



Unlock the Future of CyberStorage

# HPC Facility Featuring Diverse Acceleration Engines for DNA Data Storage



A. Marelli, R. Micheloni, C. Zambelli, S.F. Schifano, and E. Calore

- Research team members
- HPC facilities @ UniFE and INFN
- HPC for DNA-based data storage
- DNAssim on heterogeneous CPUs and GPUs
- ECC and Edit distance on GPUs and FPGAs
- Conclusion

# AGENDA



**RESEARCH TEAM**



**Università  
degli Studi  
di Ferrara**



**Istituto Nazionale di Fisica Nucleare**

- Prof. Cristian Zambelli – Università degli Studi di Ferrara
- Prof. Sebastiano Fabio Schifano – Università degli Studi di Ferrara and INFN
- Dr. Enrico Calore – INFN
- Dr. Rino Micheloni – Avaneidi
- Dr. Alessia Marelli – Avaneidi
- Ph.D. + M.Sc. students
- Advancing DNA-based data storage research together

# TEAM MEMBERS

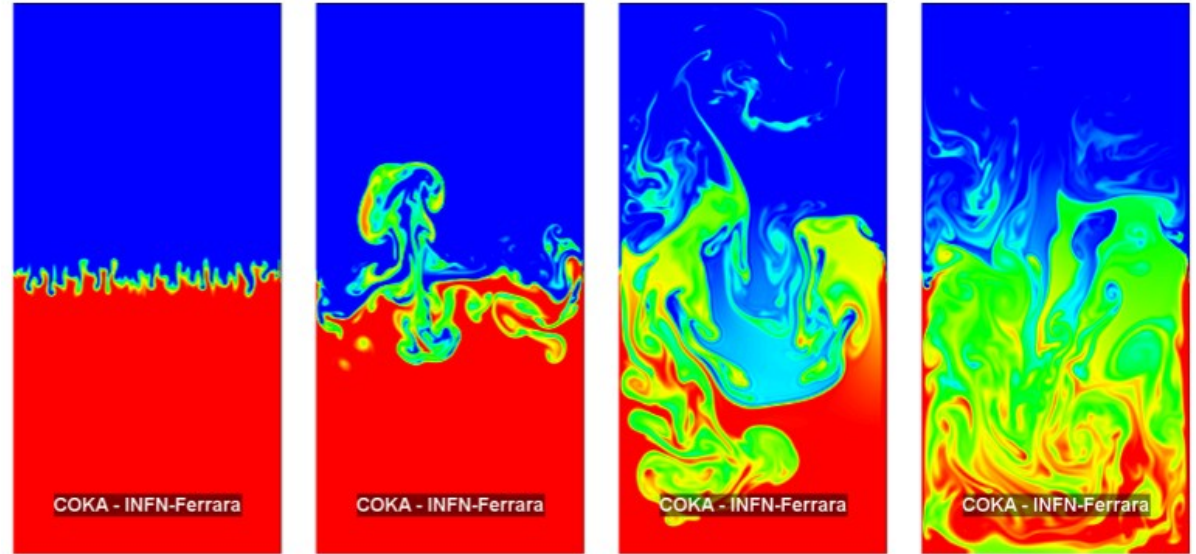


**HPC FACILITIES**





- Solution of complex physics problems (e.g., LBM, LQCD, etc.)
- AI/ML tasks (e.g., computer vision)



- The High-Performance Computing (HPC) facility in Ferrara (Italy)
- Computing on Kepler Architectures (COKA)
- Started in 2014 by the University of Ferrara and the INFN Ferrara

