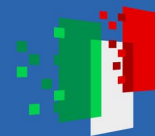




Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



Future
Artificial
Intelligence
Research

Streaming Readout Workshop

SRO-XII

December 2-4, 2024
University of Tokyo, Japan

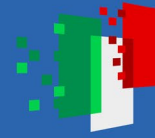


Real-Time data reduction with Artificial Intelligence for SRO.

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

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

Gagik Gavalian
Jefferson Lab
Newport News (Virginia)



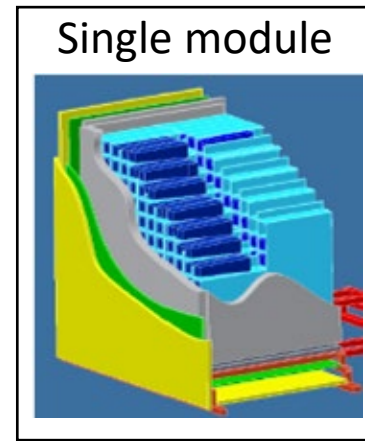
High Energy Physics Experiment: Beam Dump eXperiment (BDX)

Jefferson Lab

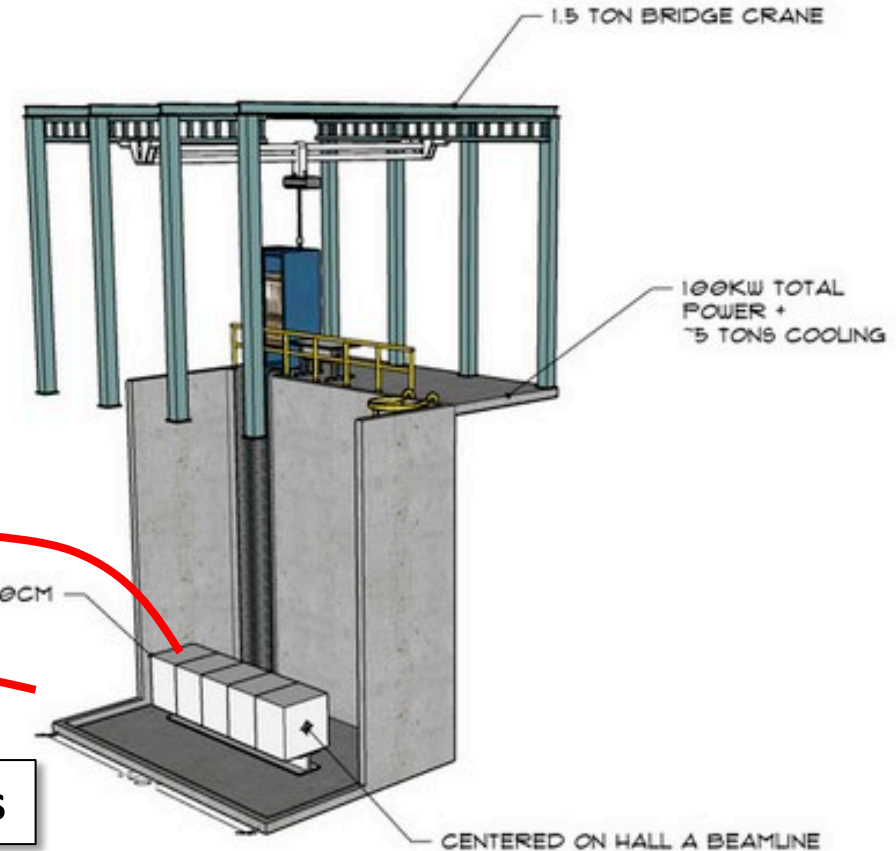


≈ 1000
Calorimeter channels
(30MB/s)

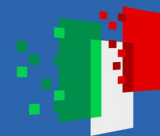
≈ 300
Veto channels
(500MB/s)



