



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA



Future
Artificial
Intelligence
Research

October 19 - 25, 2024

**CHEP
2024**

Conference on Computing in High Energy and Nuclear Physics

October 19-25, 2024, Kracow, Poland



Real-Time implementation of Artificial Intelligence compression algorithm for High-Speed Streaming Readout signals.

Fabio Rossi (presenter), **Marco Battaglieri**
Istituto Nazionale di Fisica Nucleare
Genova (Italy)

Edoardo Ragusa, Paolo Gastaldo
SEALab Università di Genova
Genova (Italy)

Gagik Gavalian
Jefferson Lab
Newport News (Virginia)



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA



Future
Artificial
Intelligence
Research

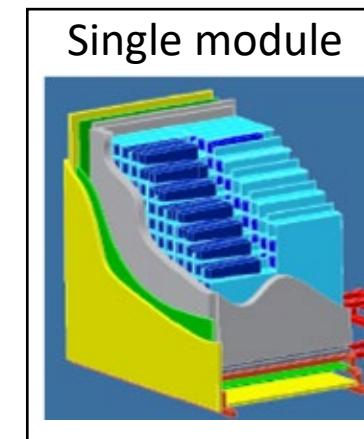
High Energy Physics Experiment: Beam Dump eXperiment (BDX)

Jefferson Lab

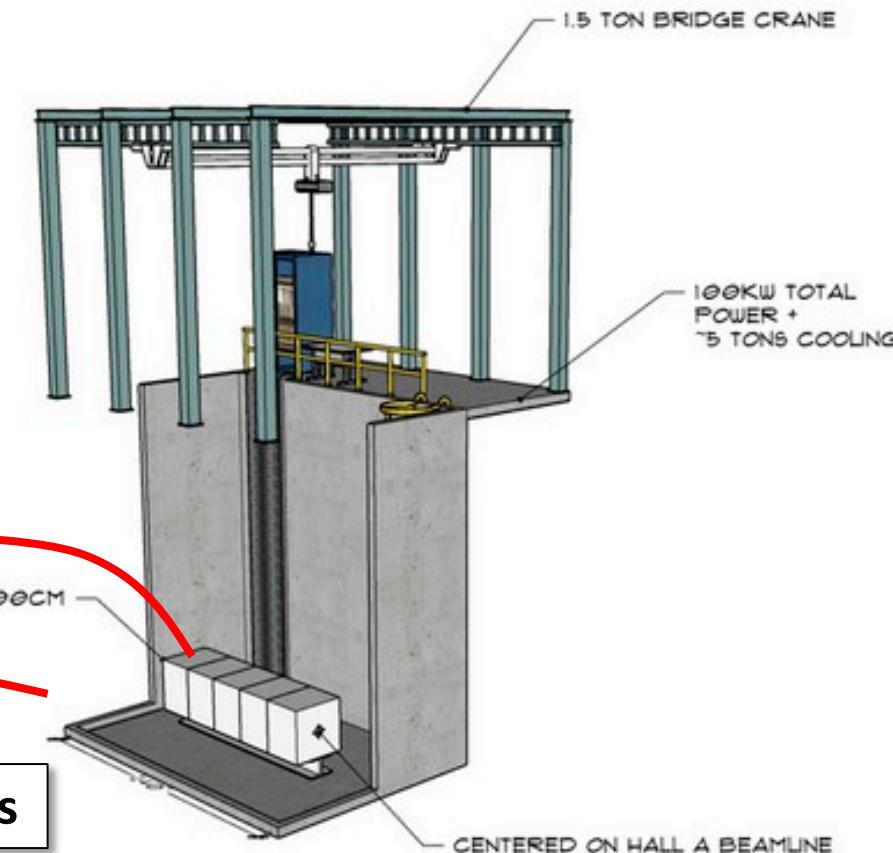


≈ 1000
Calorimeter channels
(30MB/s)

≈ 300
Veto channels
(500MB/s)

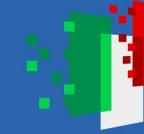


5 MODULES AT 80CM X 100CM X 100CM

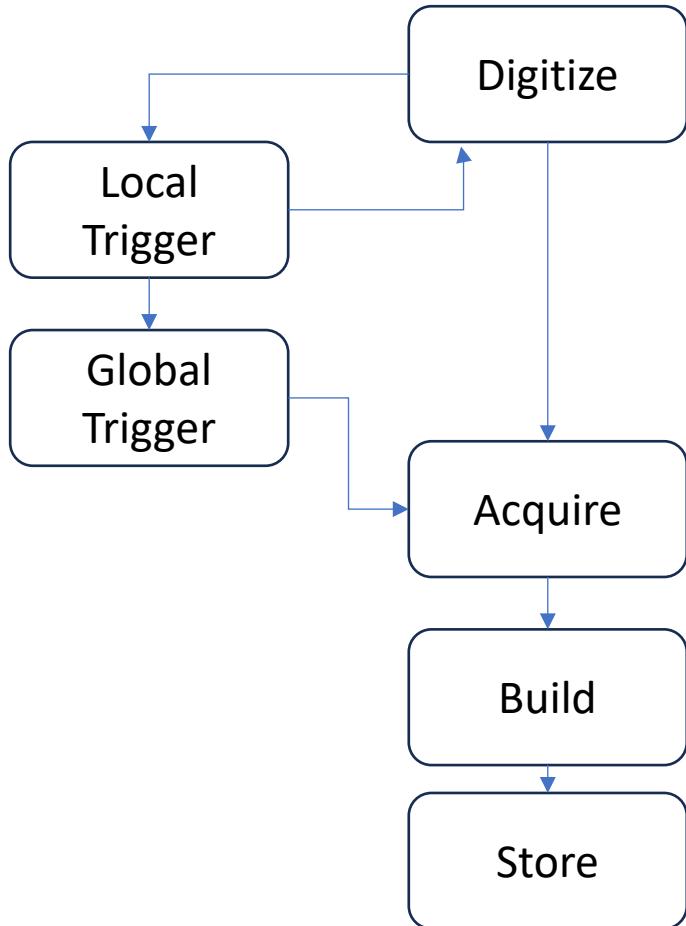


Very rare occurrence of Dark Matter events

Retrieved from: Battaglieri, M., et al. "Dark matter search in a Beam-Dump eXperiment (BDX) at Jefferson Lab." arXiv preprint arXiv:1607.01390 (2016).



Traditional triggered DAQ VS Streaming Readout SRO



Cons:

Only few information form the trigger.
Trigger logic difficult to implement and debug.
Not easy to adapt to different condition.

Pros:

It works reliably.

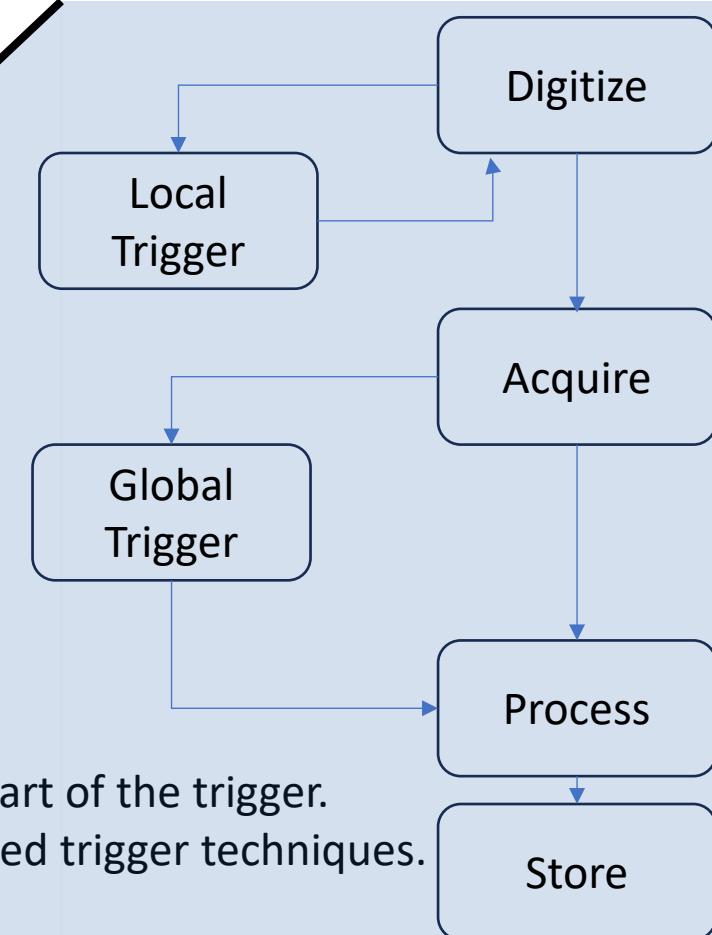
Triggered
Streaming

Cons:

