



HAL
open science

Image coding algorithm for DNA data storage combining JPEG and autoencoders

Xavier Pic, Marc Antonini

► **To cite this version:**

Xavier Pic, Marc Antonini. Image coding algorithm for DNA data storage combining JPEG and autoencoders. Munich Workshop on Coding and Cryptography, Jun 2022, Munich, Germany. hal-03710257

HAL Id: hal-03710257

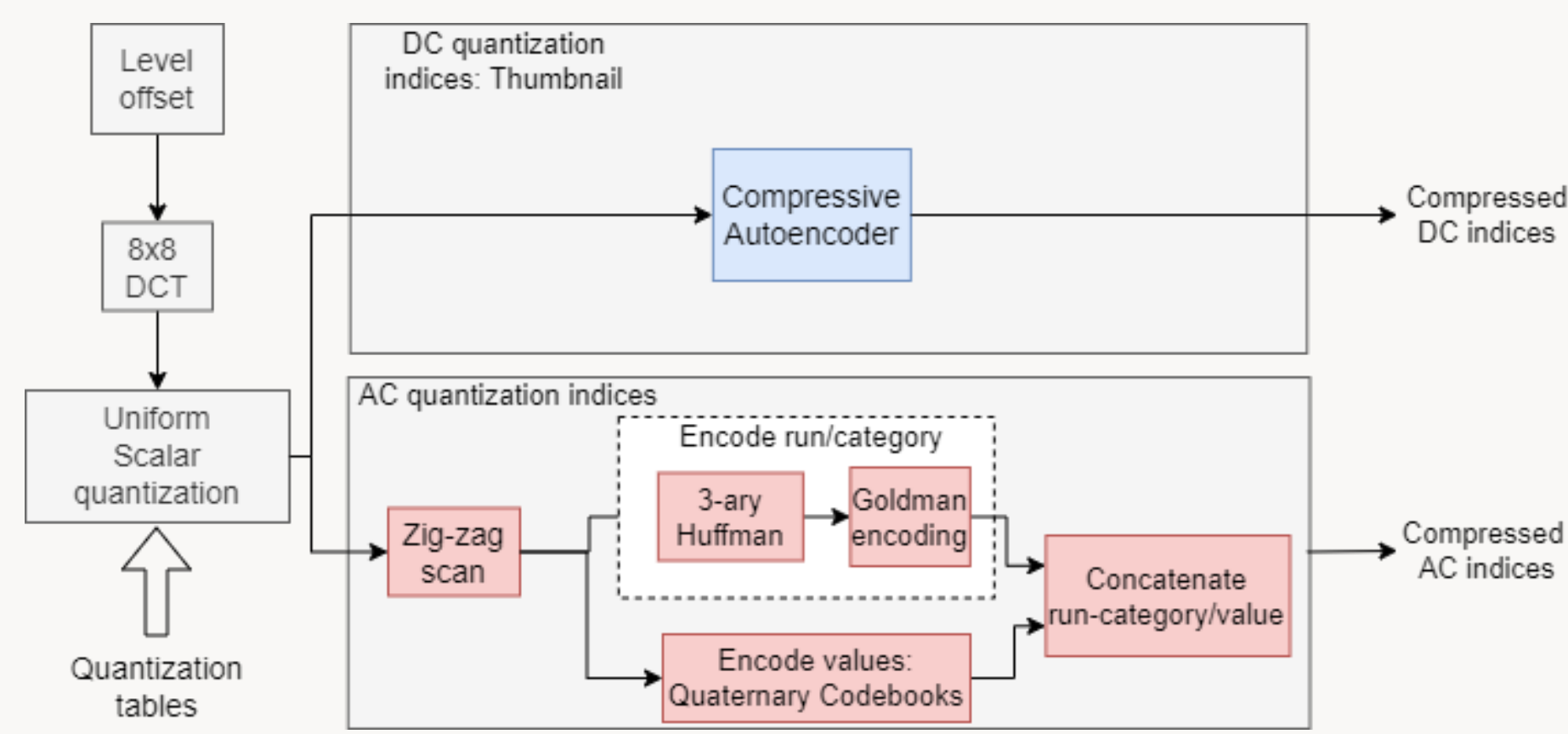
<https://hal.science/hal-03710257>

Submitted on 30 Jun 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Proposed workflow