# HPC @ UNIFE AND INFN

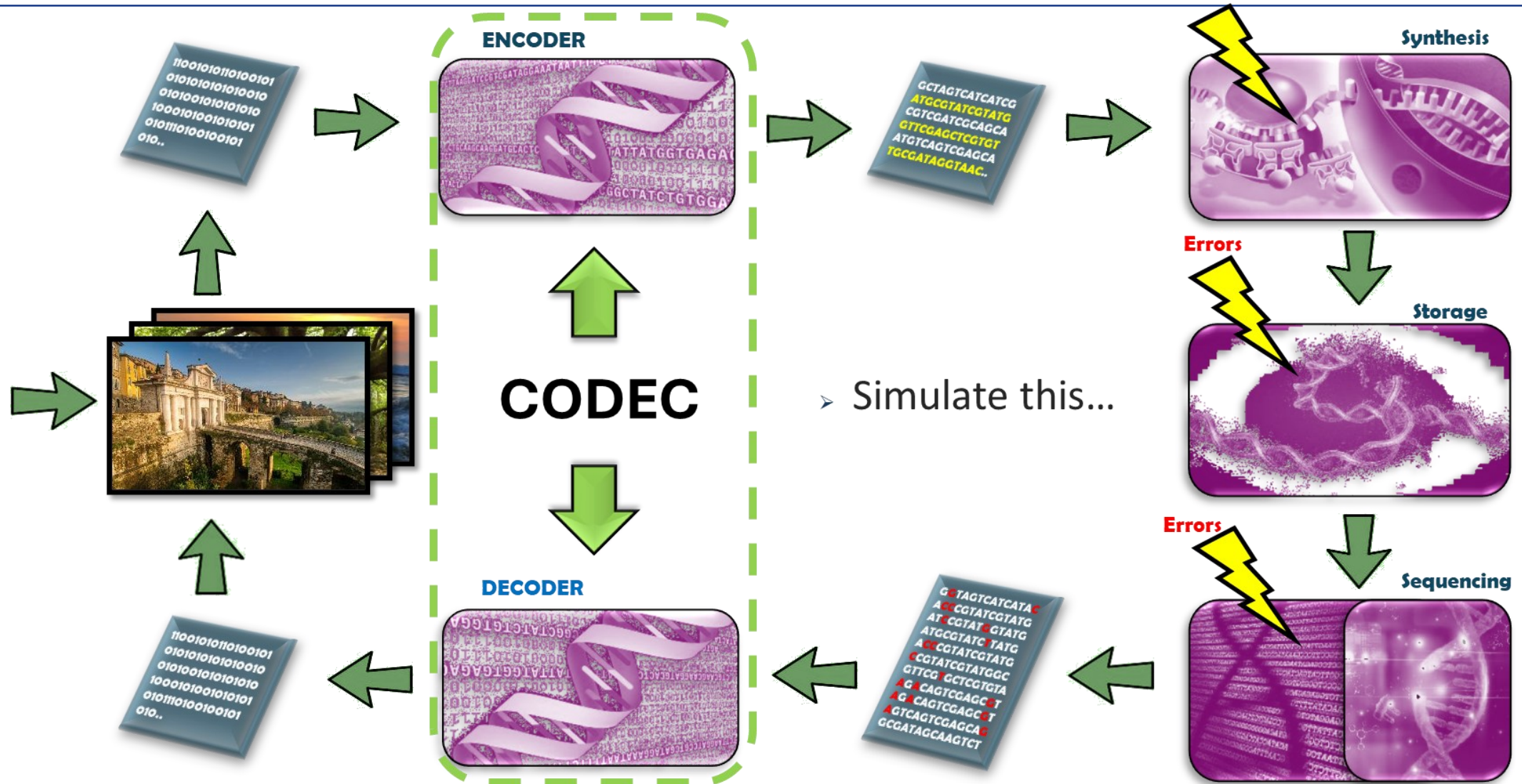


**HPC FOR DNA-BASED  
DATA STORAGE**

- When validating an algorithm, it is necessary to simulate a large amount of data
  - With a software-only approach (...or pure CPU...) we may not go too far!
  - CPU-only? Not for us! Work with GPUs and FPGAs
- If the simulated algorithms or error models are highly complex
  - Have a facility capable of handling computationally intensive workloads
- In the case of DNA data storage, the goal is to emulate a real model
  - The model includes the replication process introduced by PCR and the IDS channel
  - This generates a large amount of data after sequencing

# WHY HPC IN DNA?





# THE DNA-BASED DATA STORAGE



- **COKA GPU nodes**

- ECC sims. + DNAsim

- **NVIDIA GH200 nodes**

- Edit distance
- Alignment filters

- **Skylake FPGA node**

- Edit distance



- The HPC facility has been revamped in 2023
- Providing a playground for DNA data storage research

# EXTENDING THE FACILITY





- HPC machine with 3 computing nodes
- NVIDIA Tesla T4 GPUs are coupled with AMD Xilinx Alveo U50
- Including Custom Avaneidi Computational Storage Drives for performance boost