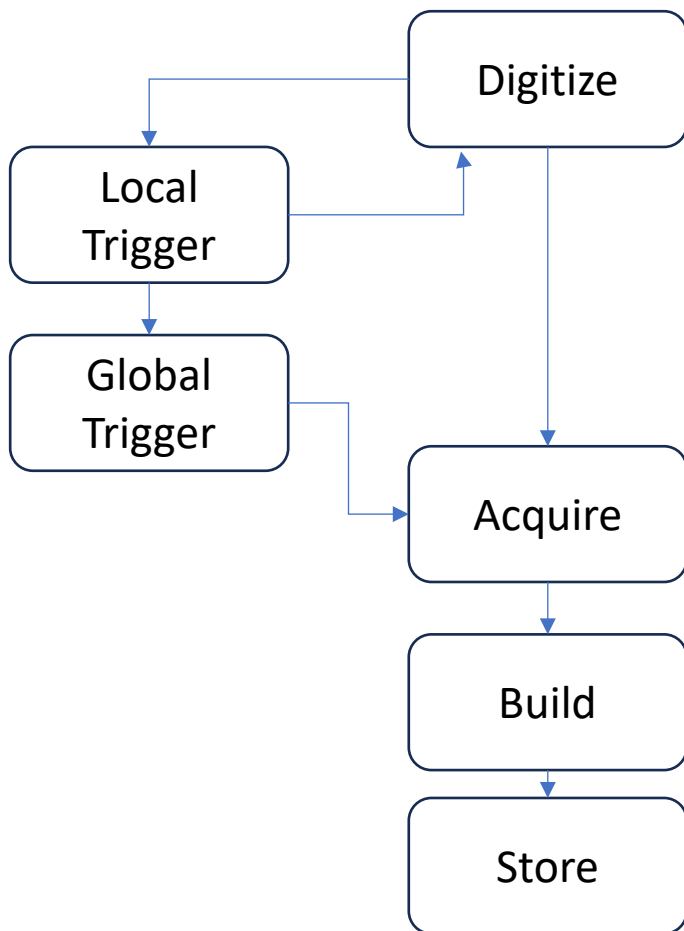
5 MODULES AT 80CM X 100CM X 100CM



Very rare occurrence of Dark Matter events



Traditional triggered DAQ VS Streaming Readout



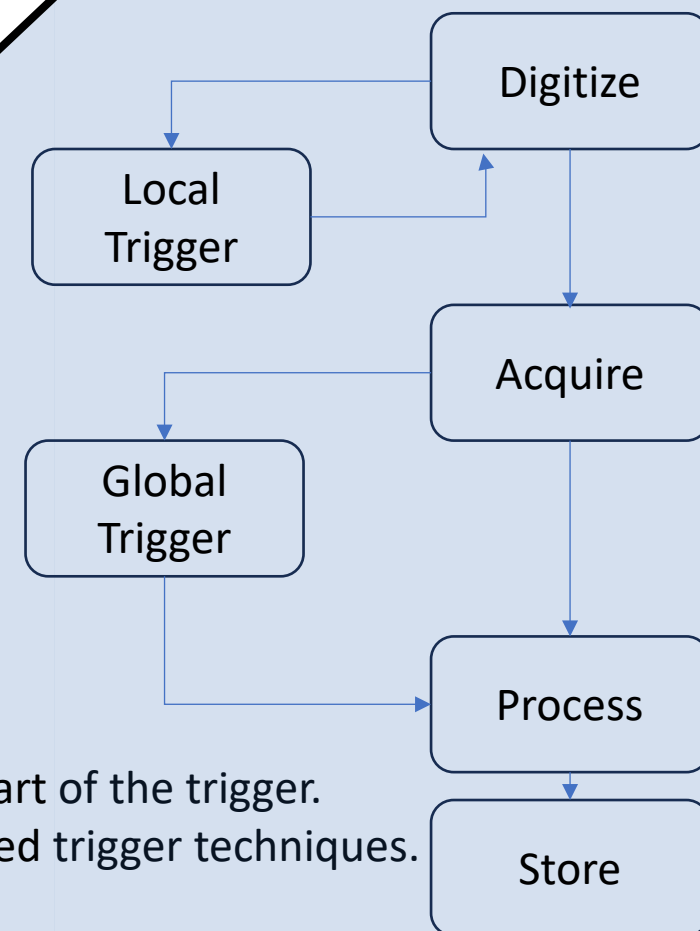
Cons:

