

Robust image coding on synthetic DNA: Reducing sequencing noise with inpainting



Eva Gil San Antonio

Mattia Piretti

Melpomeni Dimopoulou

Marc Antonini



Université Côte d'Azur, CNRS, I3s, France



THE CHALLENGE:

- Data explosion → 80% of the worldwide data is rarely accessed ("cold")
- Conventional storage devices have a limited life-span (15-20 years)

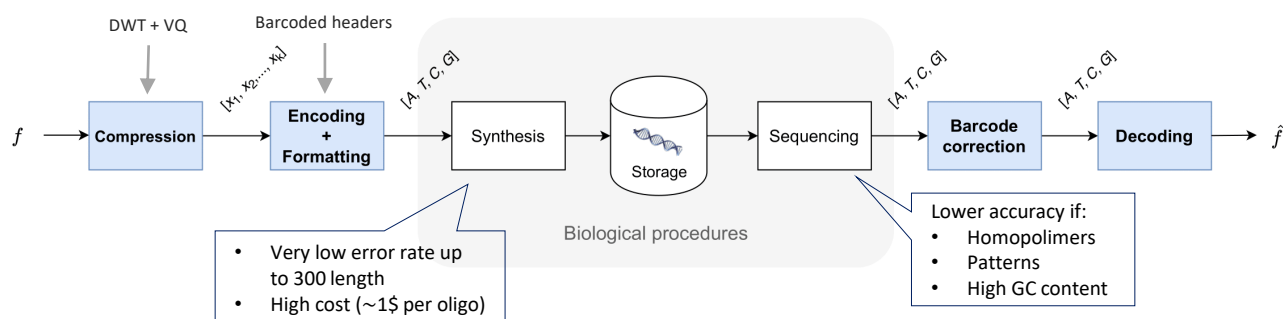
DID YOU KNOW THAT DNA...?

- Allows the storage of information at a high density (500 GB/mm³)
- Allows to preserve any information for hundreds of years without any loss if stored under the optimal conditions.

OUR CONTRIBUTION

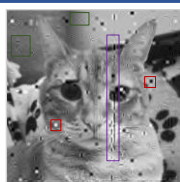
Despite being a promising solution, **DNA data storage** faces two major obstacles: the large cost of synthesis and the high error rate introduced during **nanopore sequencing**. While most of the works focus on adding redundancy for error correction, this work combines noise resistance to minimize the impact of the errors in the decoded data and post-processing based on **Texture Synthesis** to further improve the quality of the decoding.

Proposed workflow for DNA data storage:

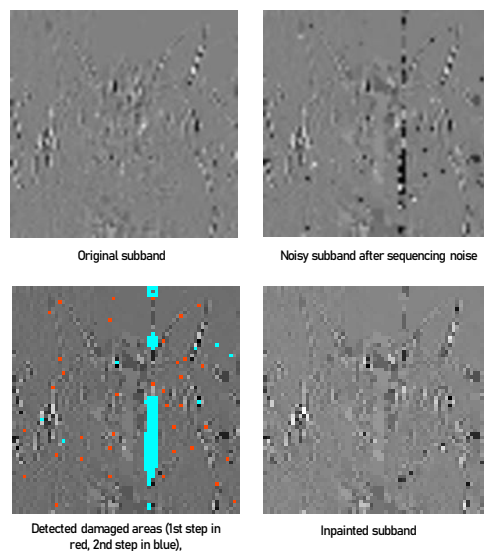


UNLIKE IN MOST INPAINTING SCENARIOS...

- DNA noise involves substitutions, insertions and deletions of nucleotides
- Noise applied on each subband separately → Heterogeneous damage
- Large **spots** (Low freq. subbands), **patterns** (High freq. subbands) and **lines** (lost DNA strands due to errors in the header)



Visual results on a DWT subband



AUTOMATIC DETECTION OF THE DAMAGED AREAS

2-step algorithm (performed in each subband separately):

- Detection of errors in single pixels (substitutions) → *Deviation of each pixel and its neighbours*

$$\begin{aligned} \tau_1 &: \text{phase 1 threshold} \\ N_p &: \text{neighborhood of the pixel } p \\ I &: \text{damaged image} \\ O_1 &: \text{Output mask phase 1} \end{aligned} \quad \forall p_{(x,y)} \in I, \text{ if } \begin{cases} \sqrt{\frac{(p_{(x,y)} - N_p)^2}{S_p}} \geq \tau_1, O_1(x,y) = 1 \\ \leq \tau_1, O_1(x,y) = 0 \end{cases}$$

- Detection of damaged neighbourhoods (indels) → *Internal variance of the neighbourhoods*

$$\begin{aligned} \tau_2 &: \text{phase 2 threshold} \\ S &= \sigma(I) \\ O_2 &: \text{Output mask phase 2} \end{aligned} \quad \forall p_{(x,y)} \in I, \text{ if } \begin{cases} \frac{S_p}{S} \geq \tau_2, O_2(x,y) = 1 \\ \leq \tau_2, O_2(x,y) = 0 \end{cases}$$

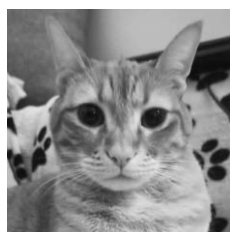
The result of the 2-step algorithm is a binary mask. The detected damaged regions are then corrected using Texture Synthesis*

*A. Criminisi, P. Pérez, and K. Toyama, "Region filling and object removal by exemplar-based image inpainting," *IEEE Transactions on Image Processing*, vol. 13, no. 9, pp. 1200–1212, 2004.

THE RESULTS

Simulation of sequencing noise*

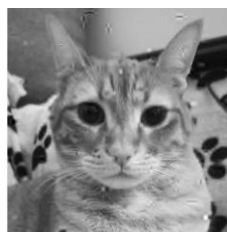
- Substitution and indel rates:
 - 2.3% deletions
 - 1.01% insertions
 - 1.5% substitutions
- 80% of the noise concentrated in the first and last 20nt of each oligo
- 200 noisy copies of each input oligo
- Headers encoded using barcodes
- Consensus based on majority voting



Quantized image without sequencing noise
Compression ratio = 4.9708 bits/nt
PSNR = 48.12 dB
SSIM = 0.991



Visual impact of the sequencing noise in the decoded image
PSNR = 36.2 dB
SSIM = 0.92



Post-processed image
PSNR = 38.7 dB
SSIM = 0.94

*Nanopore noise rates adapted from bibliography

CONCLUSIONS

While most works up to the date focus on Illumina sequencing due to its higher accuracy, we introduce nanopore sequencing in our workflow, attempting to speed up the process as well as reduce its cost. The results of this study are very promising given the high error-rates imposed by the nanopore sequencers showing significant visual improvement. However, since this study provides results on simulated nanopore noise, a wet-lab experiment is a priority future step to verify in practice the efficiency of the proposed encoding.