# AVANEIDI POWER UP!

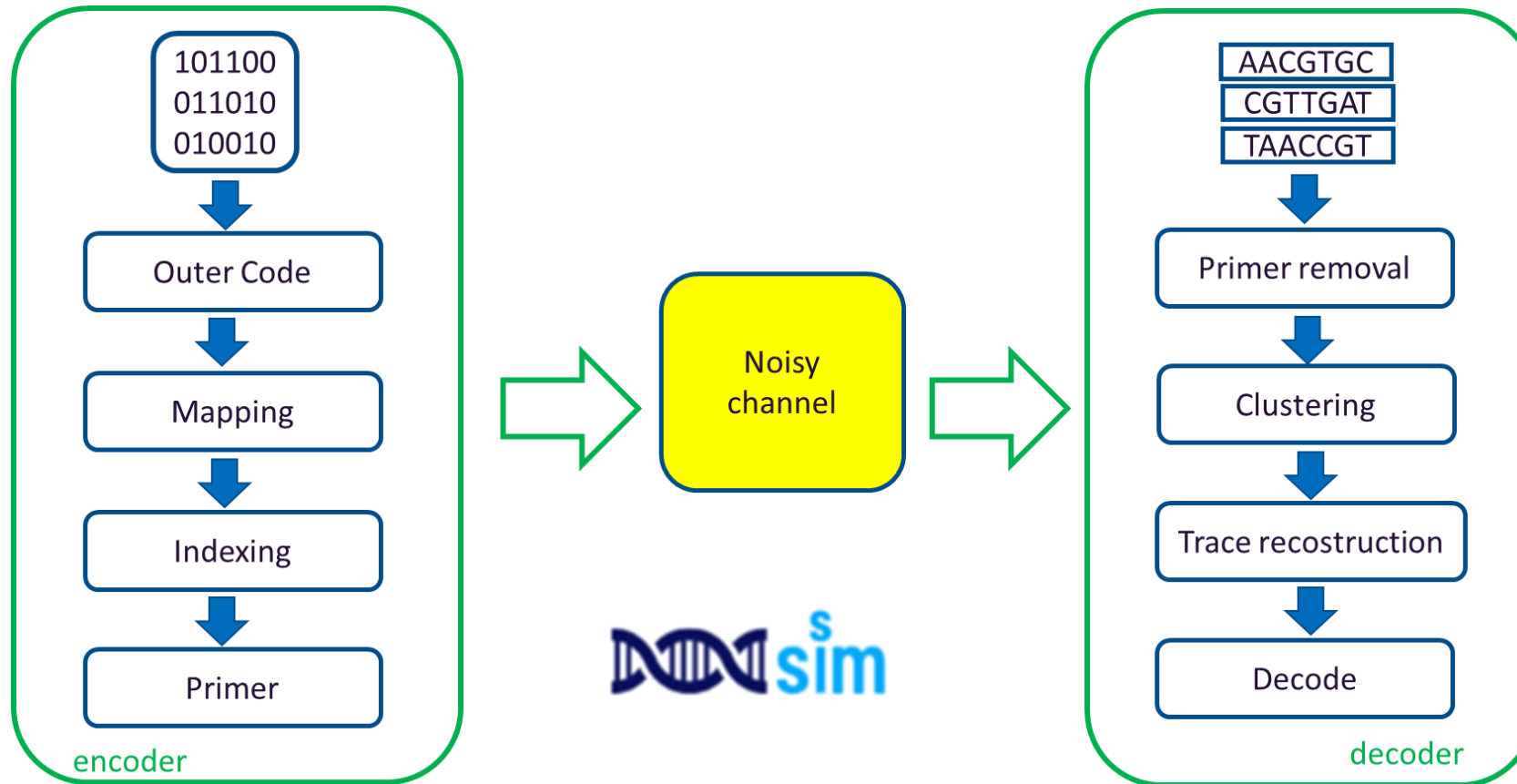
- AMD FPGA devices in Ferrara
  - Thanks to the AMD University Program (AUP)
- 2x Alveo U50 → This one is the target for DNA studies (ECC and Edit distance)
- 1x Alveo U55C
- 1x Alveo U250
- 1x Versal VCK5000
- 1x Samsung SmartSSD
- Future plans to include also the Alveo V80 in our facility

# RECONFIGURABLE COMPUTING



**DNASSIM®**

- Because of the intrinsic statistical behavior of the noise, a simulator is required
  - Figuring out the impact of encoding/decoding algorithms (CODEC noise-model tailored)

