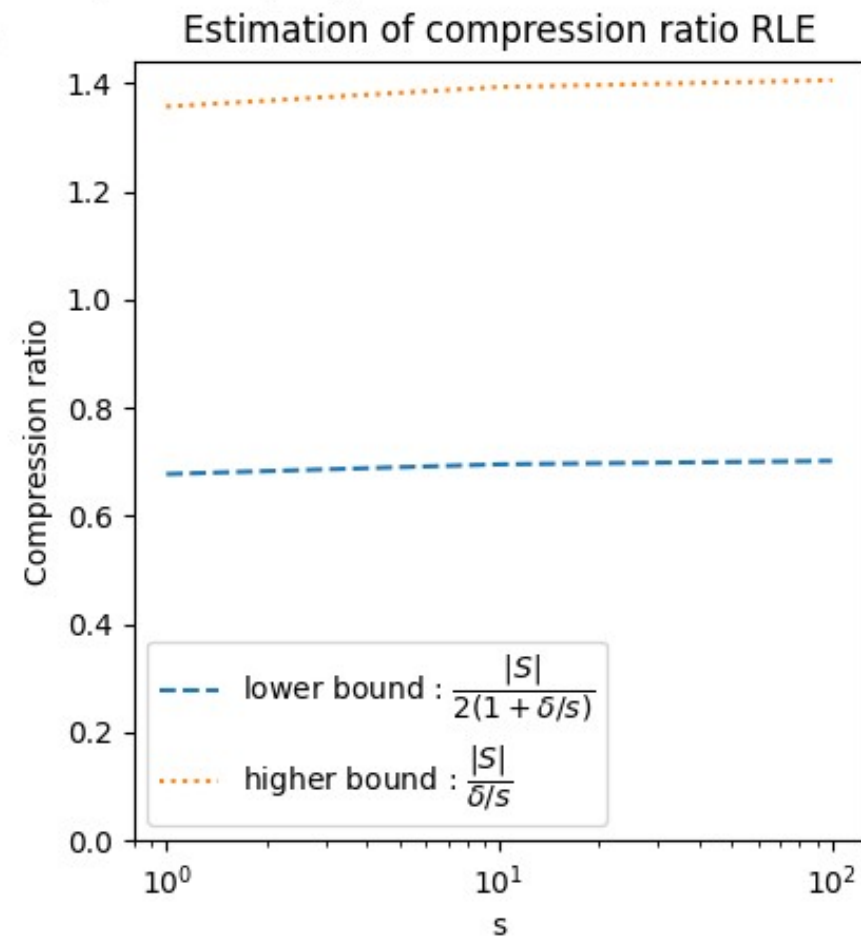
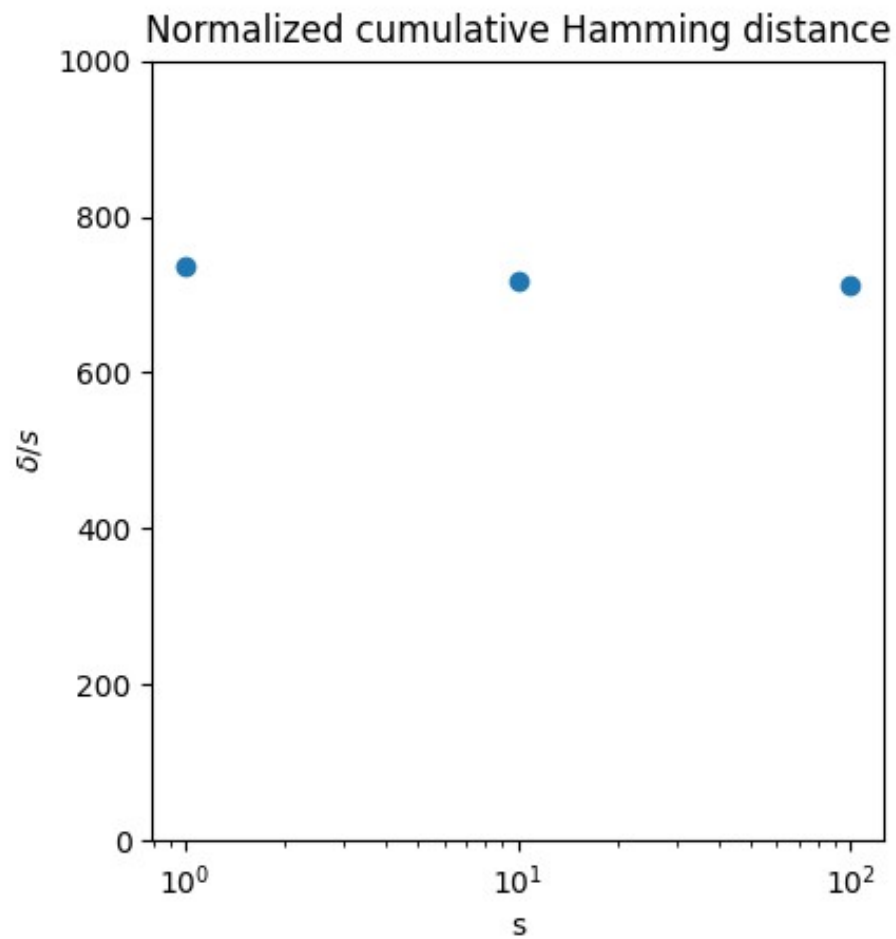> 70% density