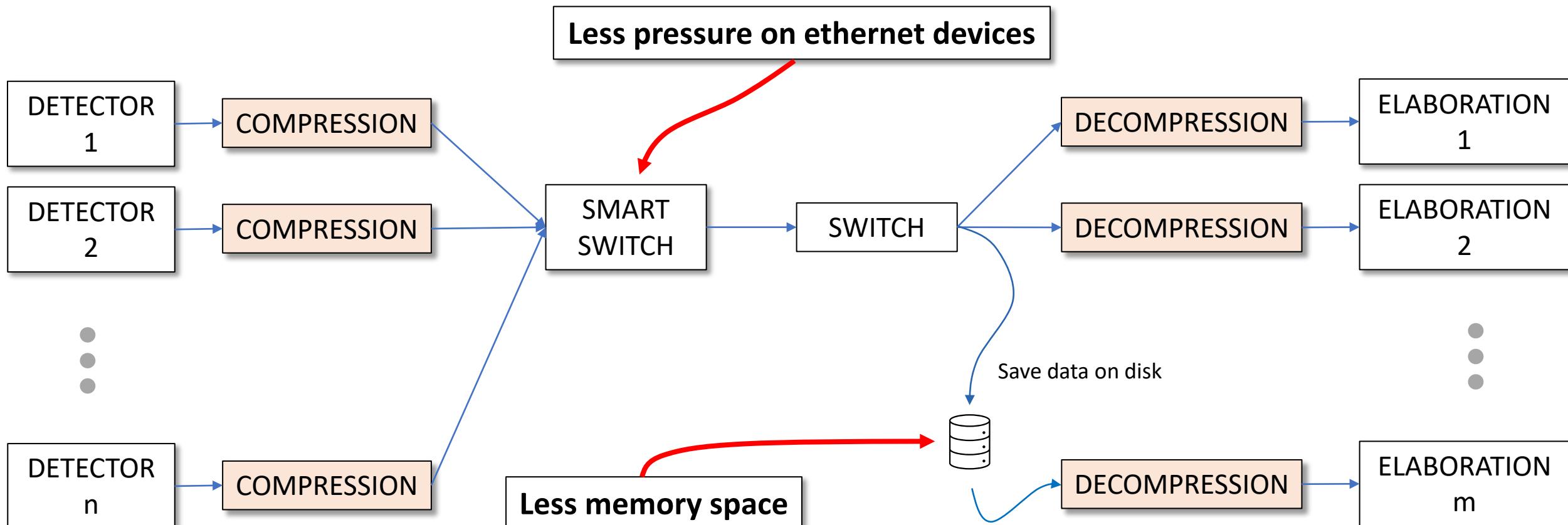
High data rate.
New design.

Pros:

All channels can be part of the trigger.
High level sophisticated trigger techniques.
Software trigger.



Block scheme of data flow





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA

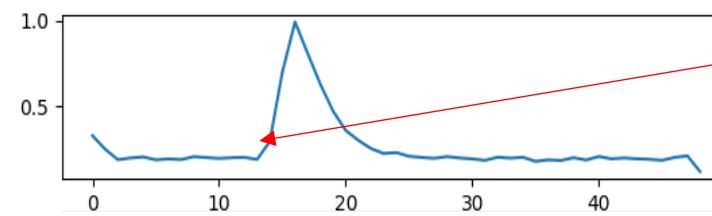


Future
Artificial
Intelligence
Research

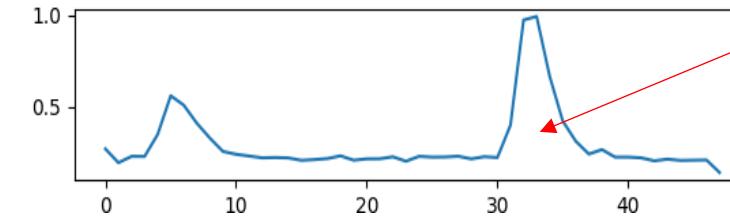
Data from physical Experiment

High
Event Probability
Low

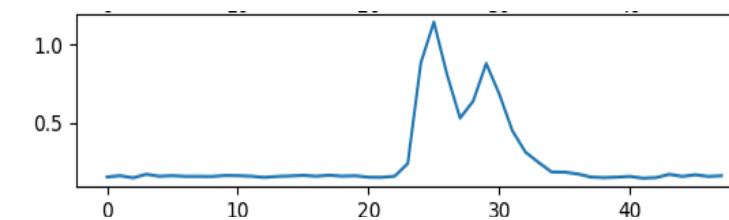
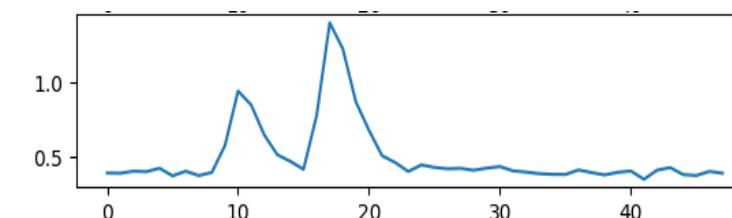
Very-Low probability signals
could be sent uncompressed