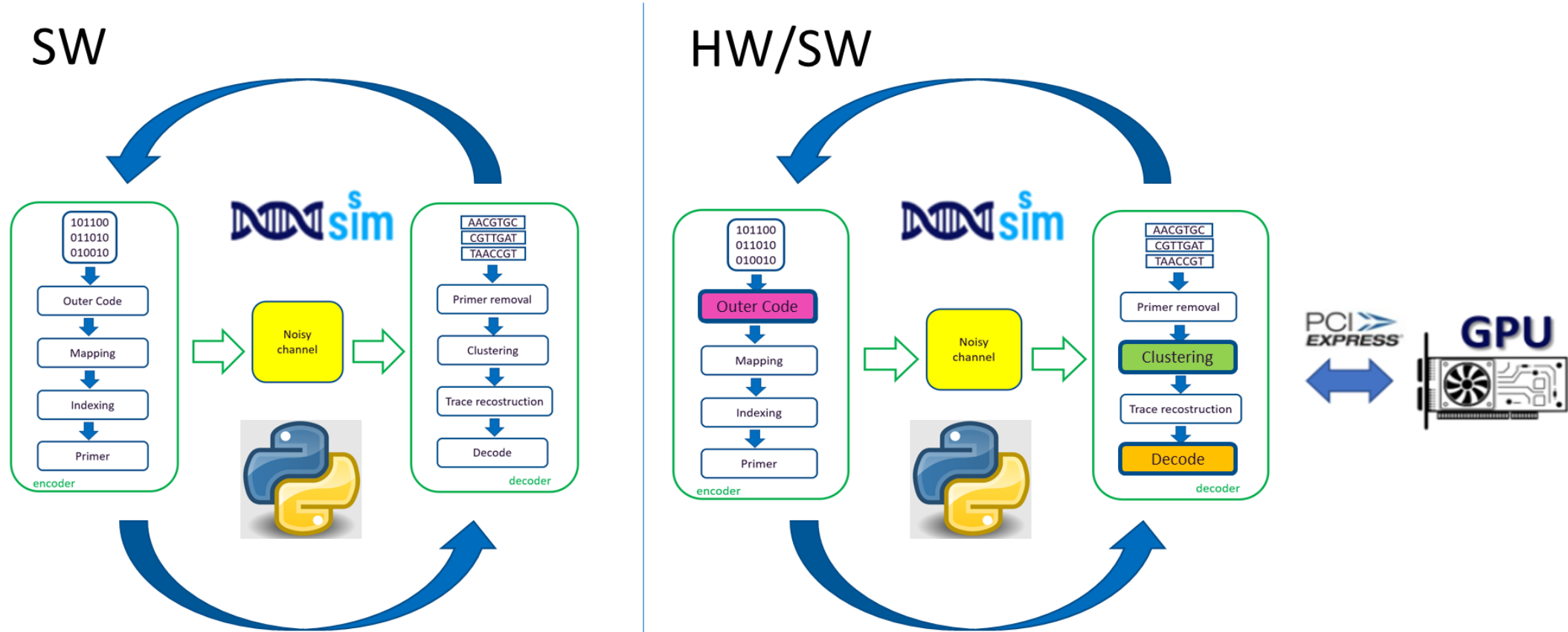


- Reprinted with permission from A. Marelli et al. paper with DOI. 10.3390/electronics12122621 under Creative Commons License (CC-BY-4.0)

# INTRODUCING DNASSIM®

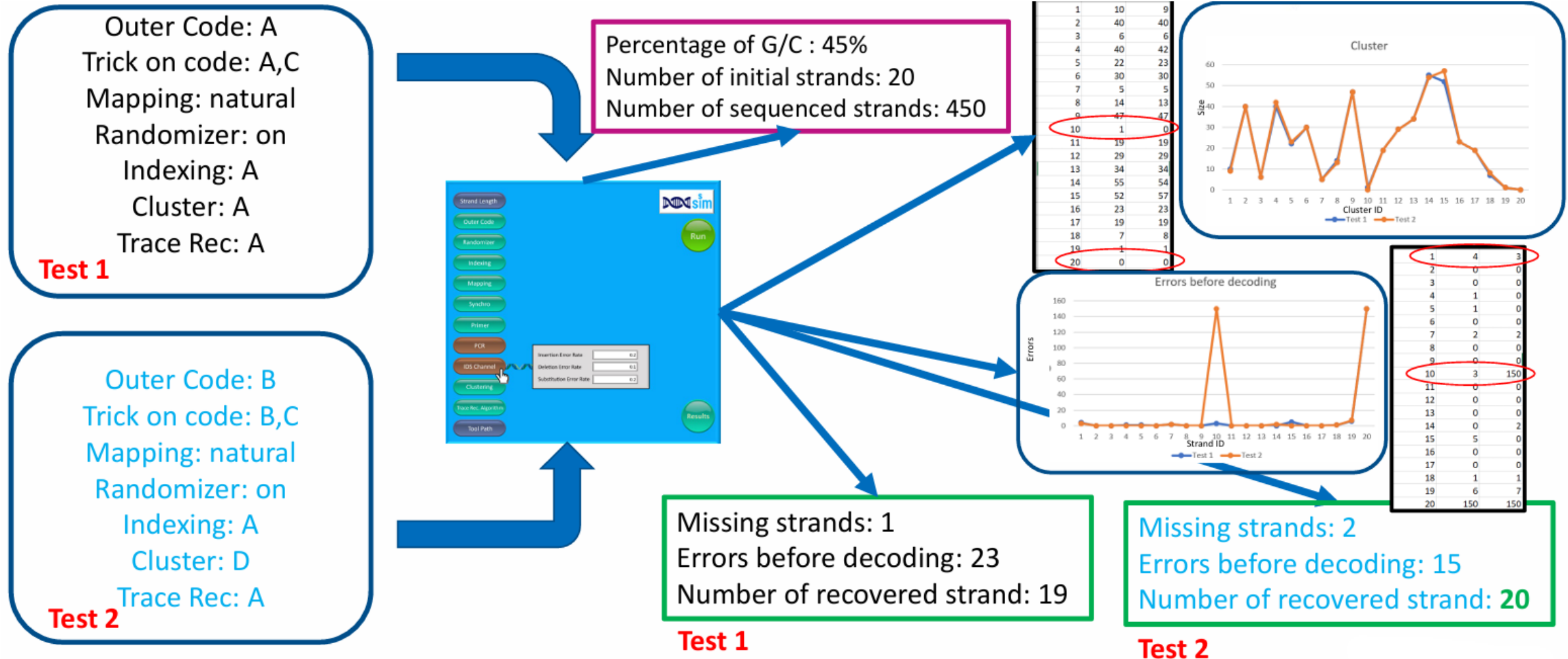


- The number and complexity of the steps involved in the DNA storing process is high
  - Number of simulations is huge, and a “pure software” simulator can easily run out of gas