normalized cumulative Hamming distance : 2.56

Presence-absence matrix for dustbin - 35 sketches - s = 25



~ 3% density

normalized
cumulative Hamming
distance : 33.08

8

# Compressibility of dustbin

- |S| = 1000 (number of genomes sketched)
- s in [1, 100] (still computing for s = 1000)
- Type of genomes sketched :
  - Dustbin (part 24, n°23)

Dustbin sketches compression (RLE)

Normalized cumulative Hamming distance

Estimation of compression ratio RLE

lower bound : $\dfrac{|S|}{2(1 + \delta/s)}$

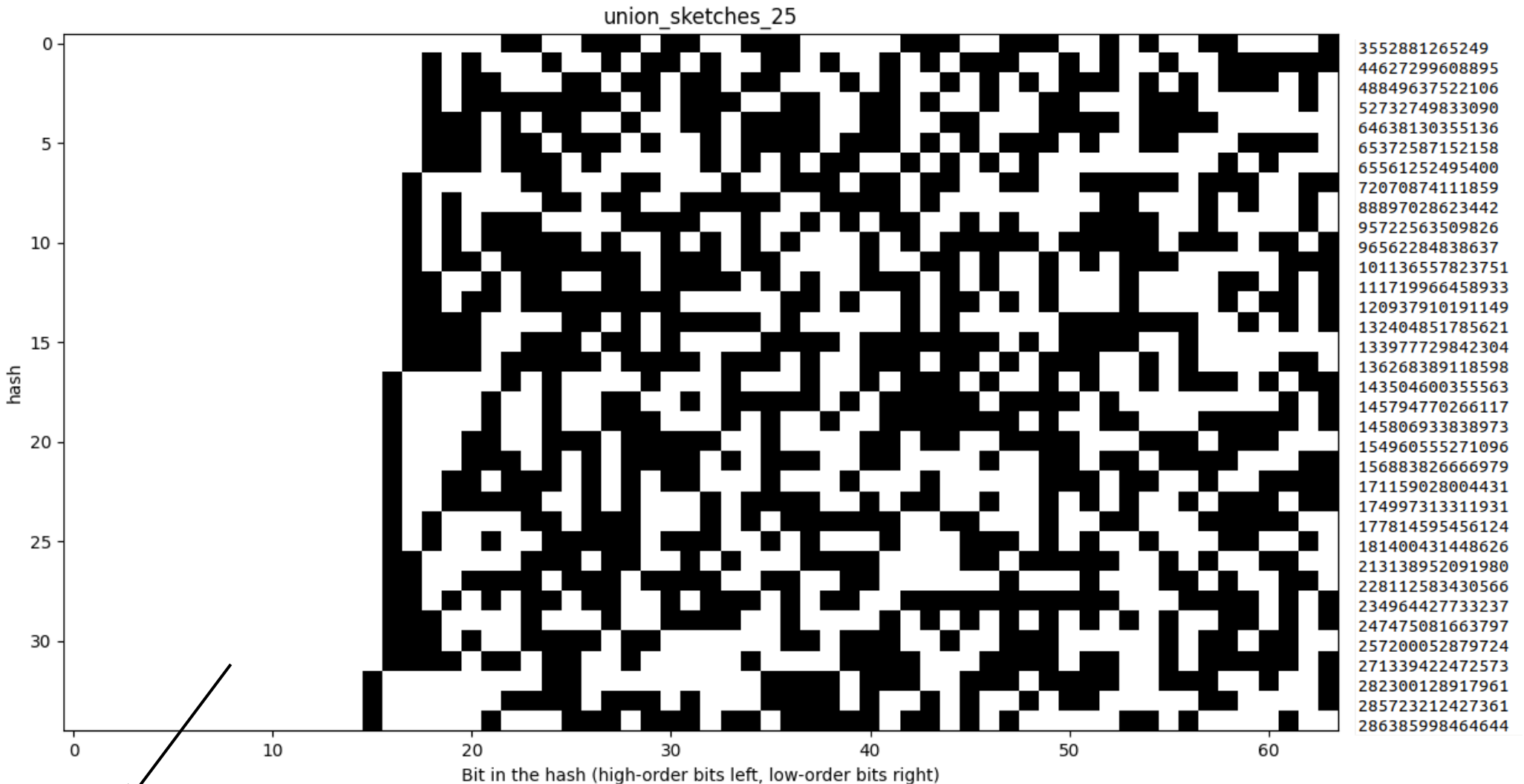higher bound : $\dfrac{|S|}{\delta/s}$

# What can we take away ?

- Factor 2 depending on how well the union of sketches is compressed

- For the binary matrix :
  - Better compression scheme than RLE
  - Is the order optimal ?
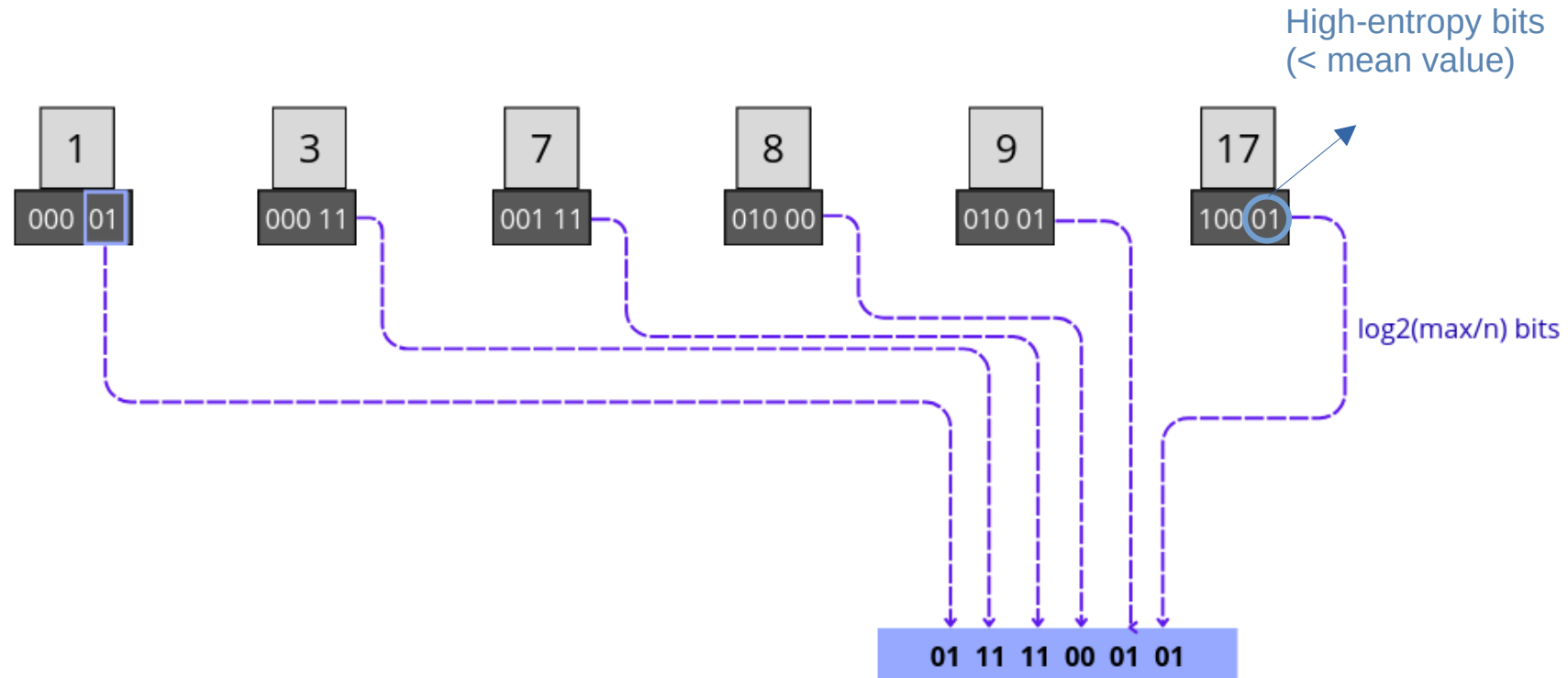
# Representing union of sketches as a binary matrix