Leading edge



Integral of impulse





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIANZA



Future
Artificial
Intelligence
Research

Autoencoder Definition

Machine Learning Algorithm

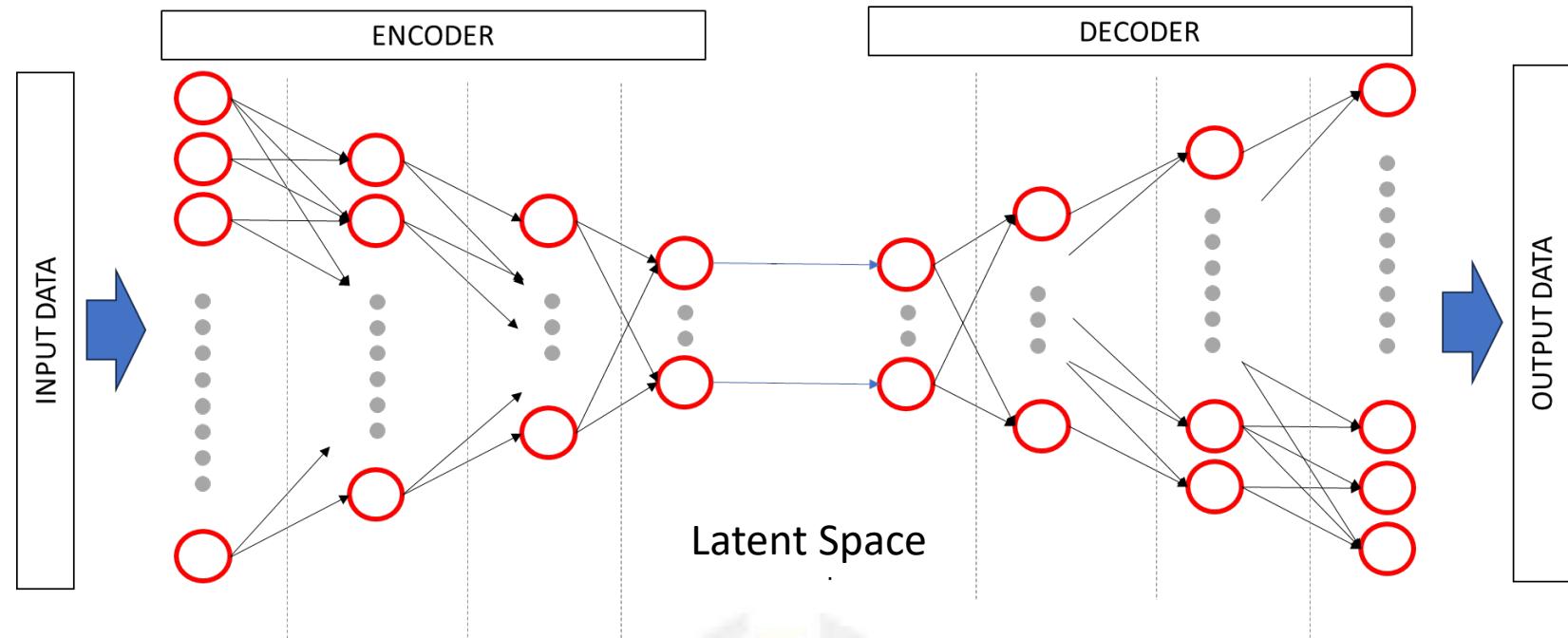
Artificial Neural Network

Unsupervised learning

Dimensionality reduction

Composed of two function:
- encoding
- decoding

FULLY CONNECTED AUTOENCODER WITH DENSE LAYER





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca

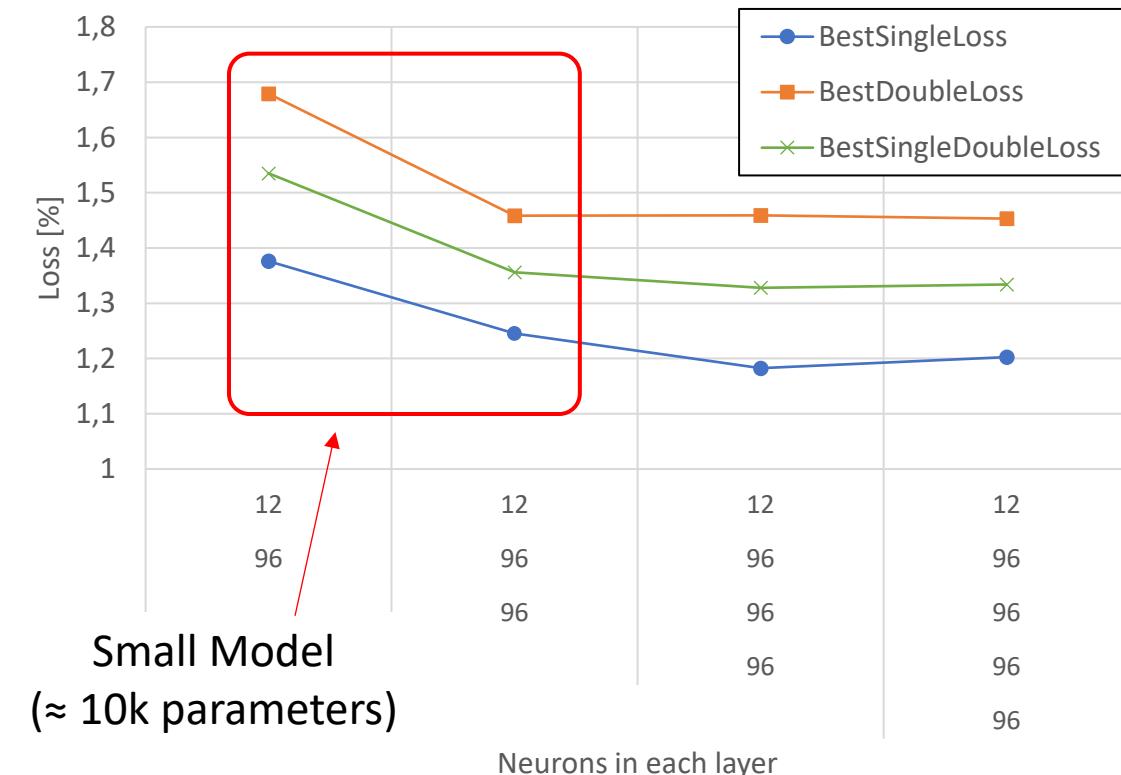
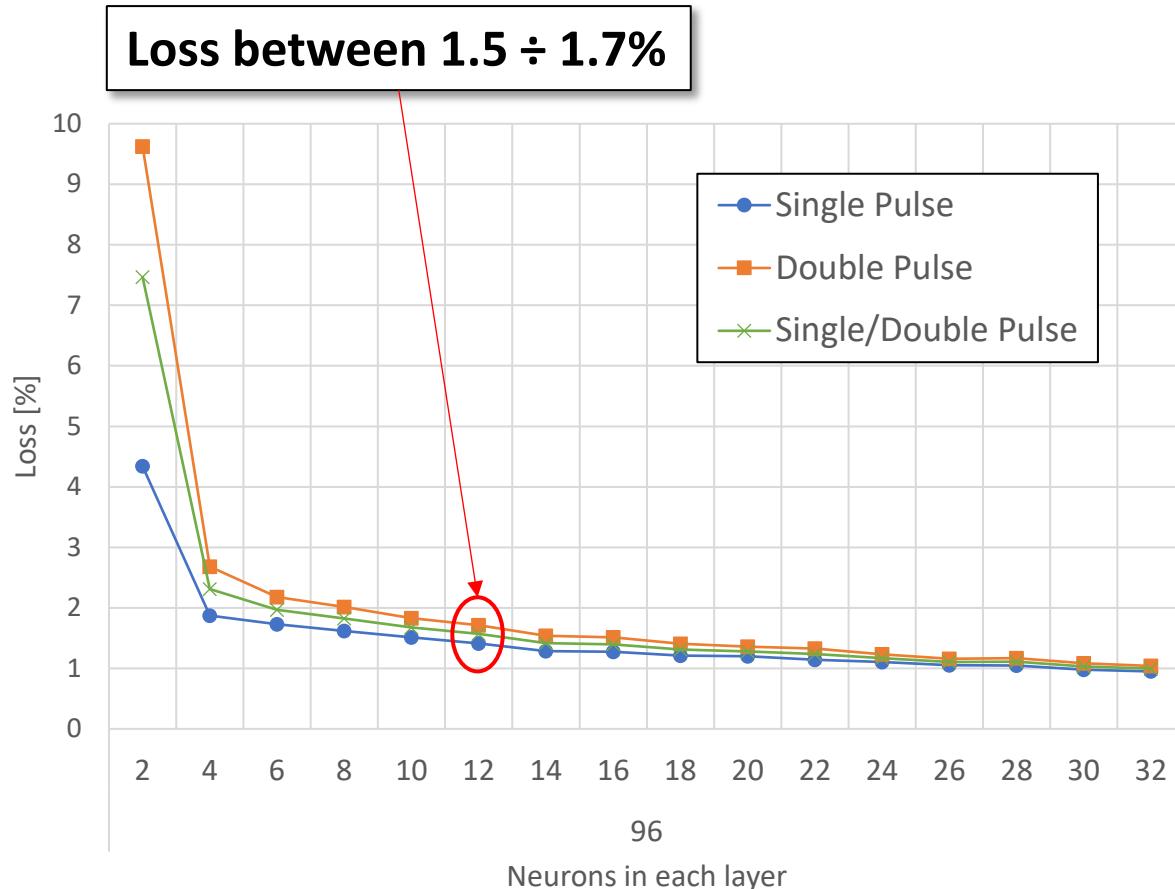


Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIANZA



Future
Artificial
Intelligence
Research

Training of the Autoencoder model



More detail in: Rossi, F., et al. "Artificial intelligence data reduction algorithm for streaming readout in high energy physics experiment." 2024 Congress in Computer Science, Computer Engineering, & Applied Computing (CSCE)