# OFFLOADING SIMULATIONS ON GPU

➤ From A. Marelli et al., SDC 2022

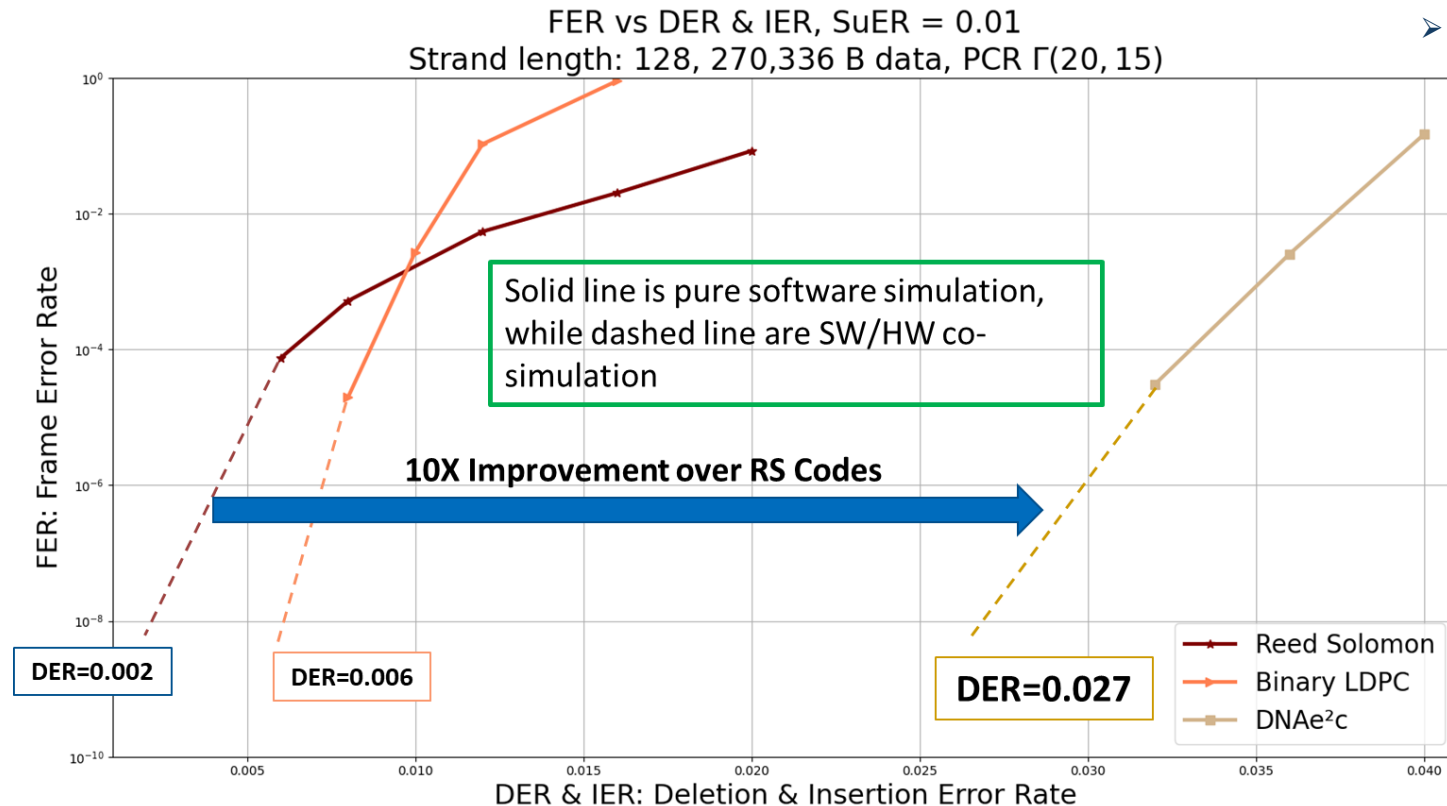


# SIMULATIONS CAPABILITY



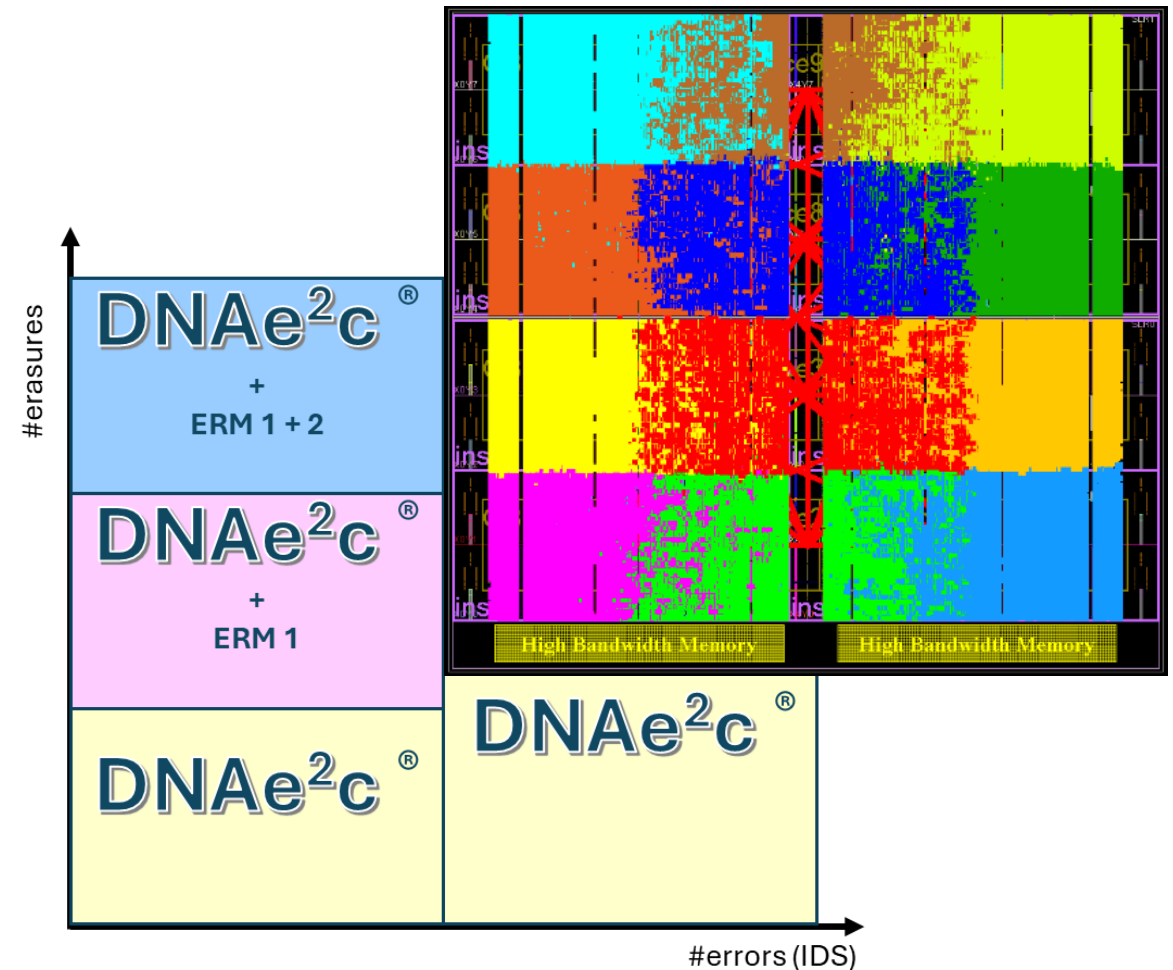
**ECC AND EDIT-DISTANCE**

- To evaluate different ECC it is necessary to go to very low Frame Error Rate
  - Simulate a huge amount of codewords.
- In the ECC case we use the GPU approach

