

On the compressibility of sketches

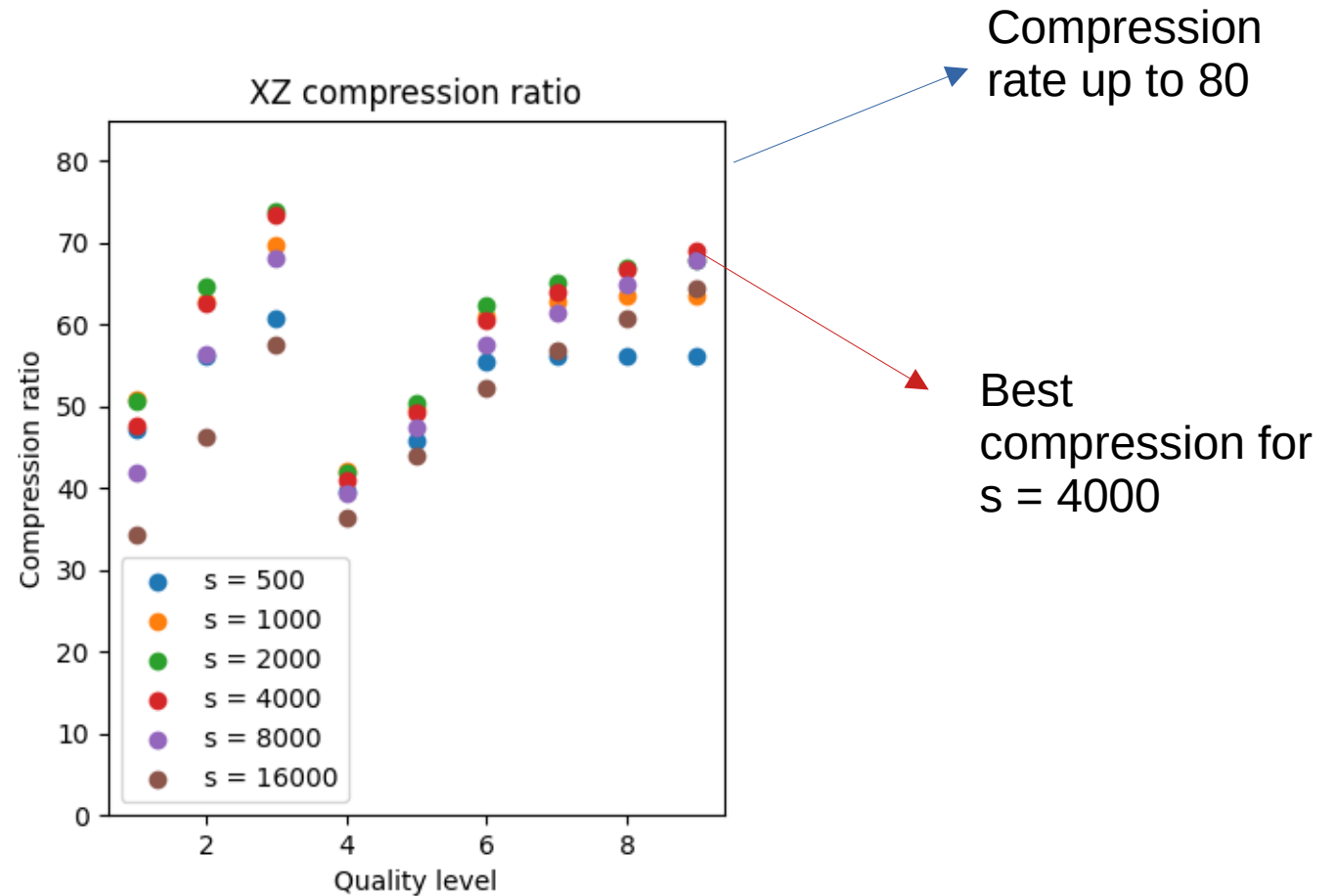
RLE-compression & Hamming distance

Marie Picard

Supervisors : Karel Břinda,
Leo Ackermann

November 20, 2025

Quick recap of previous results



H (informal)

There is a value of sketch size s , around 2000 to 4000, such that the compressibility of sketches is optimal.