- $|S| = 35$ (number of genomes sketched)
- $s = 25$ (number of hashes per sketch)
- Type of genomes sketched :
  - Neisseria gonorrhoeae (part 54, n°01)
- Phylogenetic order
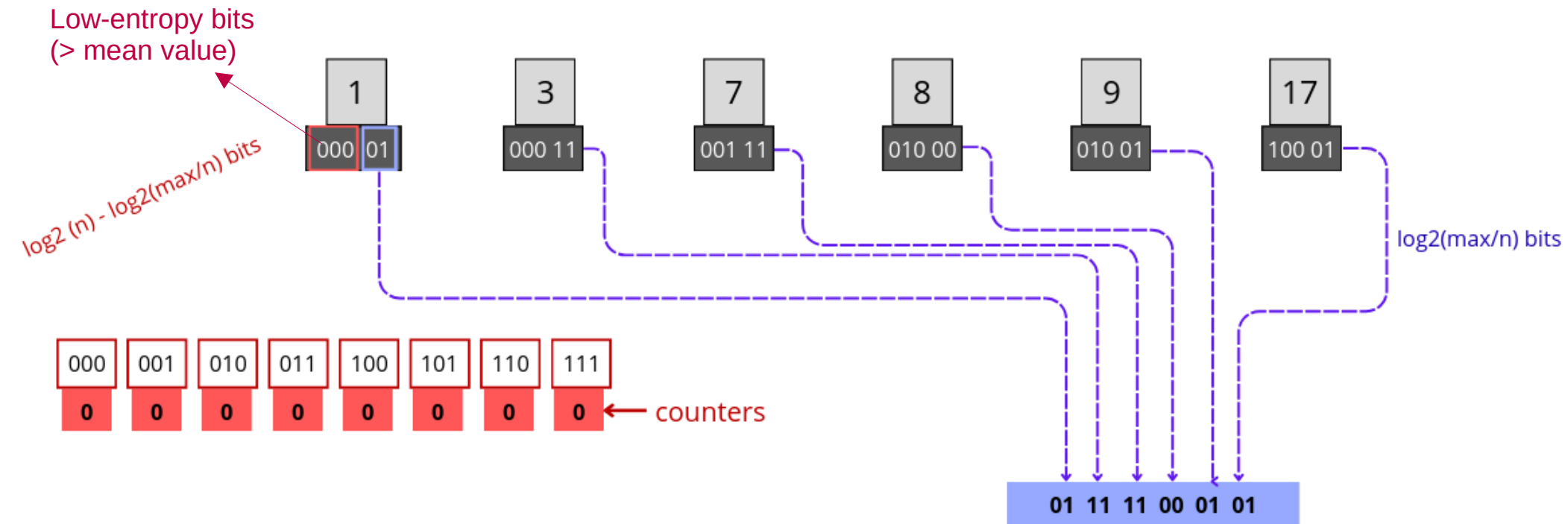
# Union of hashes as a binary matrix