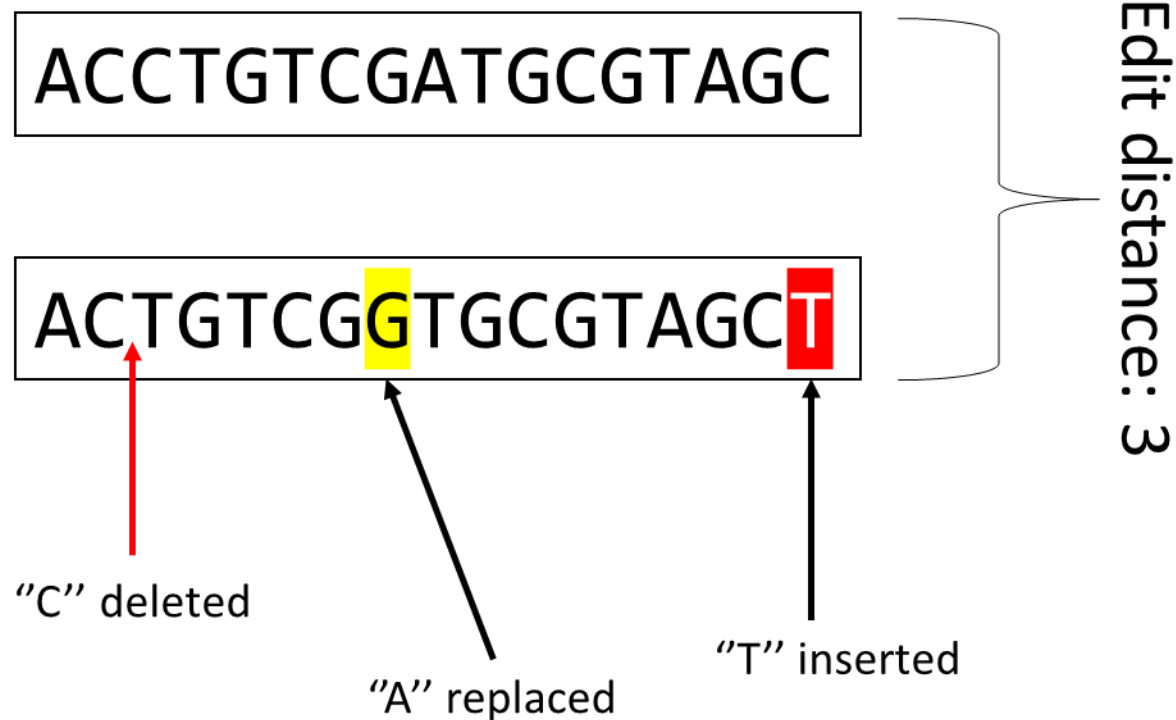


# ECC CASE STUDY

- DNAe2c<sup>®</sup> offers a complete solution suite tailored to specific error conditions
- If erasures increase, we can enhance reliability by adding ERM (Erasure Recovery Mechanism) or combining it with other correction techniques
- DNAe2c<sup>®</sup> has been simulated on AMD Alveo U50 Data Center Accelerator Cards