- Based on the work of Melpomeni et al.[1]
- Block DCT based encoding into a DNA sequence
- DNA: sequence of nucleotides A, T, C and G :
AACTCAGCATGCAGGG...
- Formatted oligos synthesized into DNA molecules
- DNA molecules sequenced to retrieve the oligos
- Decoding the oligos to retrieve the original image

JPEG thumbnails

- Critical information
- Thumbnail can be encoded separately

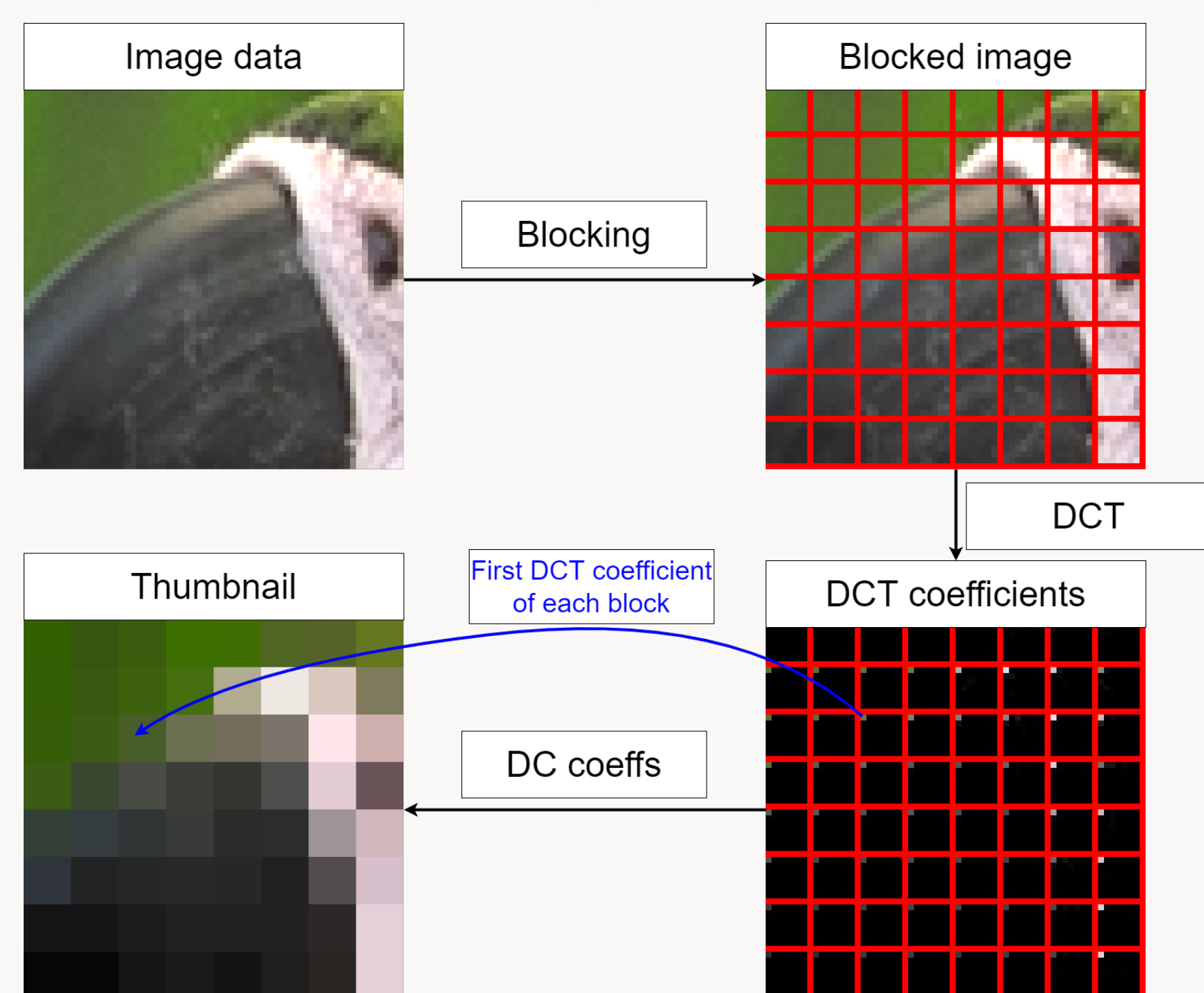
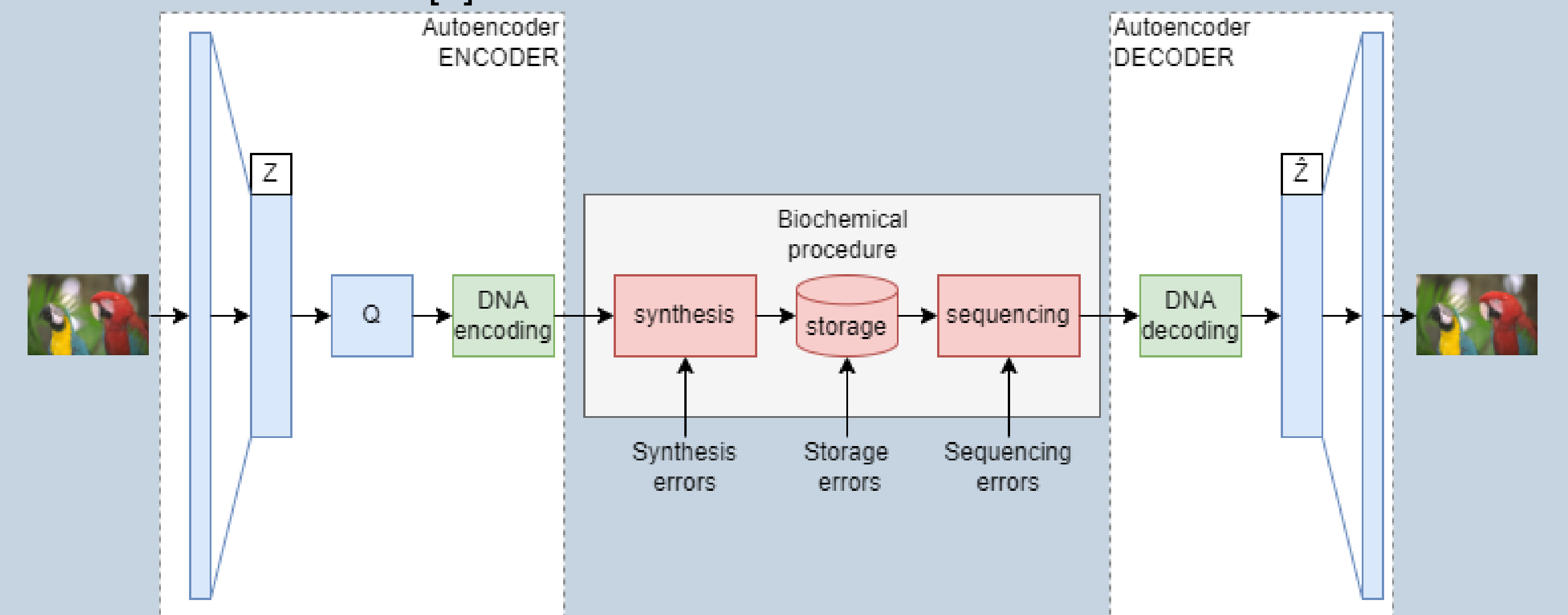


Figure 1. JPEG procedure and the thumbnail

Compressive autoencoders

- Based on Theis et al.[2]'s model



- Trained to adapt to the DNA storage channel errors:

Original	ATACTAGCT
Substitution	ATACAAGCT
Insertion	ATACGTAGCT
Deletion	ATAC_AGCT

- Reconstruct image from noised data

Data Compression

- Thumbnail (DC coefficients): Compressive Autoencoder
- AC coefficients : Goldman-based [3] Variable Length coding

Image reconstruction

Image retrieved by decoding and recombining the Thumbnail (DC coefficients) and the AC coefficients and applying inverse DCT to it.