H (formal)

$\forall \mathcal{G}, \exists s \in \mathbb{N}, \frac{1}{s} K(\text{Sketches}(\mathcal{G}, s))$ is minimal where

- ▶ K is Kolmogorov's complexity
- ▶ \mathcal{G} is the set of genomes to sketch
- ▶ $\text{Sketches}(\mathcal{G}, s)$ is the archive of all sketches of $g \in \mathcal{G}$ of size s

H (informal)

There is a value of sketch size s , around 2000 to 4000, such that the compressibility of sketches is optimal.

H (formal)

$\forall \mathcal{G}, \exists s \in \mathbb{N}, \frac{1}{s} K(\text{Sketches}(\mathcal{G}, s))$ is minimal where

- ▶ K is Kolmogorov's complexity
- ▶ \mathcal{G} is the set of genomes to sketch
- ▶ $\text{Sketches}(\mathcal{G}, s)$ is the archive of all sketches of $g \in \mathcal{G}$ of size s

Not
computable



H

$\forall \mathcal{G}, \exists s \in \mathbb{N}, \frac{1}{s} RLE(Sketches(\mathcal{G}, s))$ is minimal
where

- ▶ RLE is run-length encoding
- ▶ \mathcal{G} is the set of genomes to sketch
- ▶ $Sketches(\mathcal{G}, s)$ is the archive of all mash sketches of $g \in \mathcal{G}$ of size s

Input

List of mesh sketches :

$$\mathcal{S} = (S_1, \dots, S_{|\mathcal{S}|})$$

such that $s = |S_1| = \dots = |S_{|\mathcal{S}|}|$

Compressed representation

- ▶ Union of all sketches (sorted)

$$\mathcal{S}_U := \bigcup_{S \in \mathcal{S}} S$$

- ▶ presence matrix \mathcal{M}
 $\mathcal{M}_{i,j} = 1$ iff the i^{th} element of $\mathcal{S}_U \in S_j$

Higher Hasse values

[illegible]

RLE

Matrix of sketches

		1	1	1	1	...	1	1	1	1
←1	1	1	1	1	1	...	1	1	1	1
←1	1			1	1	...	1		1	1
←					1	...				1
					1	1	1	1
1	1	1	1	1	1	...	1	1	1	1
						...				
1	1	1	1	1	1	...	1	1	1	
1	1	1	1	1		...				
1	1	1	1	1	1	...	1	1	1	1
1	1	1	1	1	1	...	1	1	1	1
						...	1	1	1	
1	1	1	1			...			1	1
						...				
1	1	1	1	1	1	...	1	1	1	1
			1	1	1	...				
						...	1	1		1
1	1	1				...			1	1
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				
						...				

distance δ

$$\delta = \frac{1}{2} \sum_{i=1}^{|\mathcal{S}|-1} d_{Hamming}(S_i, S_{i+1})$$

- ▶ Easy to compute
- ▶ Determines the RLE compression rate

$$\delta = 0$$

Matrix of sketches

		1	1	1	1	...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
1	1			1	1	...		1		1		1
					1	...						
1	1	1	1	1	1	...		1	1	1	1	1
						...						
1	1	1	1	1	1	...		1	1	1		
1	1	1	1	1		...						
1	1	1	1	1	1	...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
						...		1	1	1		
						...					1	1
1	1	1	1			...						
						...						
1	1	1	1	1	1	...		1	1	1	1	1
						...						
			1	1	1	...						
						...		1	1			
1	1	1				...				1	1	1
						...						
						...						
						...						
						...						
\mathcal{S}_1	\mathcal{S}_2	\mathcal{S}_3	\mathcal{S}_4	\mathcal{S}_5	\mathcal{S}_6	...		\mathcal{S}_{41}	\mathcal{S}_{42}	\mathcal{S}_{43}	\mathcal{S}_{44}	\mathcal{S}_{45}