# DNA E2C<sup>®</sup>



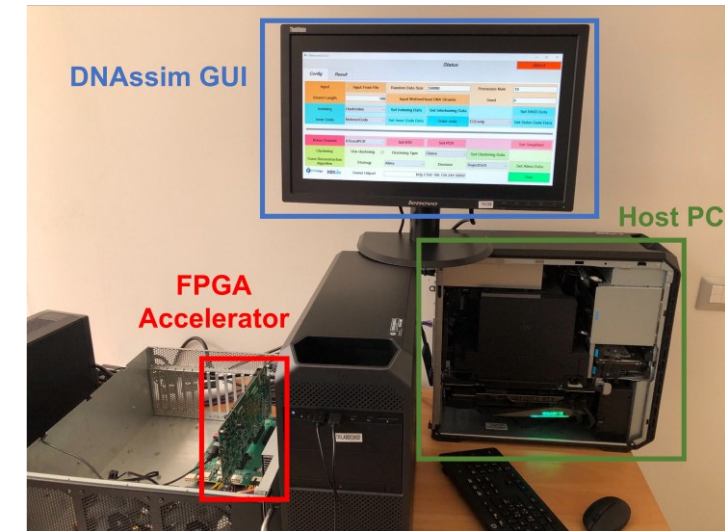
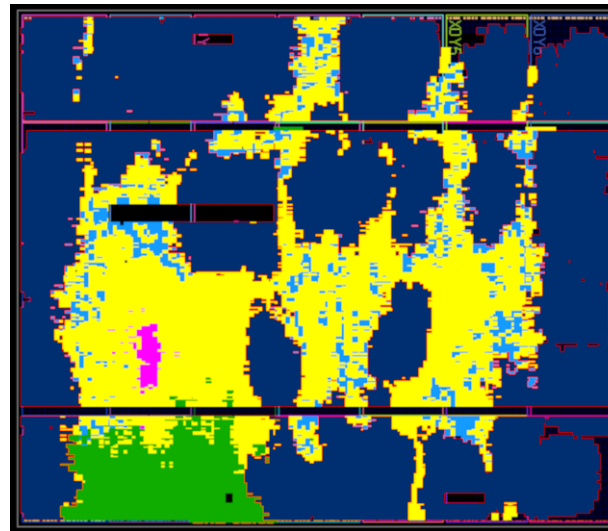
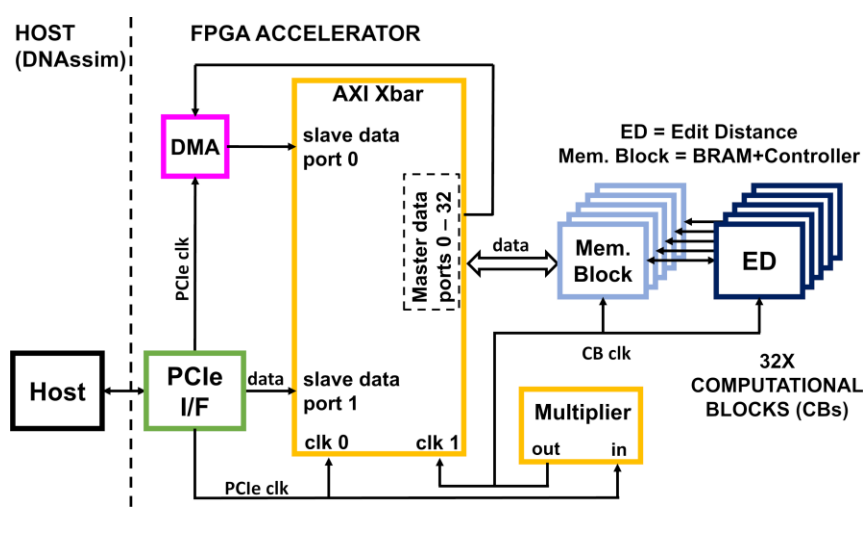
- In an IDS channel, a metric used is the Levensthein distance
- Algorithm to compute it and recovering messages can be much harder
  - When dealing with this distance

- To evaluate the edit distance, one can use a well-known dynamic programming algorithm
- ...a problem with  $O(n^2)$  complexity likely ( $n$  is the DNA string length)

# EDIT-DISTANCE



- Reprinted with permission from A. Marelli et al. paper with DOI. 10.3390/electronics12122621 under Creative Commons License (CC-BY-4.0)