Finanziato
dall'Unione europea
NextGenerationEU



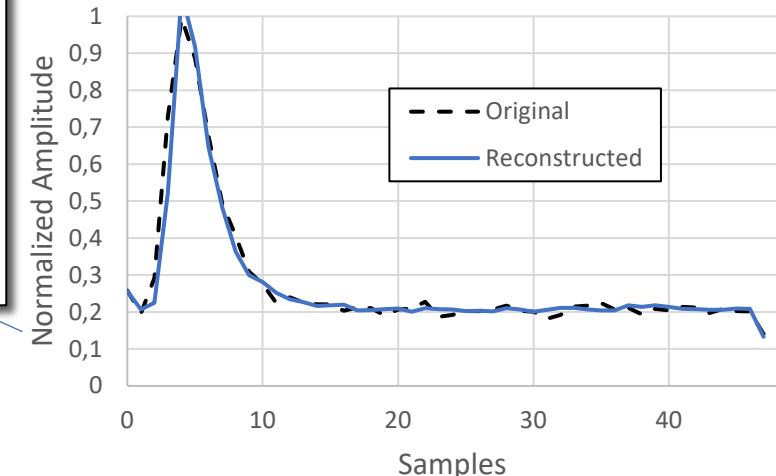
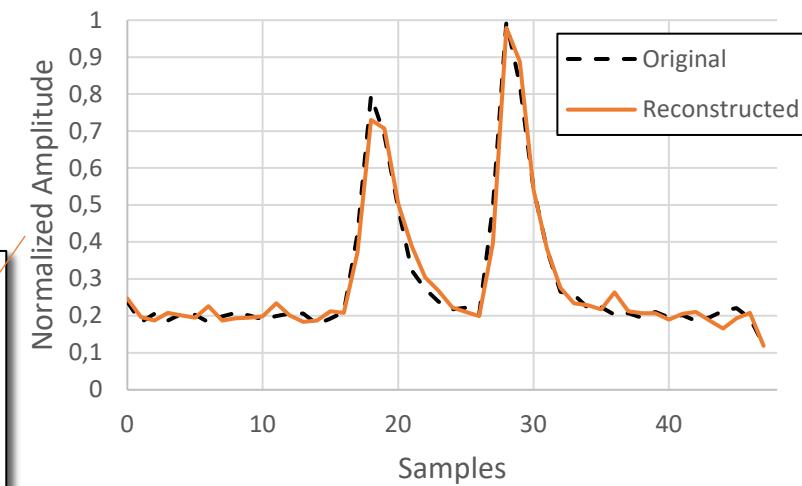
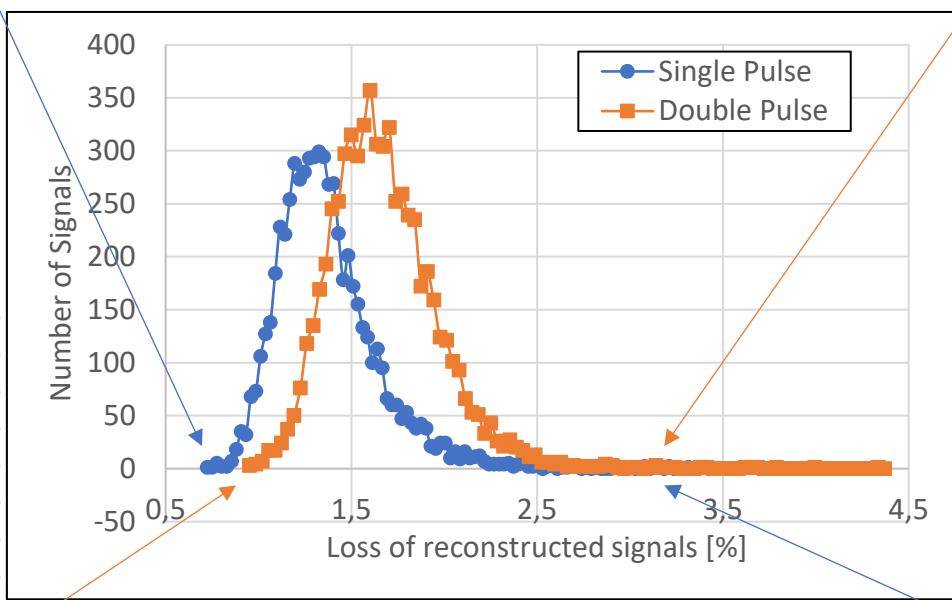
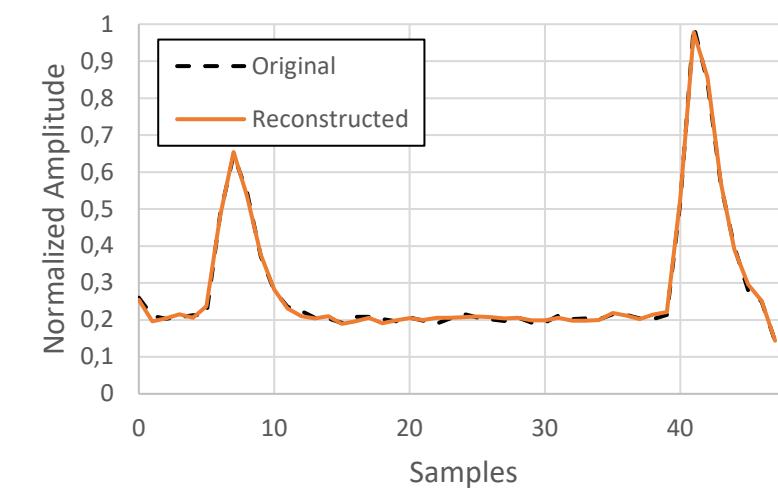
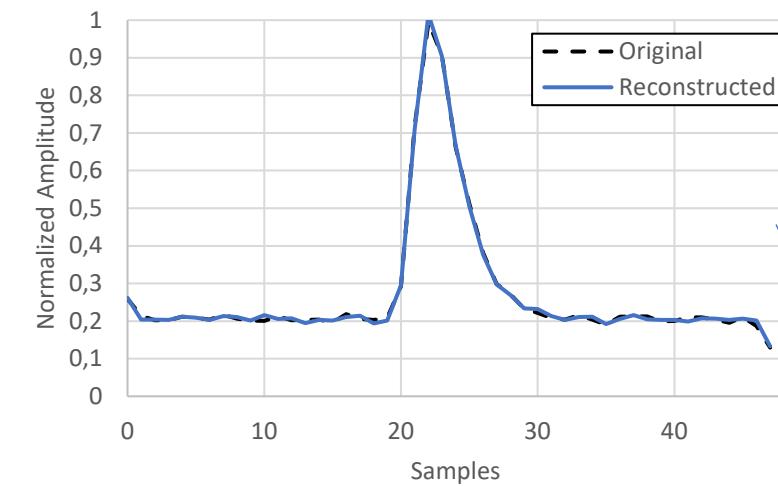
Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