$$\delta = 1$$

Matrix of sketches

		1	1	1	1	...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
1	1			1	1	...		1		1	1	1
					1	...						
						...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
						...						
1	1	1	1	1	1	...		1	1			
1	1	1	1	1		...						
1	1	1	1	1	1	...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
						...		1	1	1		
						...					1	1
1	1	1	1			...						
						...						
1	1	1	1	1	1	...		1	1	1	1	1
			1	1	1	...						
1	1	1				...		1	1	1	1	1
						...						
						...						
						...						
						...						
\mathcal{S}_1	\mathcal{S}_2	\mathcal{S}_3	\mathcal{S}_4	\mathcal{S}_5	\mathcal{S}_6	...		\mathcal{S}_{41}	\mathcal{S}_{42}	\mathcal{S}_{43}	\mathcal{S}_{44}	\mathcal{S}_{45}

$$\delta = 2$$

Matrix of sketches

		1	1	1	1	...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
1	1			1	1	...		1		1	1	1
					1	...						
						...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
						...						
1	1	1	1	1	1	...		1	1	1		
1	1	1	1	1		...						
1	1	1	1	1	1	...		1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1
						...		1	1	1		
						...					1	1
1	1	1	1			...						
						...						
1	1	1	1	1	1	...		1	1	1	1	1
			1	1	1	...						
						...						
1	1	1				...		1	1	1	1	1
						...						
						...						
						...						
\mathcal{S}_1	\mathcal{S}_2	\mathcal{S}_3	\mathcal{S}_4	\mathcal{S}_5	\mathcal{S}_6	...		\mathcal{S}_{41}	\mathcal{S}_{42}	\mathcal{S}_{43}	\mathcal{S}_{44}	\mathcal{S}_{45}

$$\delta \geq 8$$

Matrix of sketches

		1	1	1	1	...			1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1	1
1	1			1	1	...		1			1	1	1
					1	...							
						...		1	1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1	1
						...							
1	1	1	1	1	1	...		1	1	1			
1	1	1	1	1		...							
1	1	1	1	1	1	...		1	1	1	1	1	1
1	1	1	1	1	1	...		1	1	1	1	1	1
						...		1	1	1		1	1
						...					1		
1	1	1	1			...							
						...							
1	1	1	1	1	1	...		1	1	1	1	1	1
			1	1	1	...							
						...							
1	1	1				...		1	1		1	1	1
						...							
						...							
						...							
						...							
						...							
\mathcal{S}_1	\mathcal{S}_2	\mathcal{S}_3	\mathcal{S}_4	\mathcal{S}_5	\mathcal{S}_6	...		\mathcal{S}_{41}	\mathcal{S}_{42}	\mathcal{S}_{43}	\mathcal{S}_{44}	\mathcal{S}_{45}	

Line changes

Number of bit changes upon line skips :

$$x = \sum_{i=1}^{|S_U|-1} |\mathcal{M}_{i,|S|} - \mathcal{M}_{i+1,1}|$$

► $2\delta + x$ bit changes in RLE

► $x \leq 2s - 1$

So at most $2\delta + 2s - 1$ bit changes.

Size of output :

$$|S_U| + |RLE(\mathcal{M})|$$

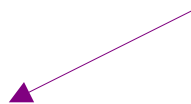
► $2\delta + x$ bit changes in RLE

► $x \leq 2s - 1$

So at most $2\delta + 2s - 1$ bit changes.