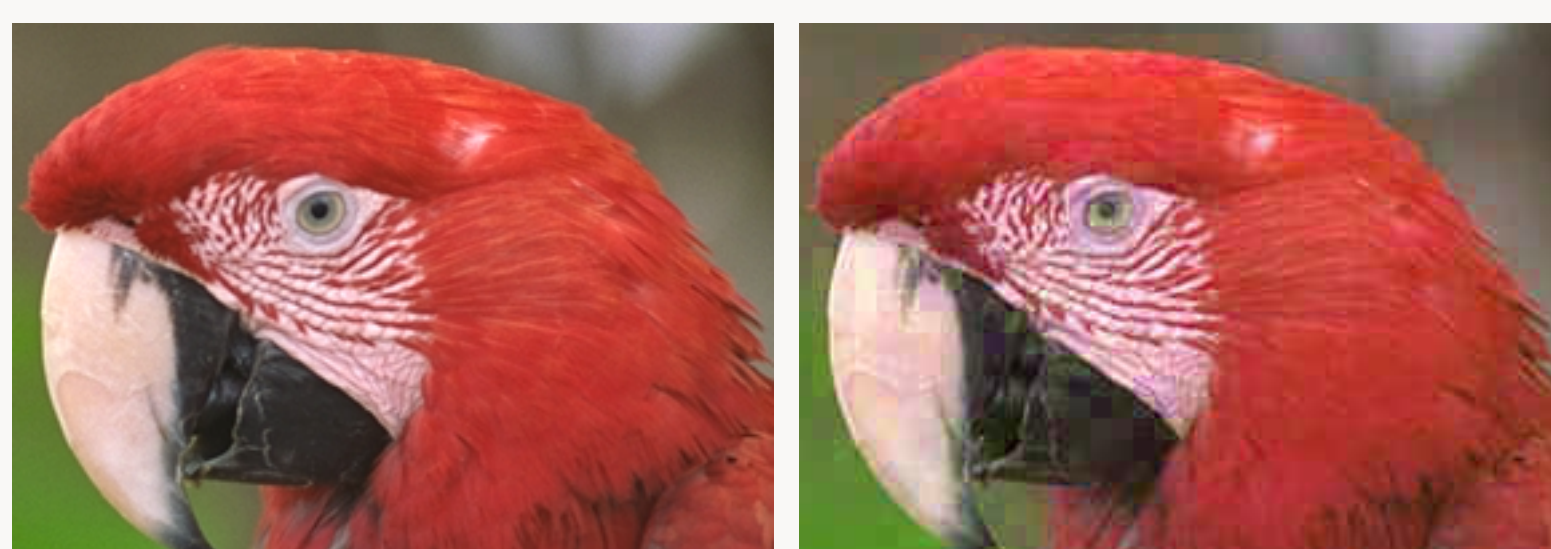


Figure 2. Original and Reconstructed image at 37.2 dB, compression rate: 33 bits/nt

References

- [1] M. Dimopoulou, M. Antonini, P. Barbry, and R. Appuswamy. "A biologically constrained encoding solution for long-term storage of images onto synthetic DNA". In: *European Signal Processing Conference (EUSIPCO)* (2019).
- [2] L. Theis, W. Shi, A. Cunningham, and F. Huszár. "Lossy image compression with compressive autoencoders". In: *International Conference on Learning Representations* (2017).
- [3] N. Goldman, P. Bertone, S. Chen, C. Dessimoz, E. M. LeProust, B. Sipo, and E. Birney. "Towards practical, high-capacity, low-maintenance information storage in synthesized DNA". In: *Nature* 494.7435 (2013), p. 77.