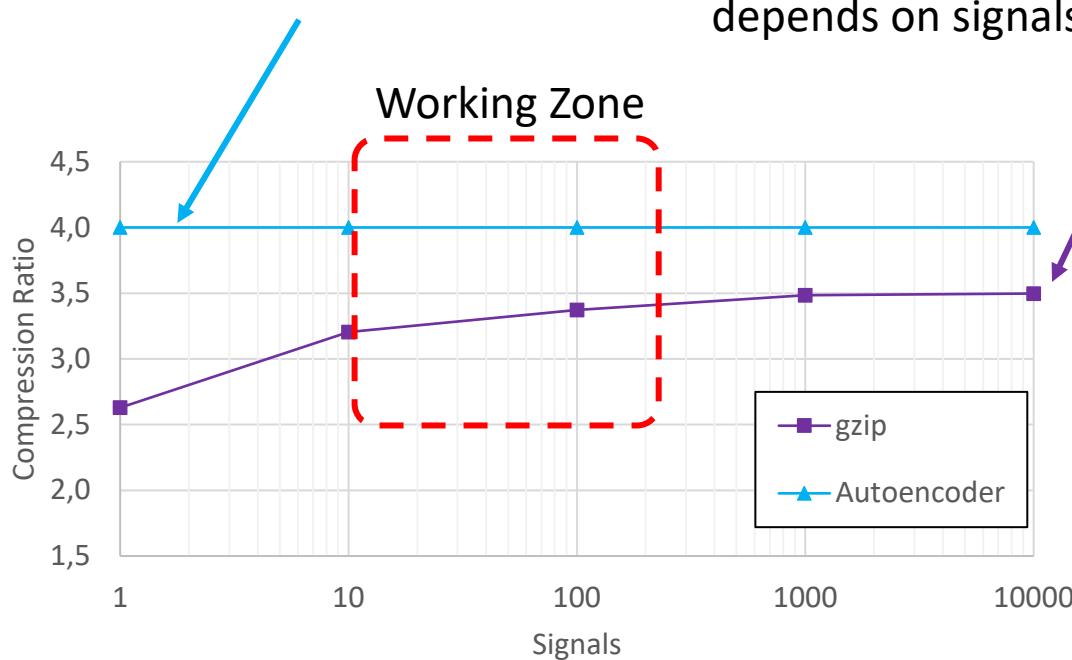
Future
Artificial
Intelligence
Research



Signals Compression

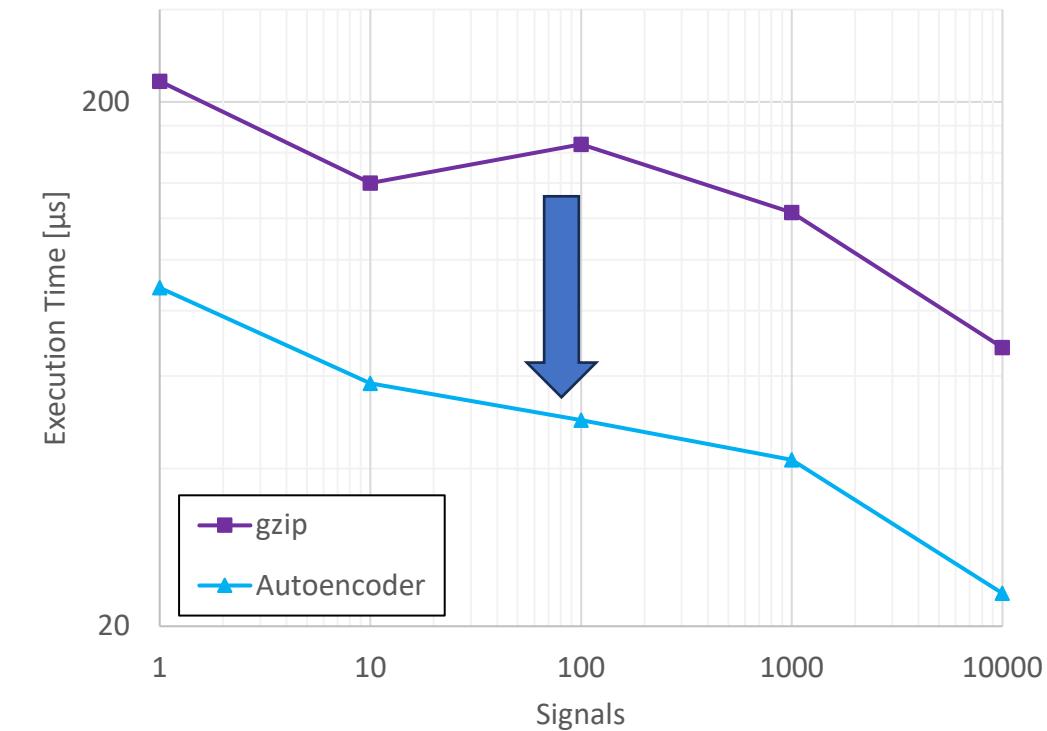
Comparison with standard lossless compression

Autoencoder compression
ratio is a Parameter



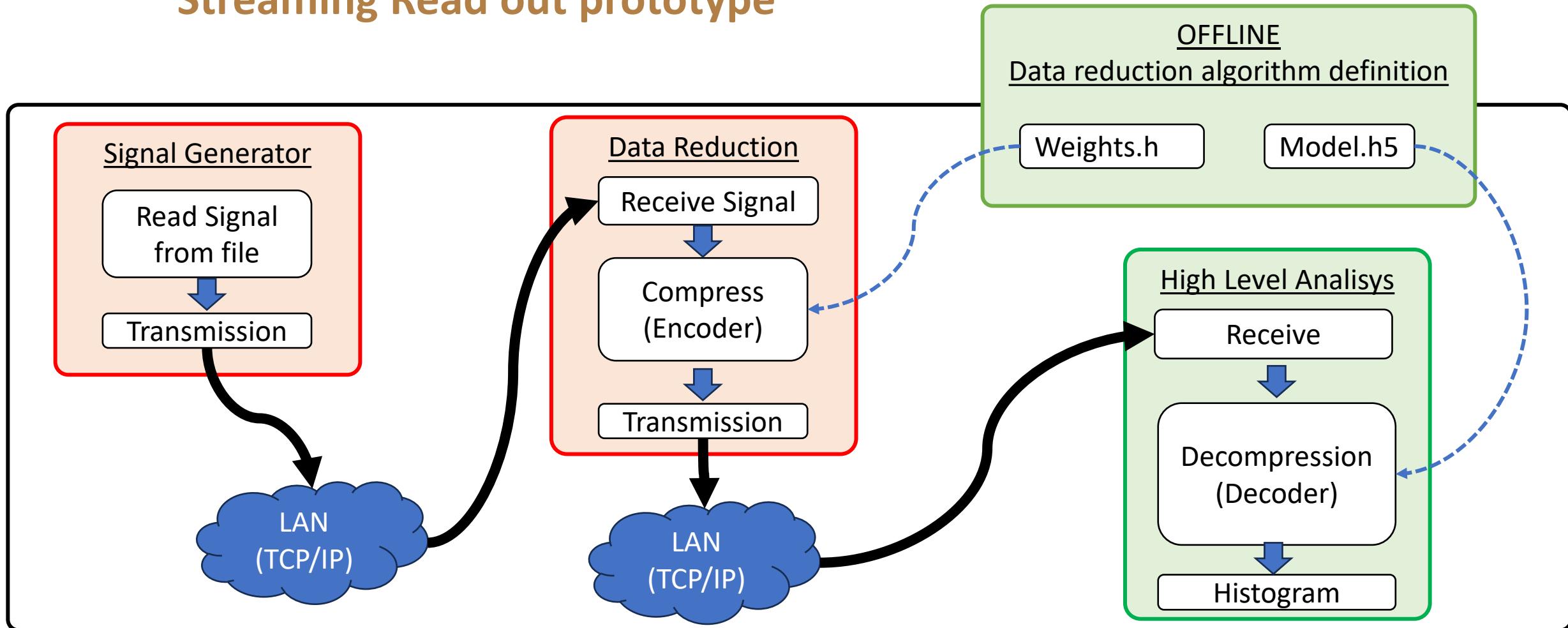
Better compression ratio

Gzip compression ratio
depends on signals number



Better also on execution time

Streaming Read out prototype





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



Future
Artificial
Intelligence
Research

Implementation of Data Reduction Node

4 x NVIDIA Tesla V100 GPU



INFN
CNAF

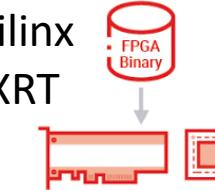
Data Reduction

Receive Signal

Compress
(Encoder)

Transmission

Xilinx
XRT



XILINX
VITIS™
AI

ALVEO V70 FPGA



Raspberry Pi 4 Rev. B



Low cost hardware

LAN
(TCP/IP)

High performance DELL C6400 server
(4 x AMD EPYC 7413 24-Core Processor)





