

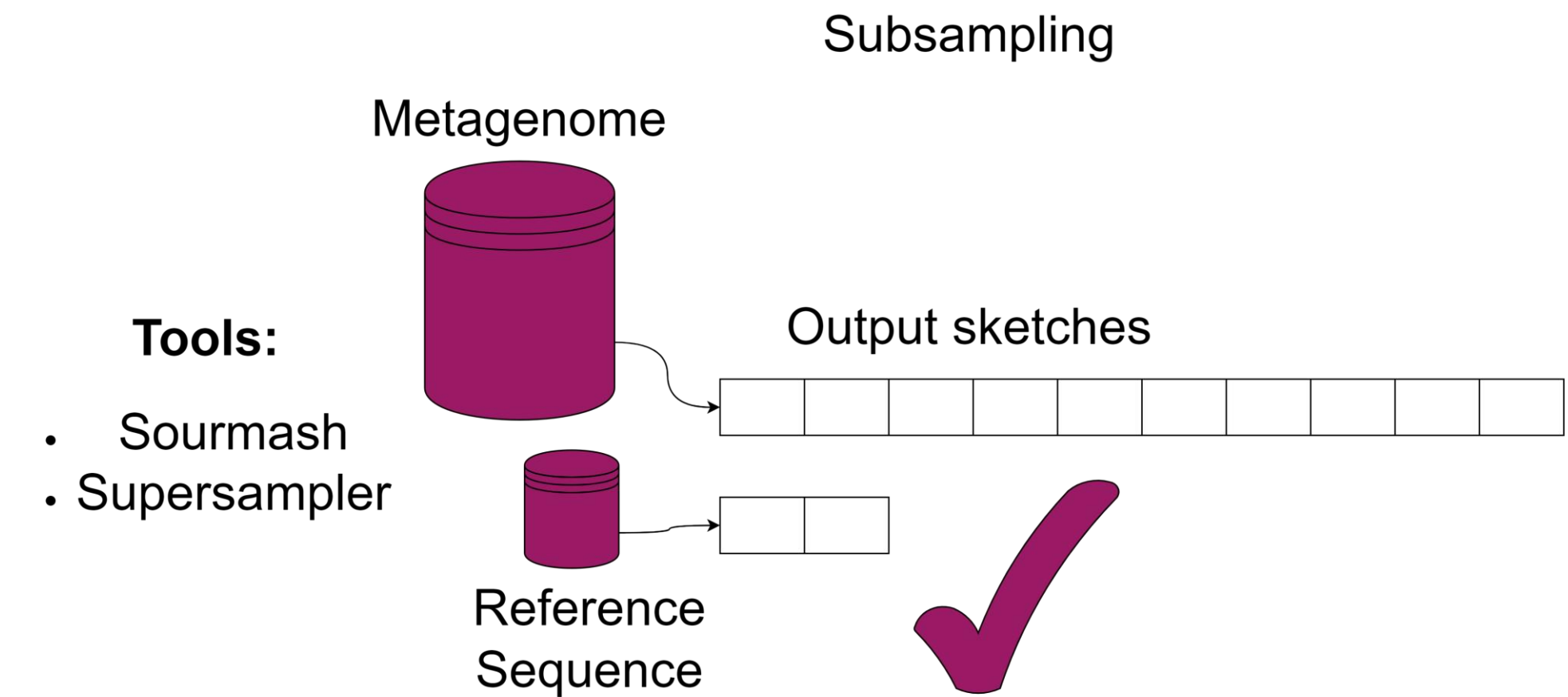
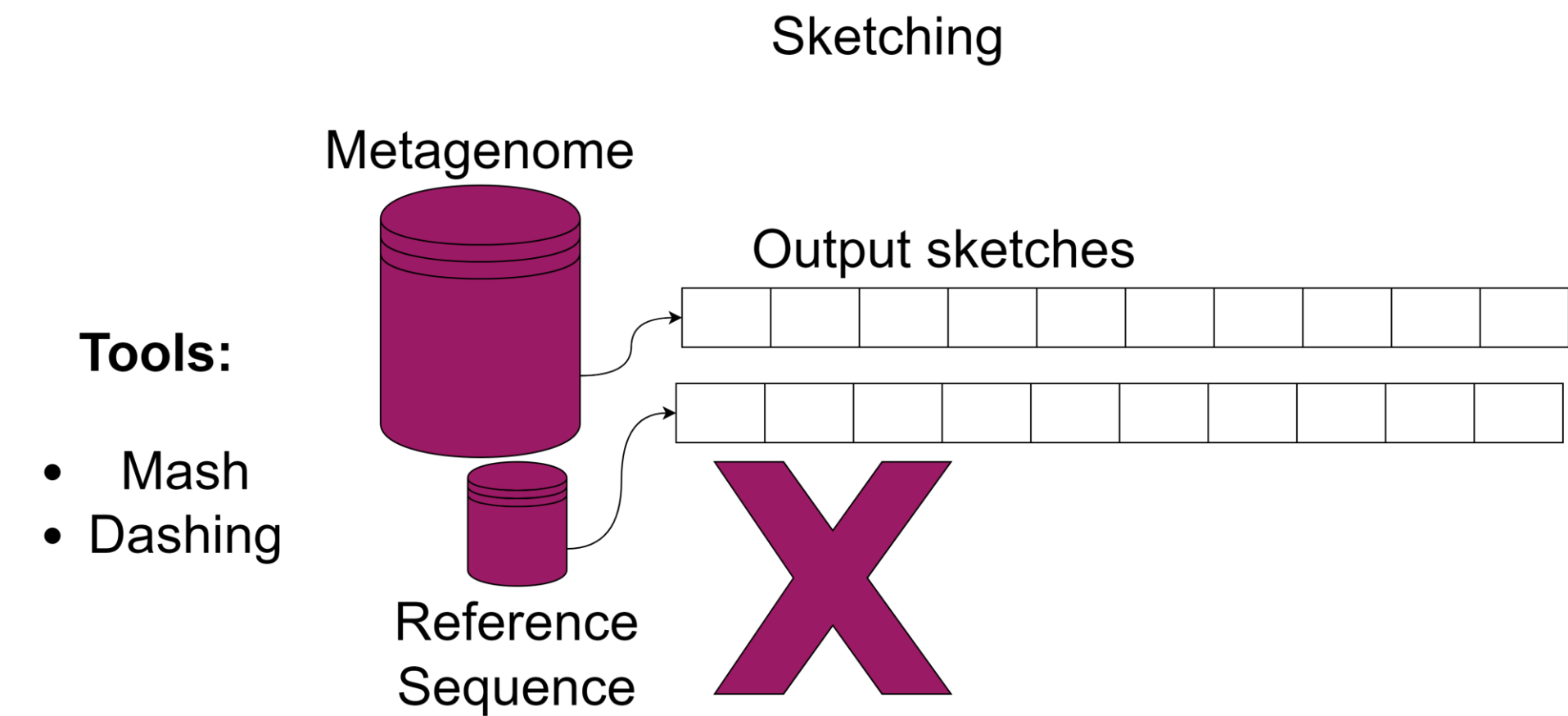
SuperSampler, a memory efficient subsampling strategy for large scale analysis of sequencing data