Model mathematical definition

$$\hat{I} = g(\overline{\beta_{DNA}}(\alpha_{DNA}(Q(f(I)) + \epsilon))) \quad (1)$$

- f : encoding part of the autoencoder
- g : decoding part of the autoencoder
- Q : quantization ($Q(z) = q \times \lfloor \frac{z}{q} + \frac{1}{2} \rfloor$)
- f_{DNA} : DNA encoding function
- g_{DNA} : DNA decoding function
- ϵ : Noise generated by the DNA storage channel

Training process

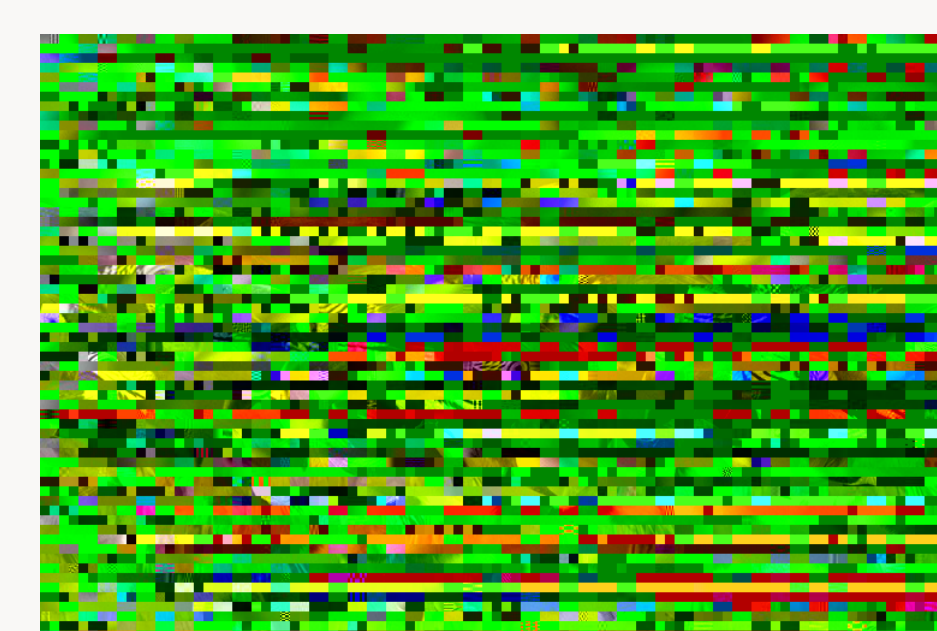
- Training Loss: $L = \|I - \hat{I}\|^2 + \lambda H(Q(Z))$
 - $\|\cdot\|^2$: Mean Squared Error
 - Z : latent space ($Z = f(I)$)
 - H : entropy $H(Q(Z)) = -\sum_{i=1}^n Pr\{Q(z_i)\} \log_2 Pr\{Q(z_i)\}$
 - λ : weight parameter
- Optimizer: Adam
- Training Dataset: Thumbnails of Flickr30k (~30k images from social media)
- Validation Dataset: Thumbnails of Kodak Image dataset (24 images)
- DNA storage channel noise model