Size of output :

$$|S_U| + |RLE(\mathcal{M})|$$



$$64s \leq |S_U| \leq 64(s + \delta)$$

► 64 bits for each hash

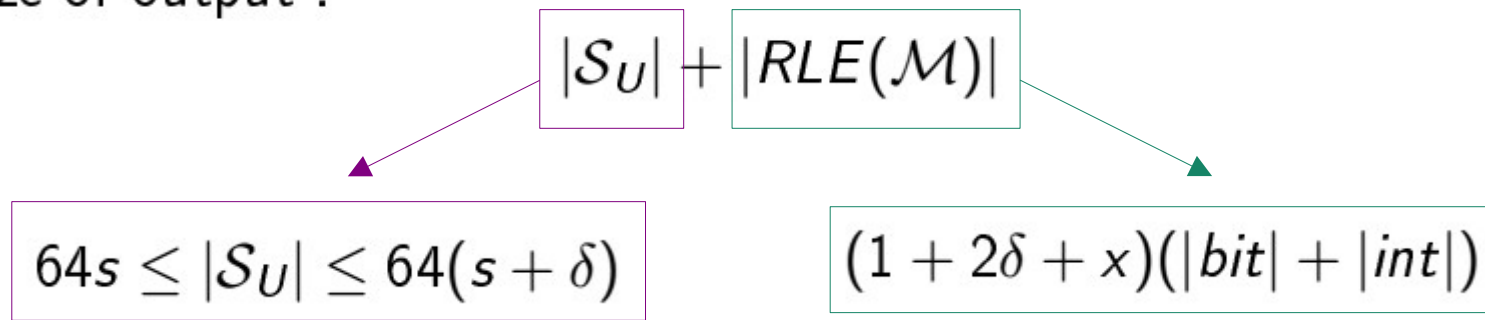
► between s and $s + \delta$ hashes

► $2\delta + x$ bit changes in RLE

► $x \leq 2s - 1$

So at most $2\delta + 2s - 1$ bit changes.

Size of output :

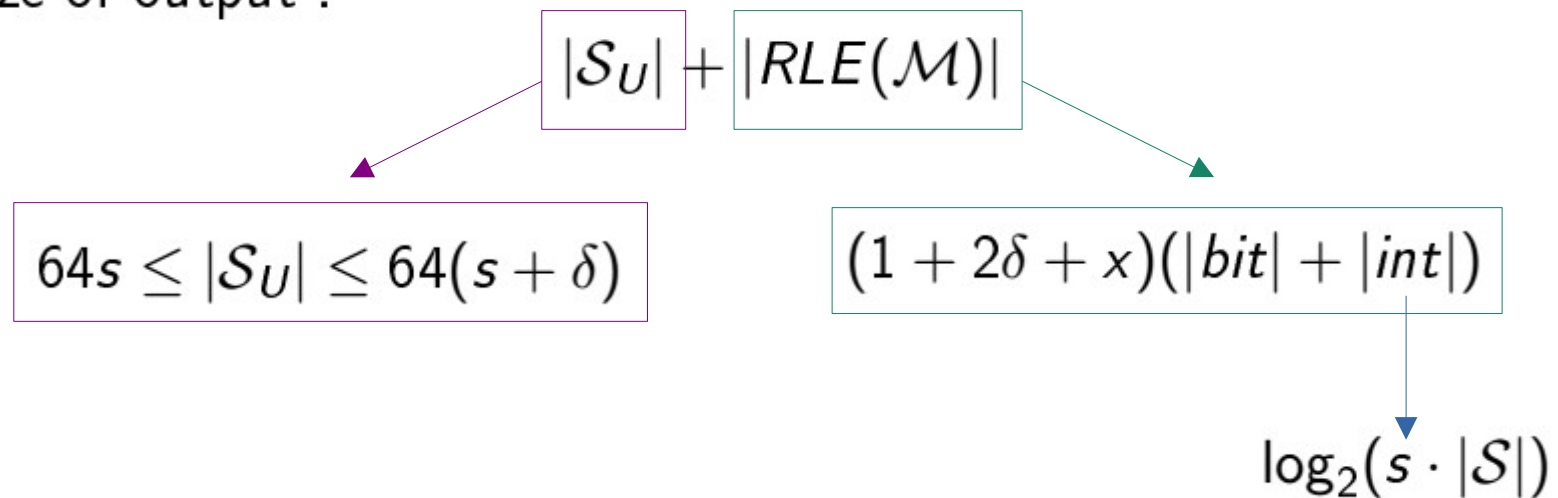
$$|S_U| + |RLE(\mathcal{M})|$$

$$64s \leq |S_U| \leq 64(s + \delta)$$
$$(1 + 2\delta + x)(|bit| + |int|)$$

► $2\delta + x$ bit changes in RLE

► $x \leq 2s - 1$

So at most $2\delta + 2s - 1$ bit changes.