union_sketches_25

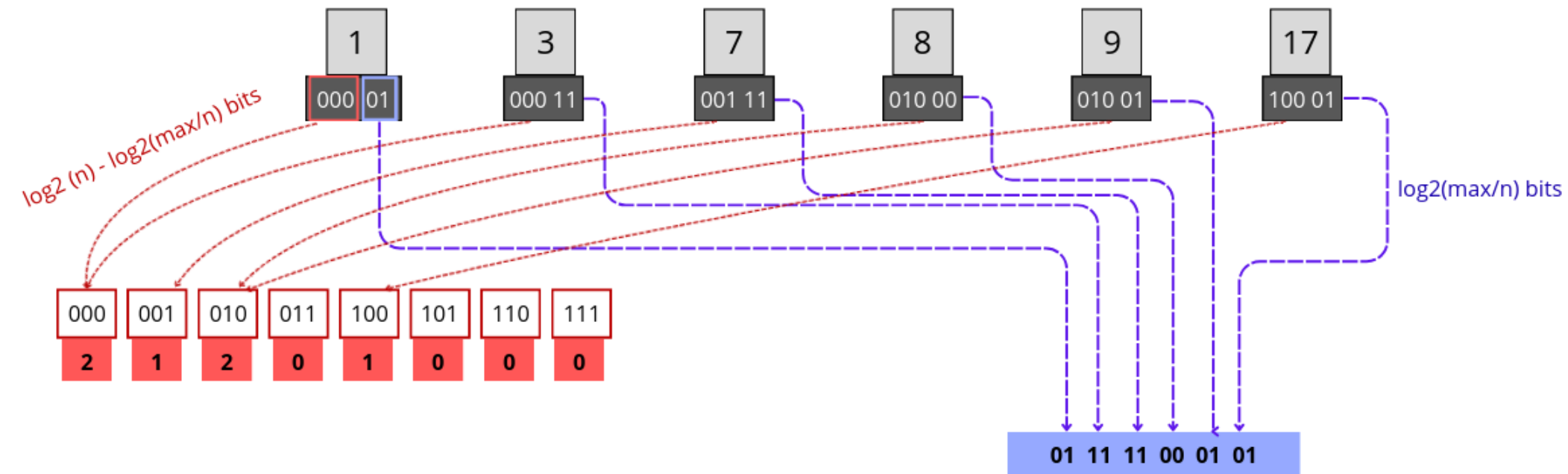**High-ordered bits are well-compressible, low-order bits not (RLE)**

# Elias-Fano compression



High-entropy bits
(< mean value)

log2(max/n) bits

| | | | | | |
|---|---|---|---|---|---|
| 1 | 3 | 7 | 8 | 9 | 17 |
| 000 01 | 000 11 | 001 11 | 010 00 | 010 01 | 100 01 |