Results

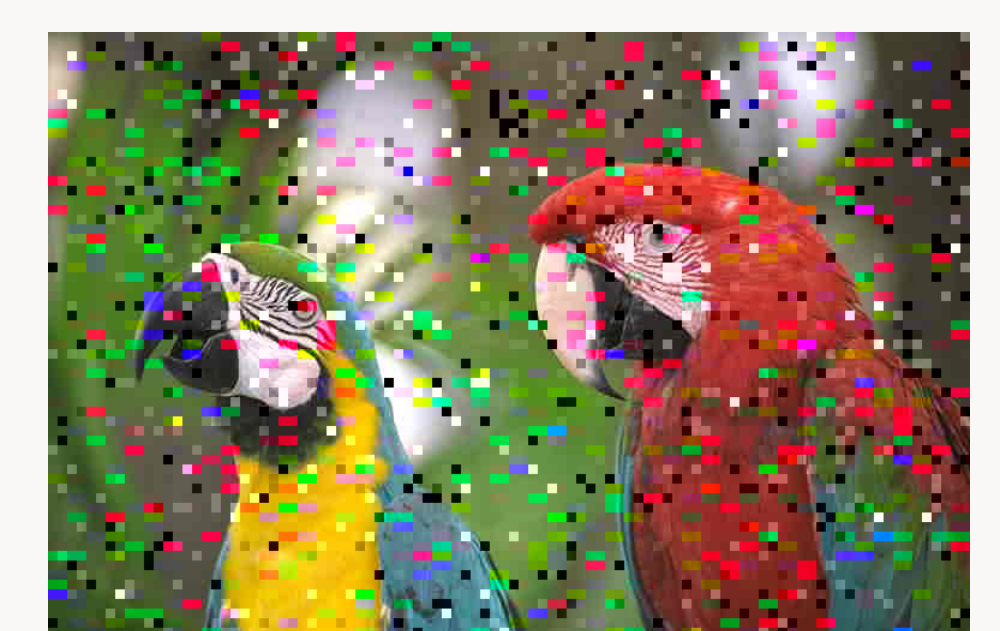
- The compression in terms of performance is lower than the jpegdna one
- Achieving lower bit-rates is harder because of the size of the compressed representation of the thumbnail
- Better performance on noisy channels
- More resistant to noise on the critical data managed by the autoencoder
- More adapted to a noisy channel like synthetic DNA data storage
- 1.5% substitution rate corresponds to the Nanopore sequencing method's substitution rate



JPEG DNA, no noise



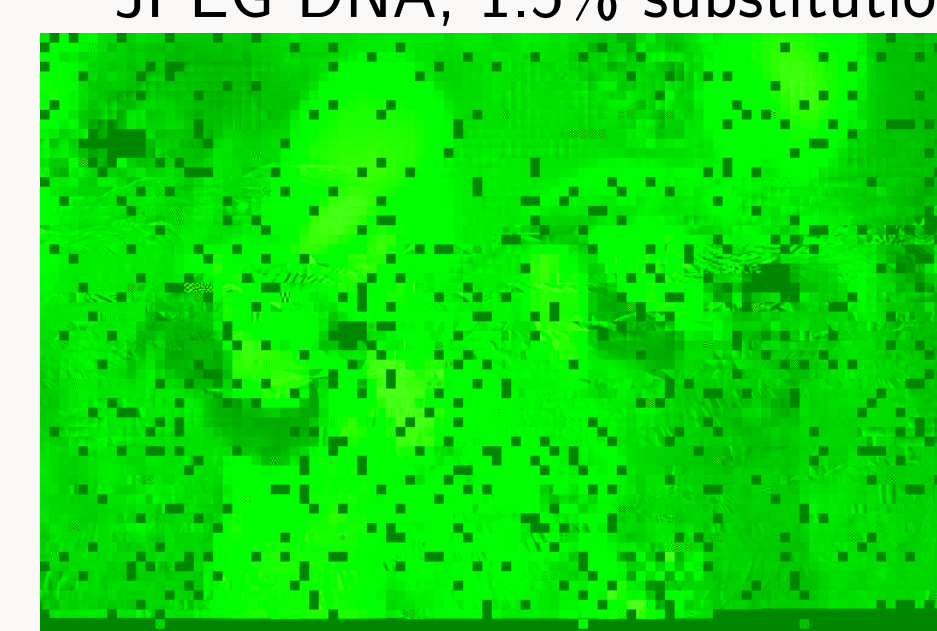
JPEG DNA, 1.5% substitution



Noised thumbnail only



Our method, no noise



Our method, 1.5% substitution



Our method, noised thumbnail only