Size of output :



For $s \leq 16000$, $|\mathcal{S}| \leq 131000$, $\log_2(s \cdot |\mathcal{S}|) \leq 31$

► $2\delta + x$ bit changes in RLE

► $x \leq 2s - 1$

So at most $2\delta + 2s - 1$ bit changes.

Size of output :

$$\begin{array}{ccc} & |S_U| + |RLE(\mathcal{M})| & \\ \swarrow & & \searrow \\ \boxed{64s \leq |S_U| \leq 64(s + \delta)} & & \boxed{\begin{array}{l} (1 + 2\delta + x)(|bit| + |int|) \\ \leq 64(\delta + s) \end{array}} \end{array}$$

Size of input :

$$\geq 64 \cdot |\mathcal{S}| \cdot s$$

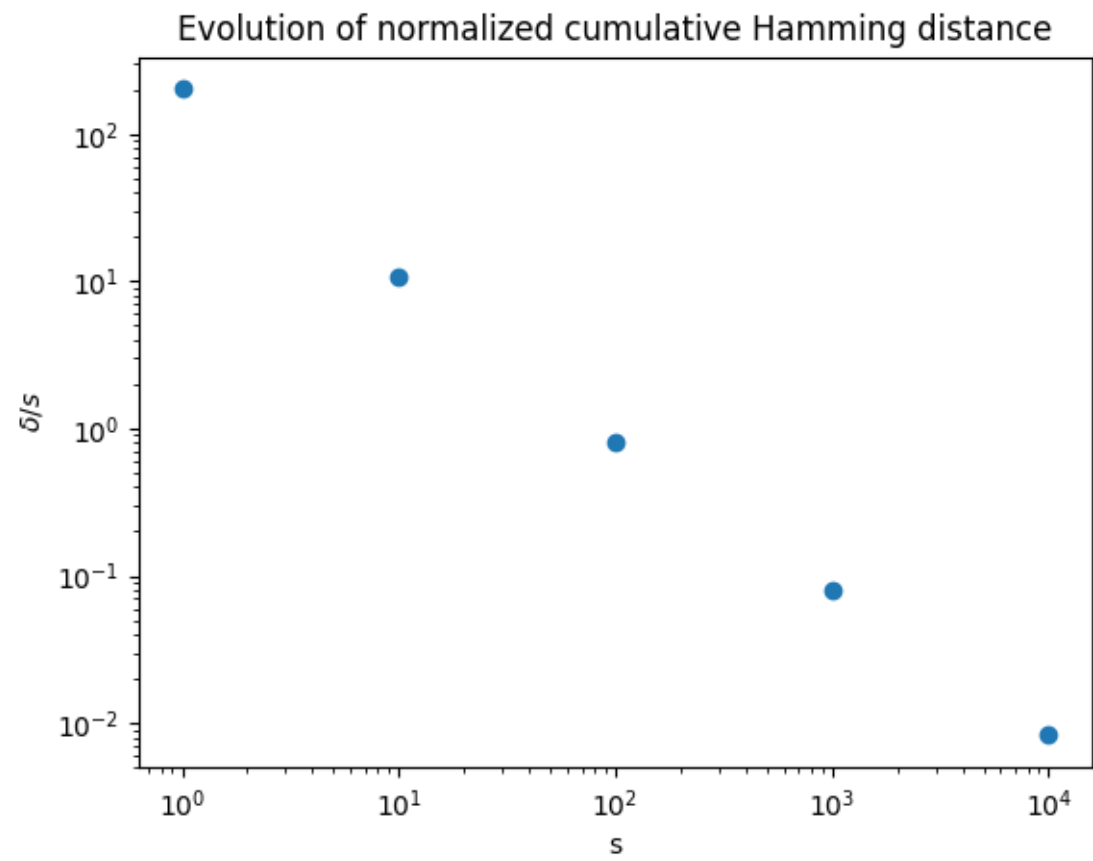
- ▶ 64 bits per hash
- ▶ s hashes per sketch
- ▶ $|\mathcal{S}|$ sketches in \mathcal{S}