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA

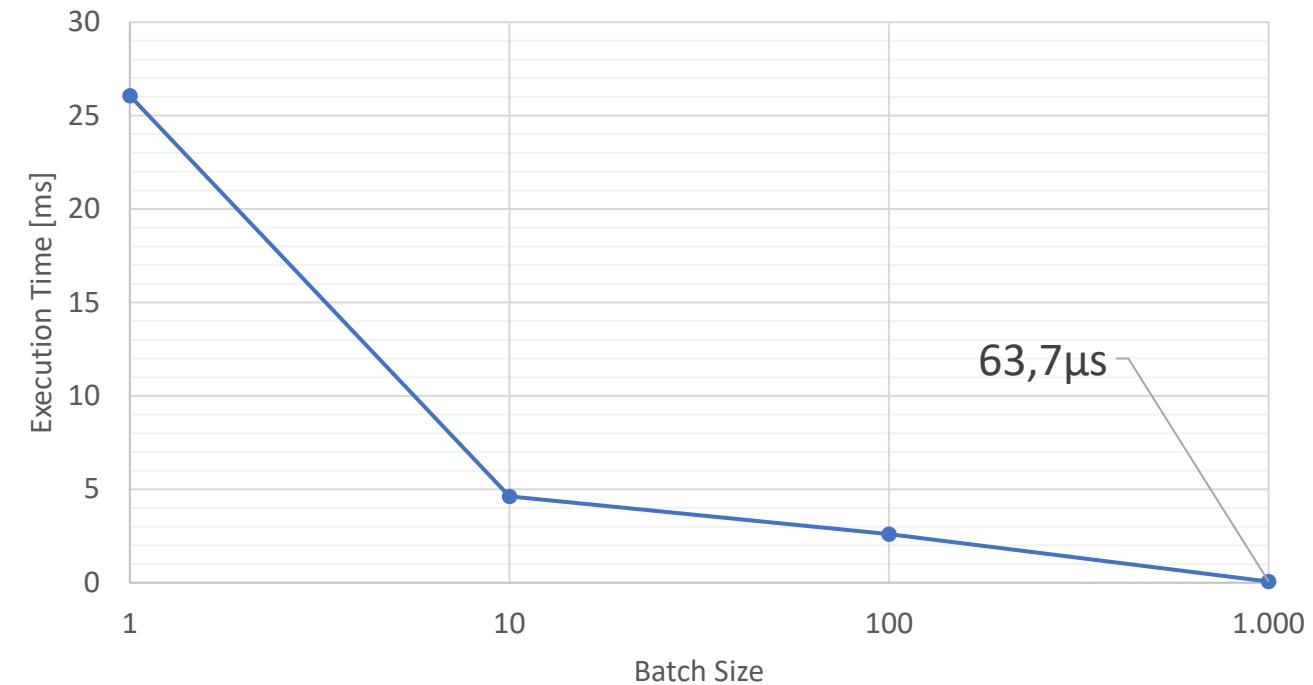


Future
Artificial
Intelligence
Research

Implementation: GPU



Execution time not enough
for the application!





Finanziato
dall'Unione europea
NextGenerationEU



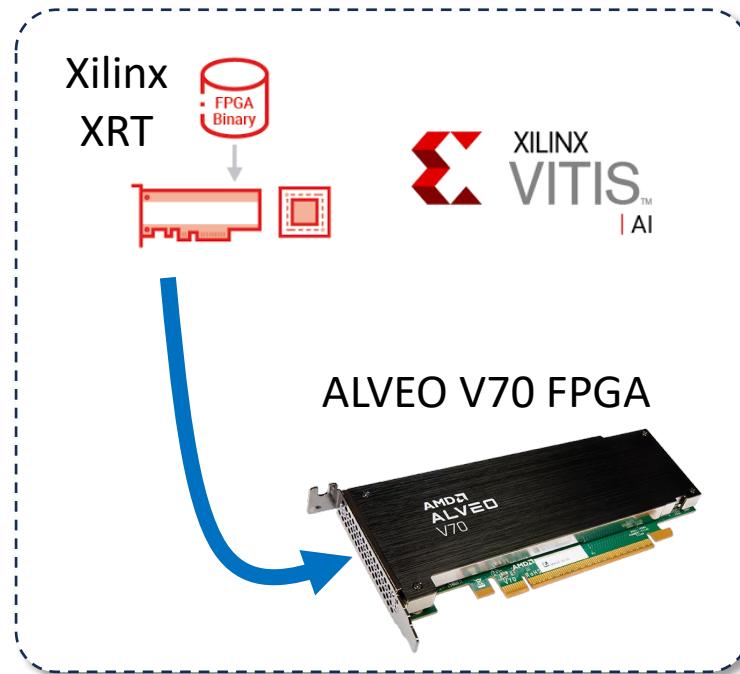
Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIANZA



Future
Artificial
Intelligence
Research



Execution time still not
enough for the application!

Implementation: FPGA

DEFINITION

- Definition
- Training
- Test
- Validation

PREPARATION

- Pruning
- Quantization



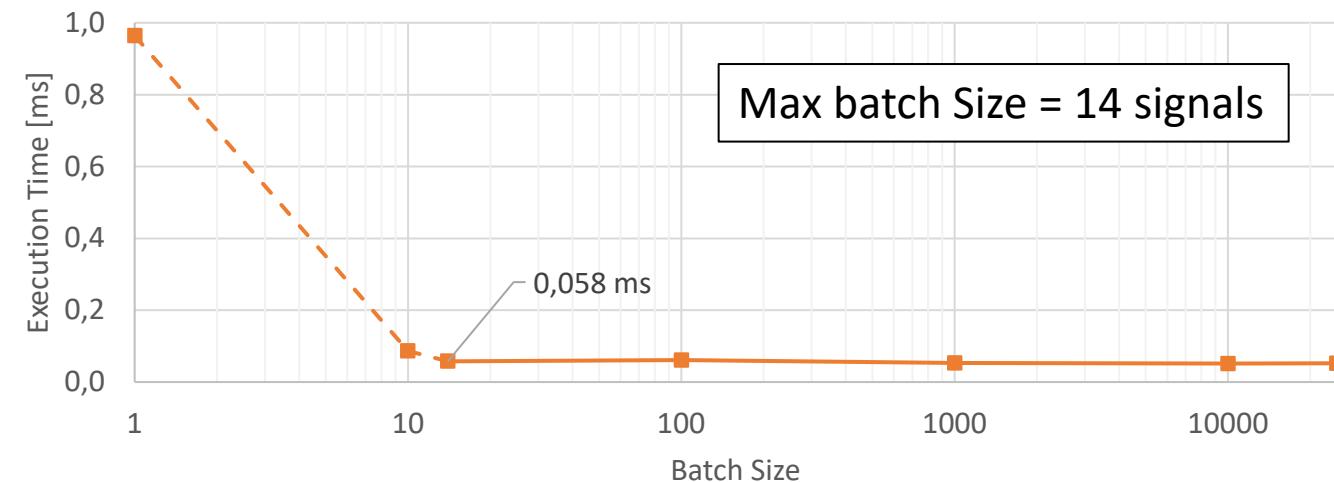
INITIALIZATION

- Compilation
- Deployment

ONLINE

- Inference

Compression time of single signal





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca

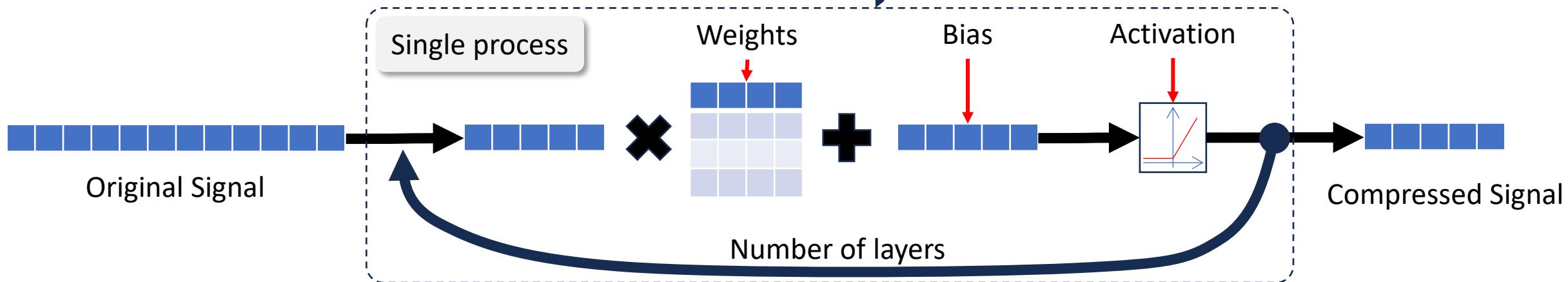
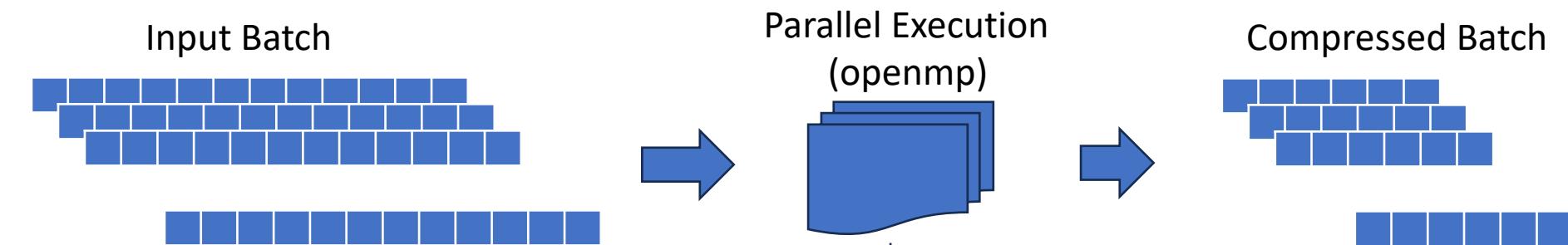
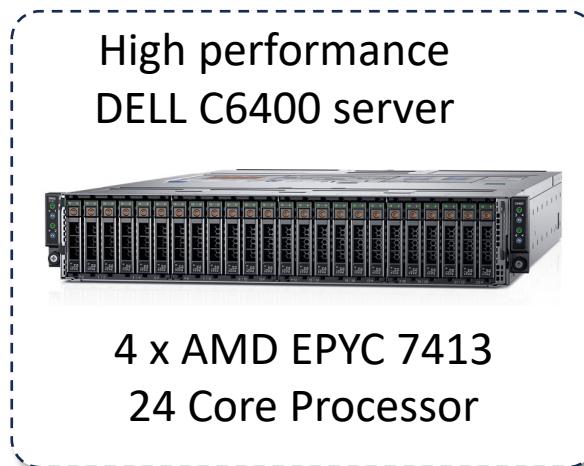


Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA



Future
Artificial
Intelligence
Research

Implementation: High performance server





Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



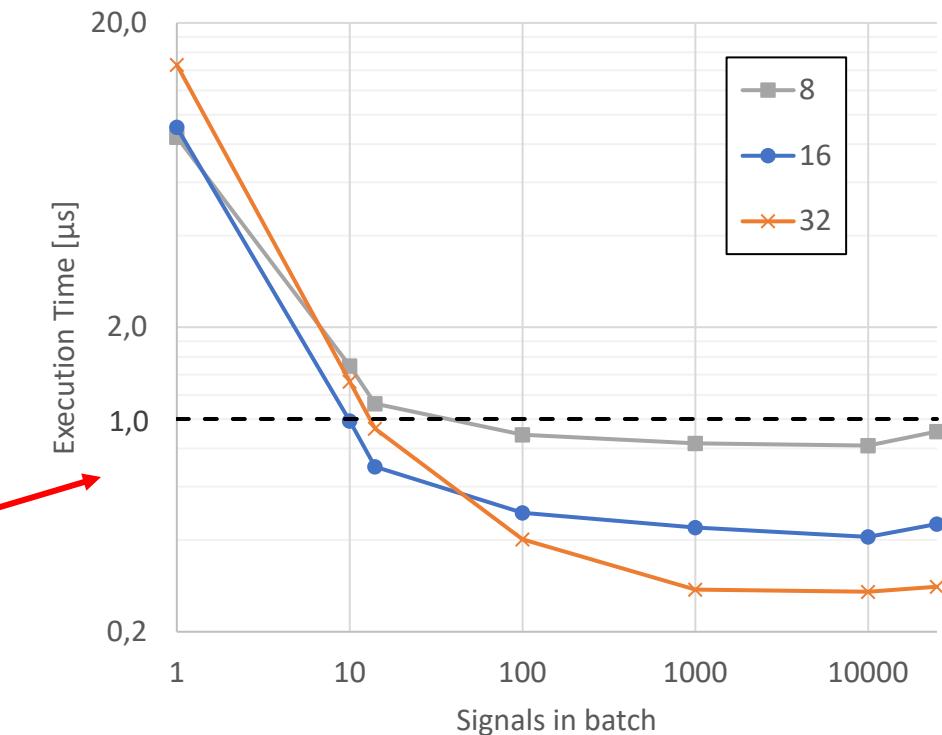
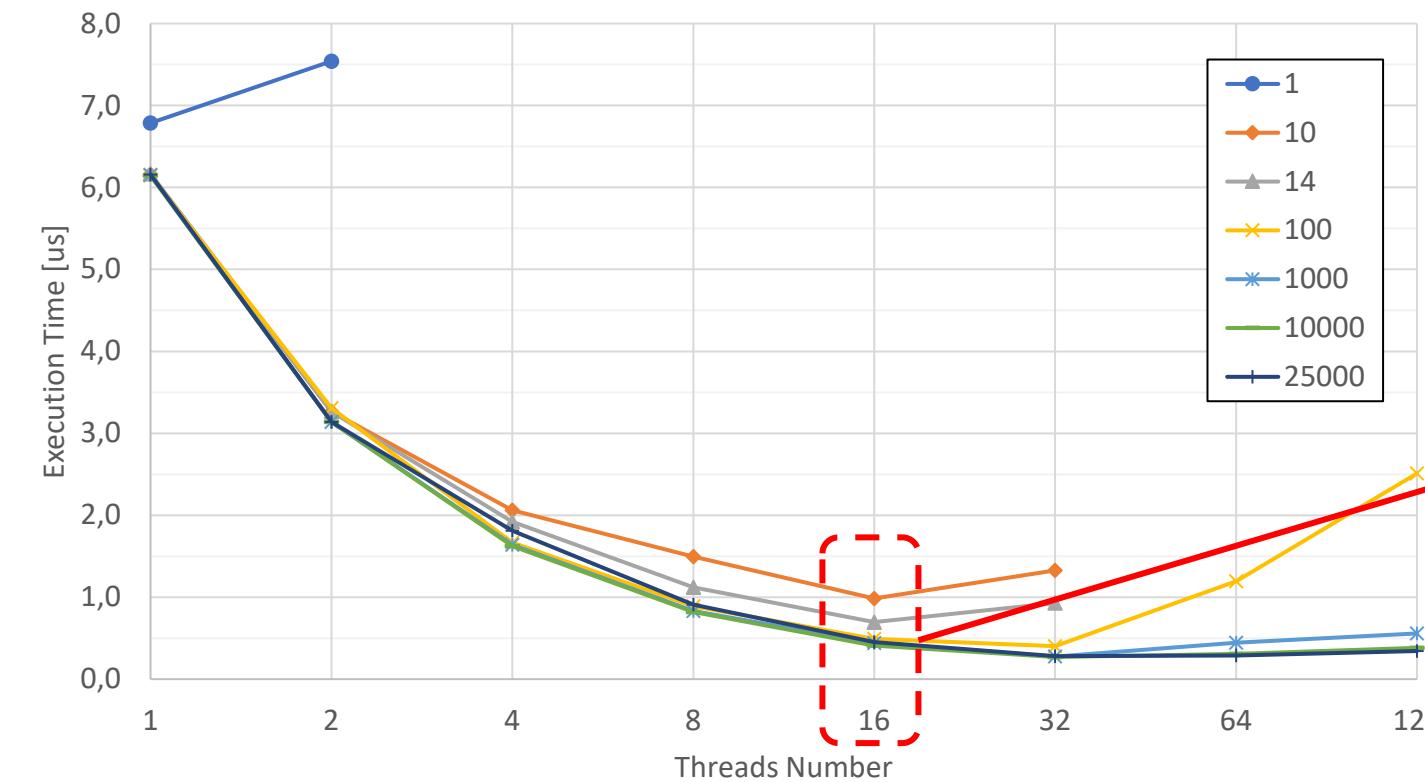
Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA



Future
Artificial
Intelligence
Research

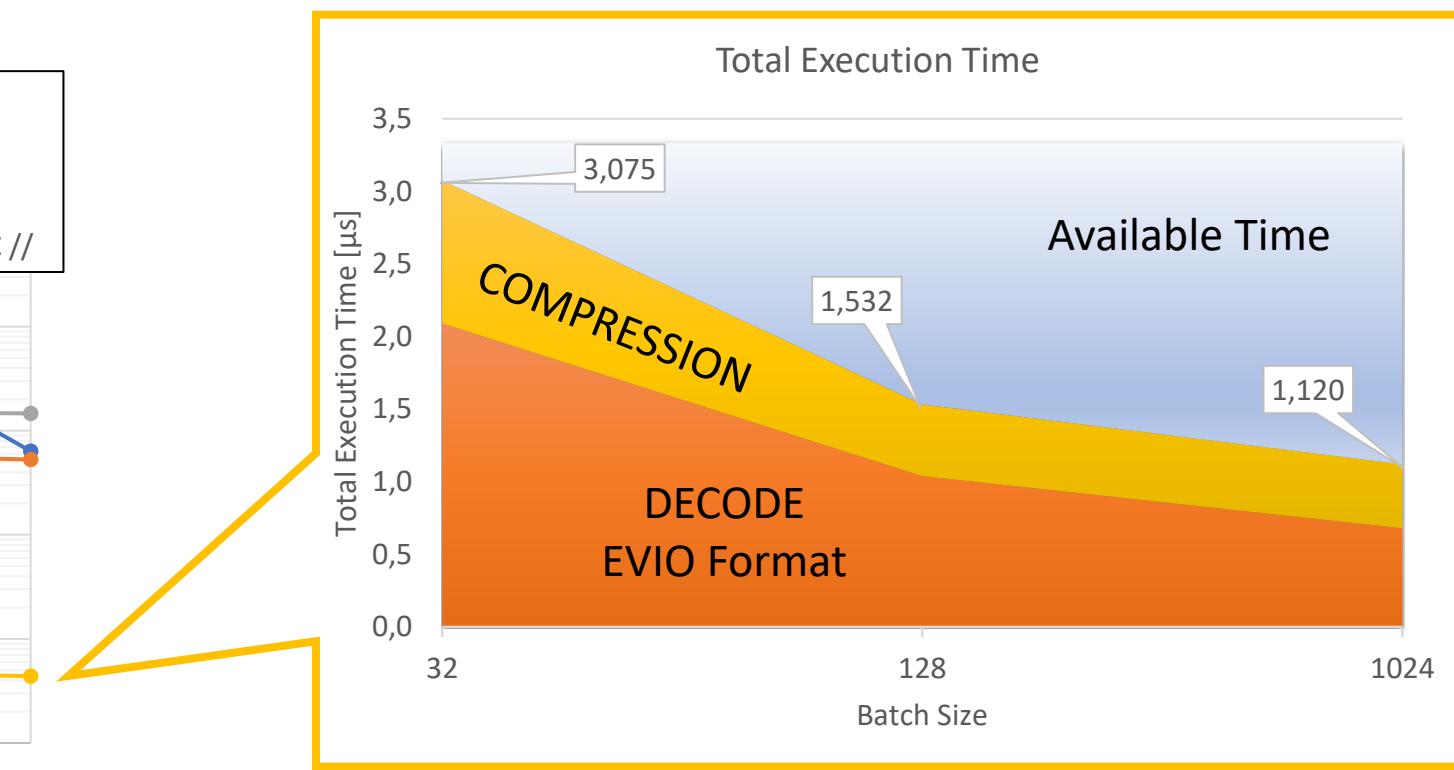
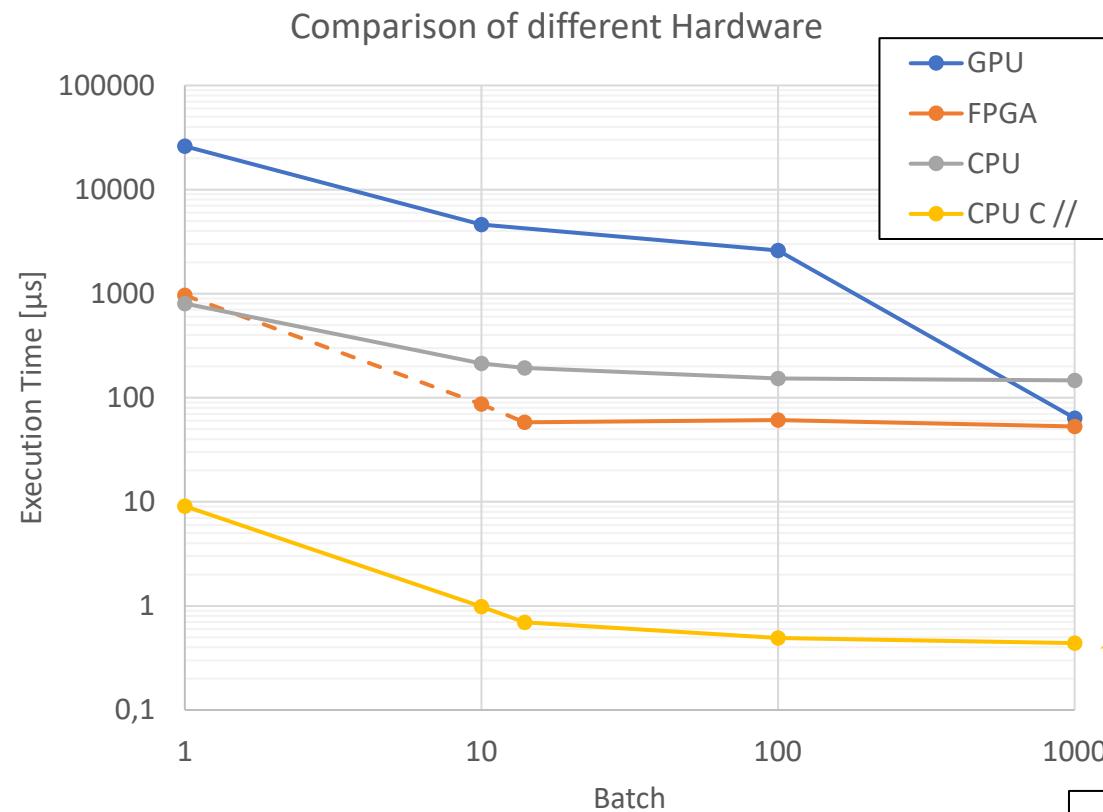
Implementation: High performance server

Execution time of different batches and threads number



Chosen 16 Threads
Reasonable execution time for the application

Conclusion



Rate can be managed for EVIO packet with at least 32 signals



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA



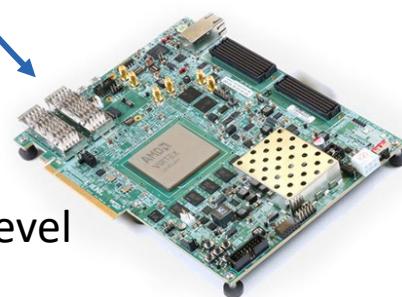
Future
Artificial
Intelligence
Research

Further Studies

Low level FPGA implementation



Dedicated connectivity
(2xQSFP28 @ 100GbE)

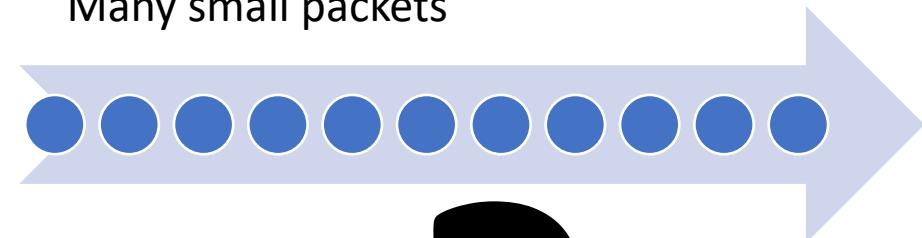


...or very Low level

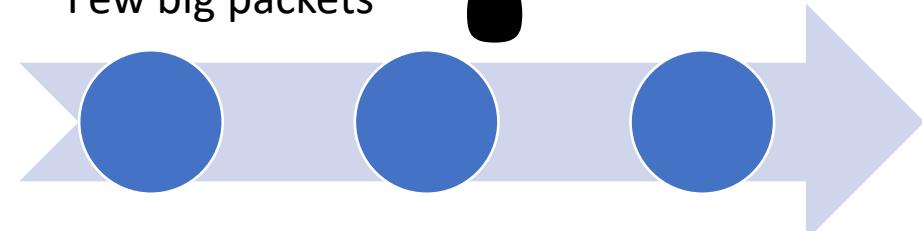
Reduce execution time and maybe save money

Statistical analysis of signals in each EVIO packet

Many small packets



Few big packets



?

TBD

Estimate performance on real acquisition



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILLENZA



Future
Artificial
Intelligence
Research

Thank you for your attention



FUTURE AI RESEARCH

<https://fondazione-fair.it>



<https://www.jlab.org>



<https://www.ge.infn.it>



<https://sealab.unige.it>

ACKNOWLEDGMENT

Authors have received support from: **FAIR - Future Artificial Intelligence Research, funded by the European Union Next-Generation EU (Italy)Research**) – spoke 6.