Timoth  Rouz ¹, Caleb Smith¹, Antoine Lefevre¹, Antoine Limasset¹

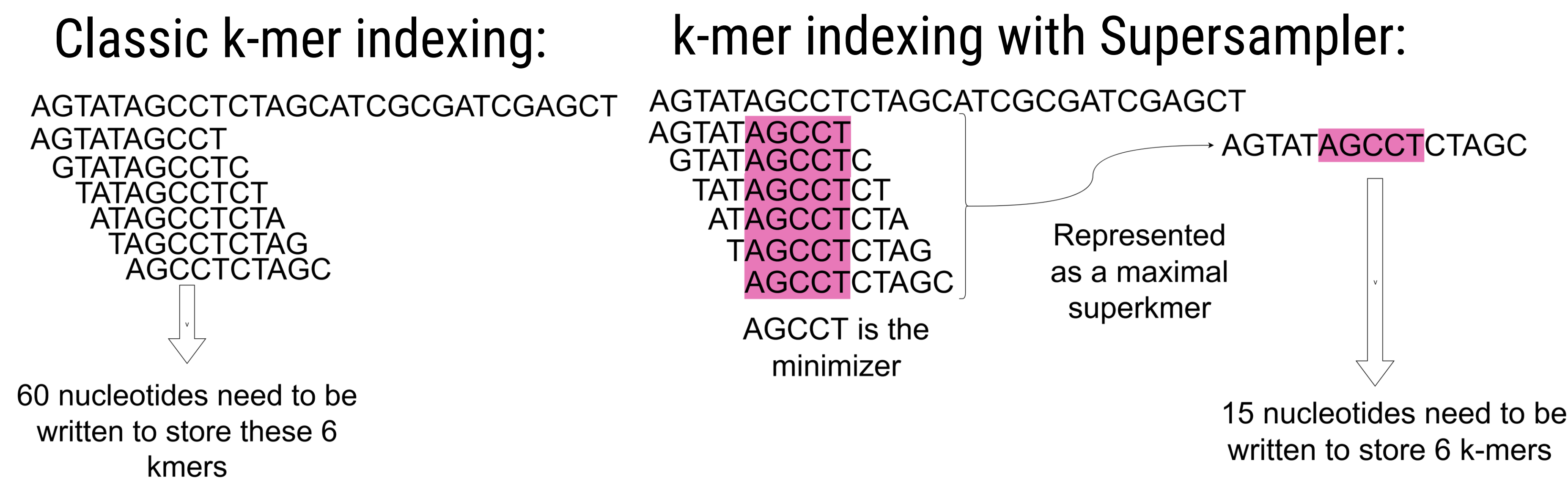
1. Univ. Lille, CNRS, Centrale Lille, UMR 9189 CRIS AL, F-59000 Lille, France

BACKGROUND

Different scaling strategies exist and among these strategies, two main approaches are prevalent. Sketching produces sketches of the same size whatever the input size, when subsampling produces sketches proportional to the input size. The sketching method causes issues when dealing with datasets that are highly different in sizes. Subsampling addresses this issue while keeping comparability between sketches.



Methods: Supersampler uses superkmers to reduce sketches size



Formulas for estimated memory cost of storing k-mers or superkmers

K-mer indexing memory cost:

$$2 \times K \times \frac{N}{S} \text{ bits}$$

Superkmer indexing memory cost:

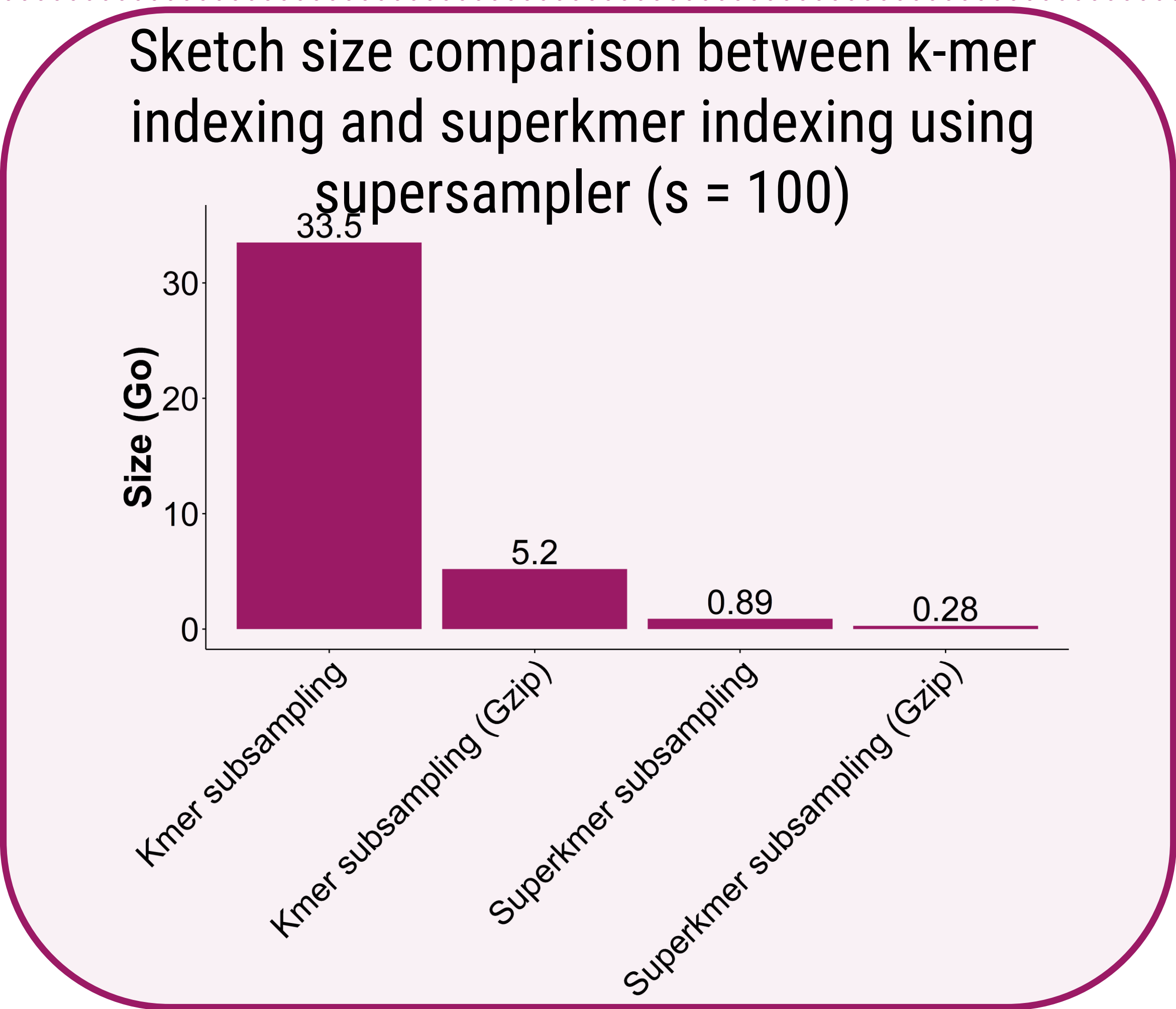
$$2 \times \frac{2K-M}{K-M} \times \frac{N}{S} \text{ bits}$$


Theoretical size of a Lung fish genome subsampled 1/100

| | K = 31 | K = 63 |
|------------|---------|---------|
| k-mers | 24,8 Go | 50,4 Go |
| Superkmers | 1,88 Go | 1,72 Go |

Lungfish genome: 40Gbp

Results on lungfish genome



Sketch sizes and computation time comparison between Sourmash and SuperSampler (s = 100)

| | Wall clock (H:m:s) | Output Size | Size Gzipped |
|--------------|--------------------|-------------|--------------|
| Sourmash | 1:26:50 | 4,2 Go | 1,6 Go |
| SuperSampler | 1:01:53 | 0,89 Go | 0,28 Go |

- Conclusion
- Supersampler is more efficient than state of the art subsampling tools both in term of computation time and in term of memory size of the output.
- What's next?
- Add multithreading to improve computation speed
 - Create a sketch comparison tool
 - Statistical comparison between k-mer and superkmer indexing

References:

- Mash: <https://mash.readthedocs.io/en/latest/index.html>
- Dashing: Baker, D.N., Langmead, B. Dashing: fast and accurate genomic distances with HyperLogLog. Genome Biol 20, 265 (2019). <https://doi.org/10.1186/s13059-019-1875-0>
- Sourmash: <https://sourmash.readthedocs.io/en/latest/index.html>