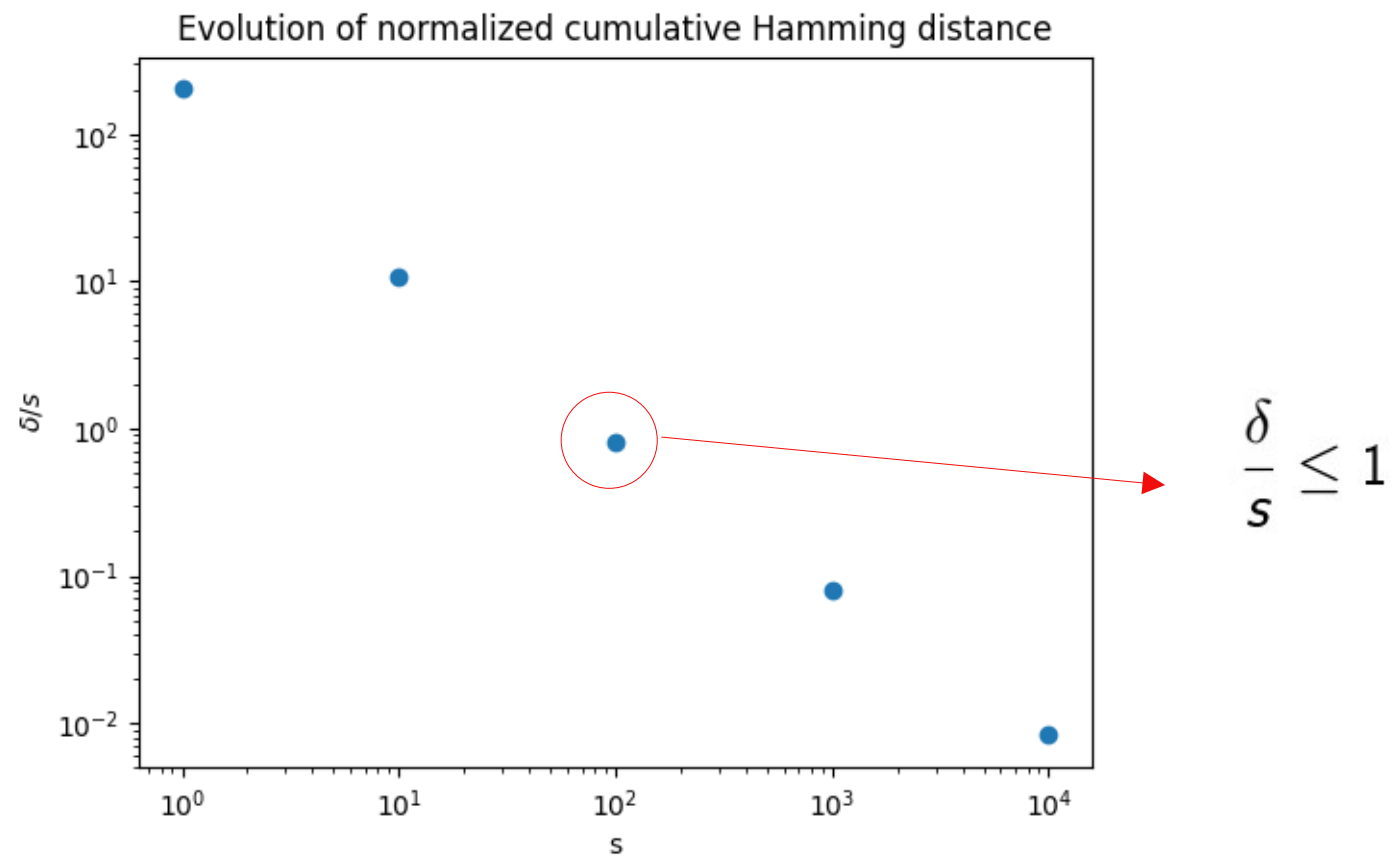
Compression ratio r :