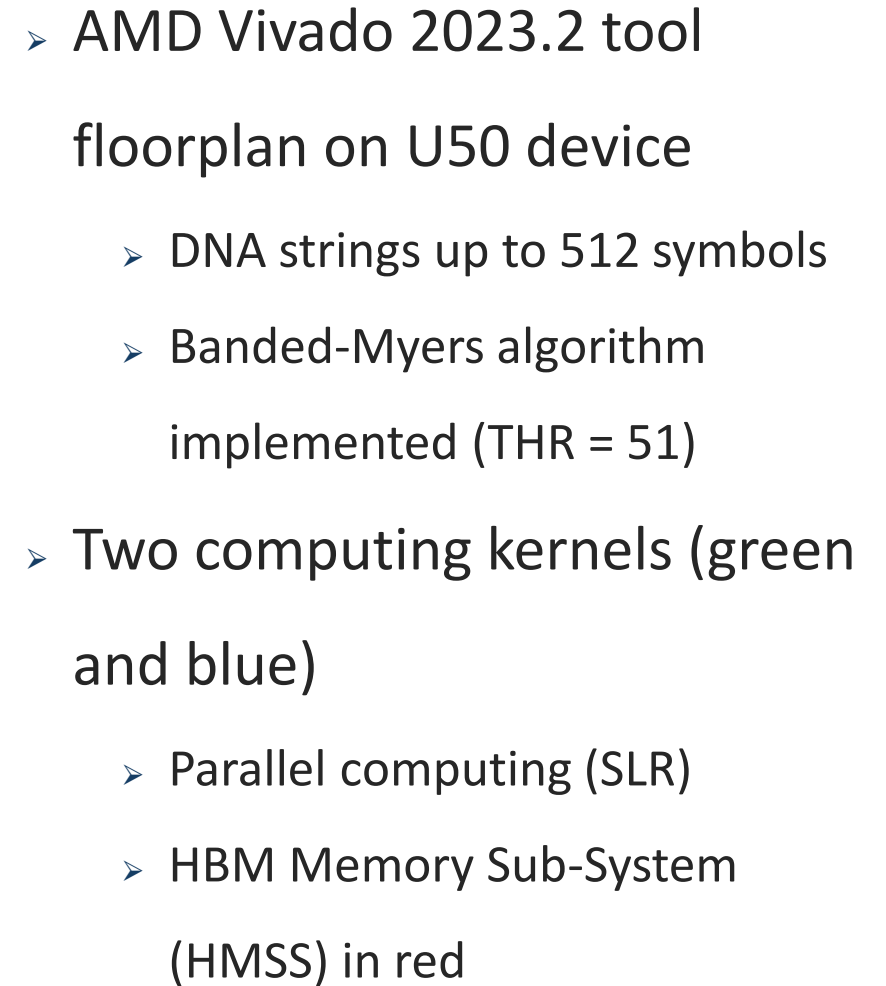


- Using a custom FPGA design to accelerate edit distance computation
  - Based on a Xilinx VC707 FPGA evaluation board (PCIe gen2 x4)
  - Up to DNA strings with 255 symbols
- Implementing the Ukkonen's version of the algorithm using HLS (slow, slow, and SLOW!)

# FIRST ATTEMPT ON FPGA

- The design parameter for an edit distance accelerator are many
  - The average DNA string length (NC)
  - The threshold used by the Edit distance algorithm (i.e., THR for Banded-Myers)
  - The number of parallel memory channels used by computing kernels
  
- ...which in turns affect
  - FPGA resource occupation (CLB, LUT, REG)
  - The achievable clock frequencies (DCLK, HCLK)
  - The acceleration performance measured in terms of Mpairs/s (MPS)
  - Energy-efficiency (MPS/J) and compute-efficiency ( $\epsilon$ )

# DESIGN PARAMETERS AND METRICS



# FPGA FLOORPLAN

- Computing- and energy-efficient edit distance algorithms
  - Combining metaprogramming and HLS to design and optimize C codes for FPGA devices
  - Achieving about 90% of computing efficiency and near 90% FPGA resources occupancy.
  - Delivering a peak throughput of 16.8 TCUPS and an energy efficiency of 46 Mpairs/Joule
- 
- S.F. Schifano et al. paper with DOI. 10.1016/j.future.2024.107591
  - Future Generations Computer Systems top-pick editor's choice paper

# RECOGNIZING TOP-CLASS METRICS



**CONCLUSIONS**

- DNA data storage is becoming a reality!
- While it presents challenges, it holds tremendous potential
- Having a simulator is essential to avoid conducting experiments that are costly in terms of money and time
- HW-SW co-simulation is essential for managing complex algorithms and large volumes of data
- An HPC facility is necessary to handle the data analysis required for enabling DNA-based data storage
- ...

# CONCLUSIONS





[avaneidi.com](http://avaneidi.com)



[marketing@avaneidi.com](mailto:marketing@avaneidi.com)



Avaneidi Spa

# THANK YOU