01 11 11 00 01 01

# Elias-Fano compression



15

# Elias-Fano compression

# Elias-Fano compression



Final compressed representation

Unary encoding of
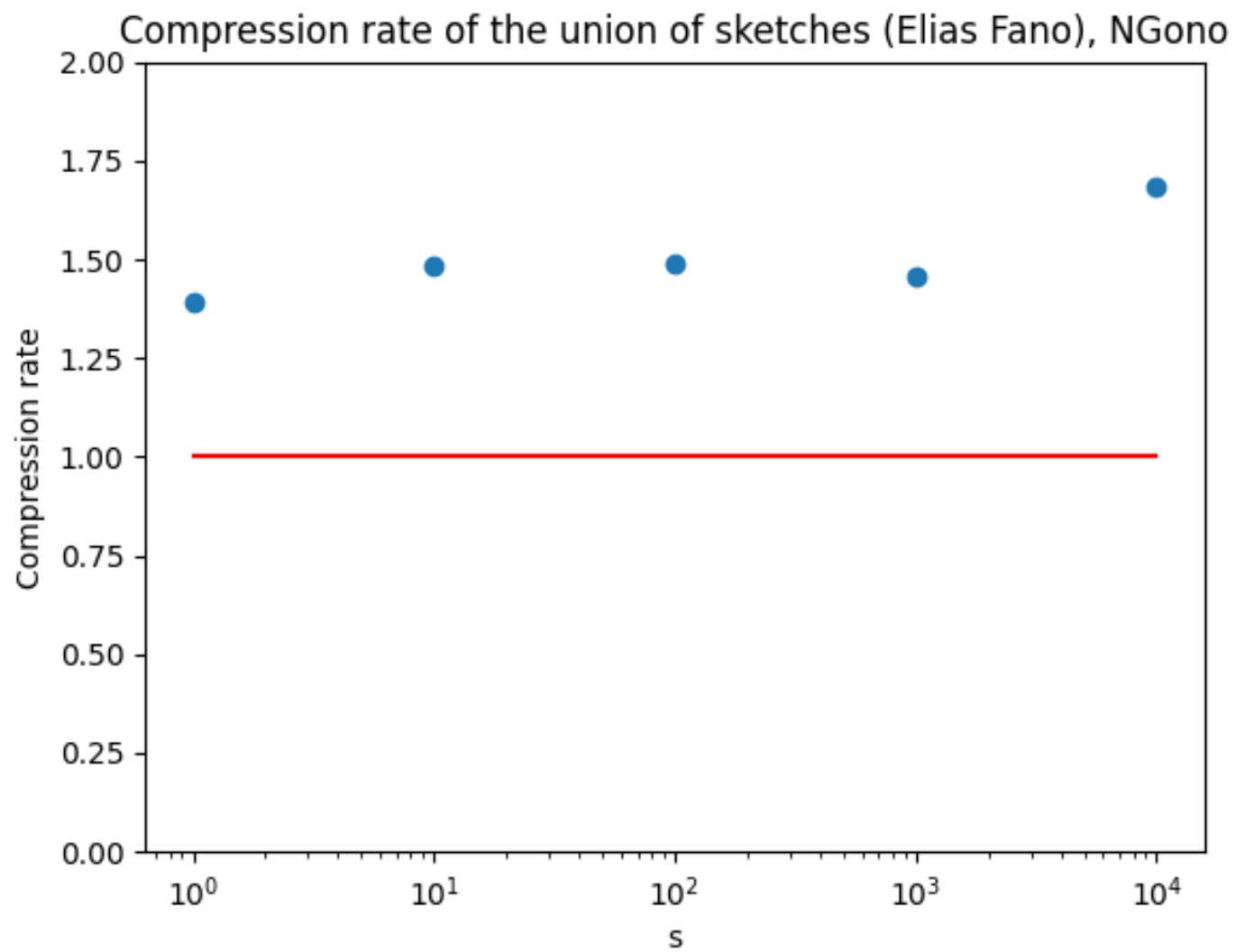the counters

# Elias-Fano compression



Final compressed representation

$$2n + n(\lceil \log_2(m) \rceil - \lceil \log_2(n) \rceil) \text{ bits}$$

# Elias Fano compression

- |S| = 4000 (number of genomes sketched)
- s in [1, 10000] (log scale)
- Type of genomes sketched :
  - neisseria_gonorrhoeae_01 (part 54)

Compression rate of the union of sketches (Elias Fano), NGono

# Conclusion

- Results :
  - Reasonable compression rate with EF
  - Phylogenetic order helps – but a suborder could be considered

- Todo :
  - Look at the total compression rate (RLE + EF)
  - Look into new compression schemes for the union