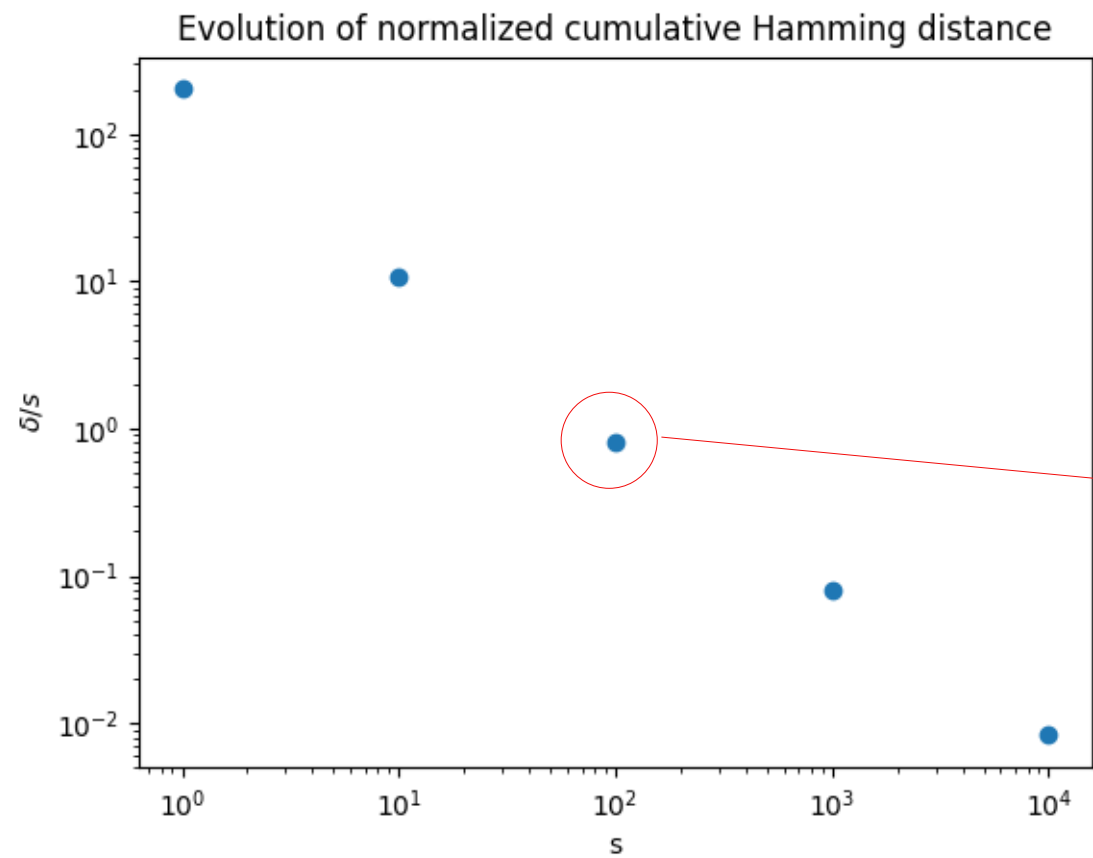
$$\begin{aligned} r &= \frac{|Input|}{|Output|} \\ &\geq \frac{64 \cdot |\mathcal{S}| \cdot s}{2 \cdot 64(\delta + s)} \\ &\geq \frac{|\mathcal{S}|}{2(1 + \frac{\delta}{s})} \end{aligned}$$

Experimental evaluation

- Genomes to sketch : neisseria_gonorrhoeae__01.tar.xz
in part 54 - <https://zenodo.org/records/15367750>
- Values of s : 1, 10, 100, 1000



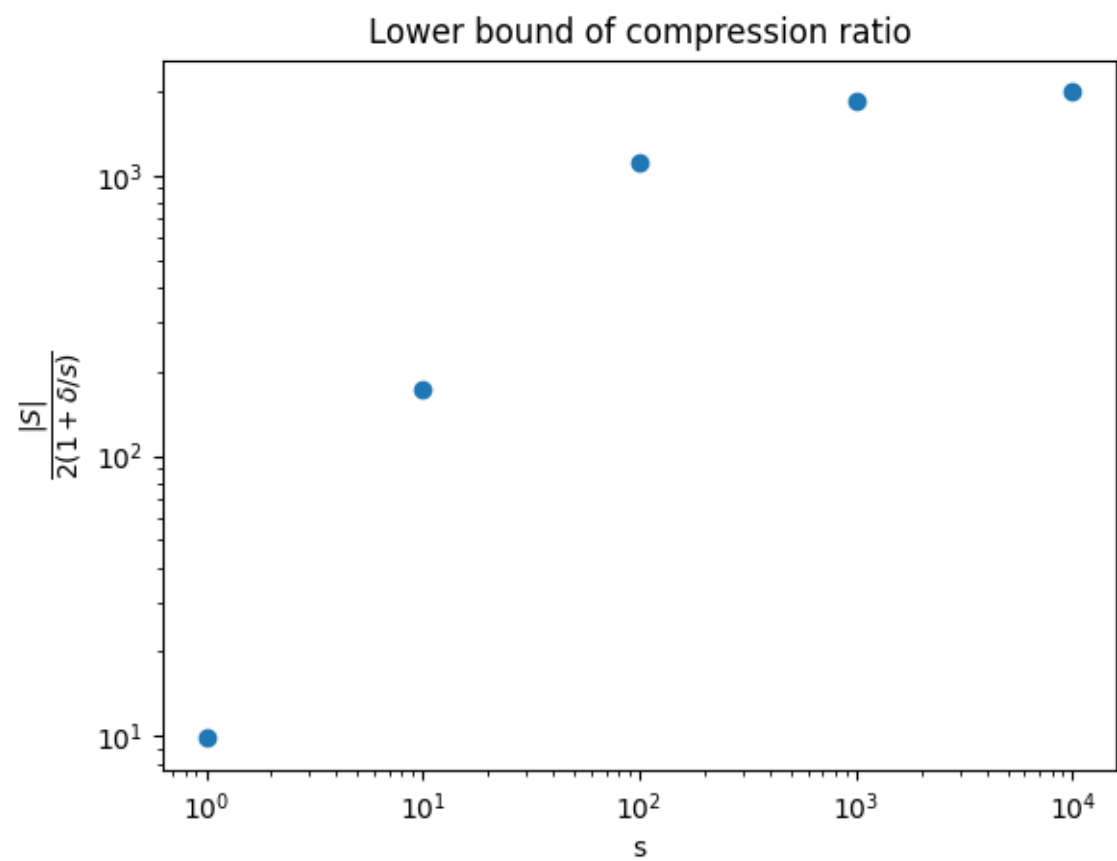




$$\frac{\delta}{s} \leq 1$$

$$r \geq \frac{4000}{4} \geq 1000$$

!!!



Conclusions :

- ▶ Compression rate of *RLE* increases with s for ngono
- ▶ Much better compression rates than with XZ/GZIP !

To do :

- ▶ Look into biological model of $\frac{\delta}{s}$
- ▶ Finer lower bound of r
- ▶ Compress \mathcal{S}_U (Elias Fano ?)
- ▶ Implement compression