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



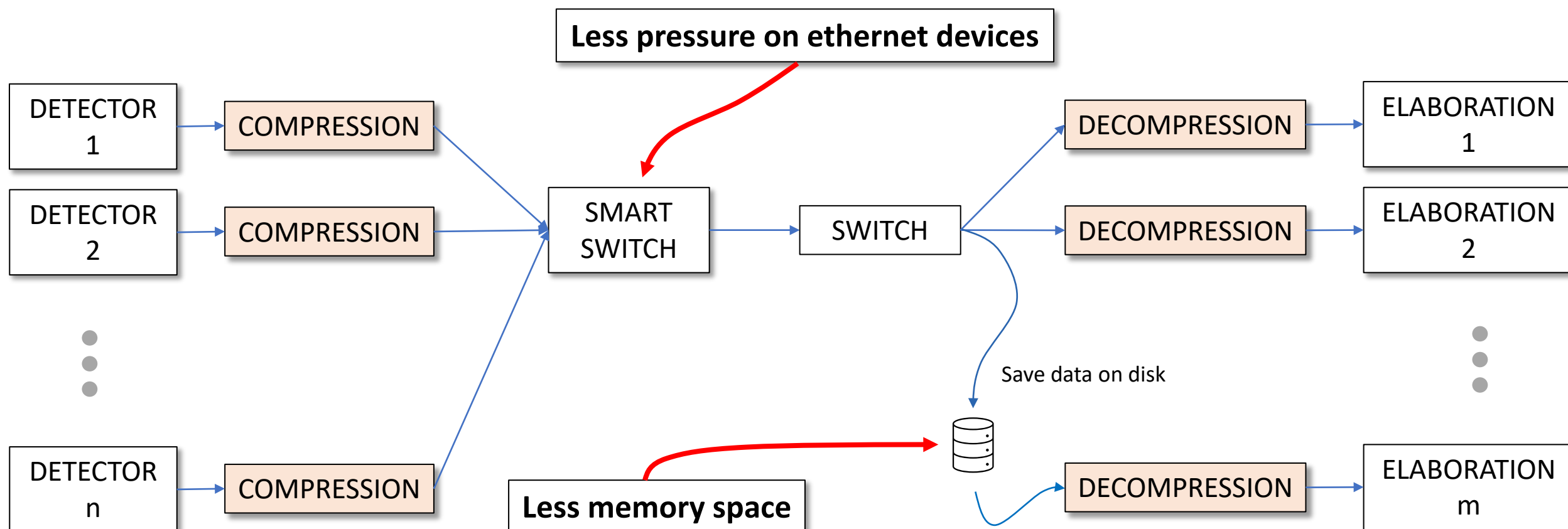
Cons:

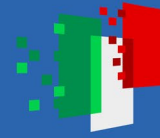
High data rate.
New design.

Pros:

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

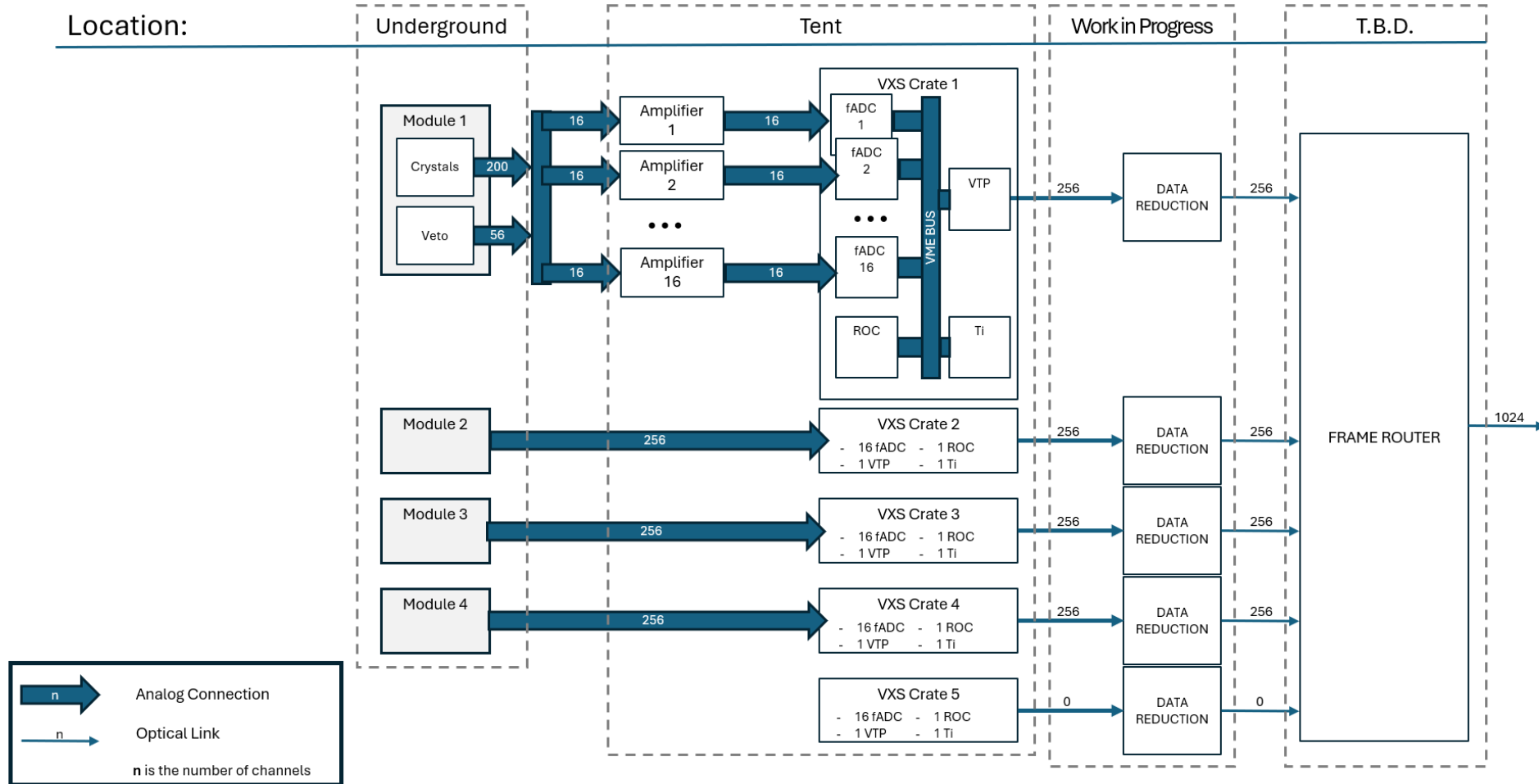
Block scheme of data flow





Detailed BDX data flow scheme

Location:

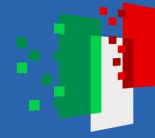




Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



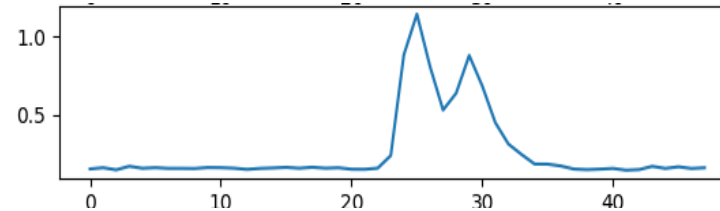
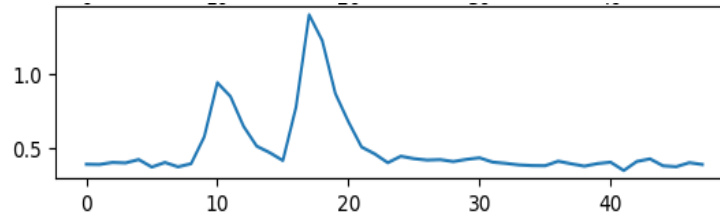
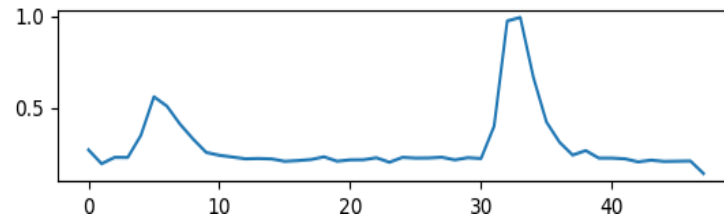
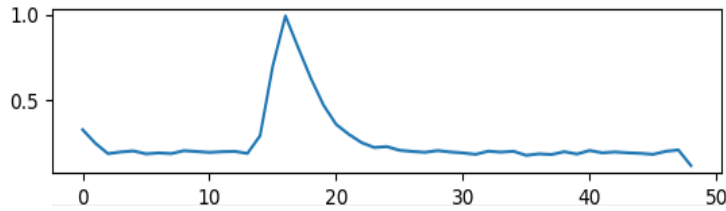
Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



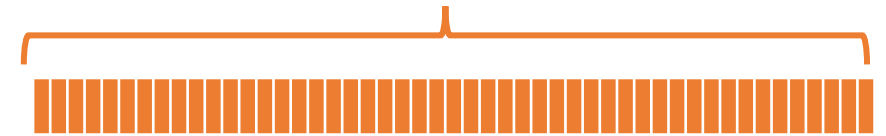
Future
Artificial
Intelligence
Research

Data from physical Experiment

High



48 samples



Event Probability

Very-Low probability signals
could be sent uncompressed

Low

Data reduction algorithm: Autoencoder

Machine Learning Algorithm

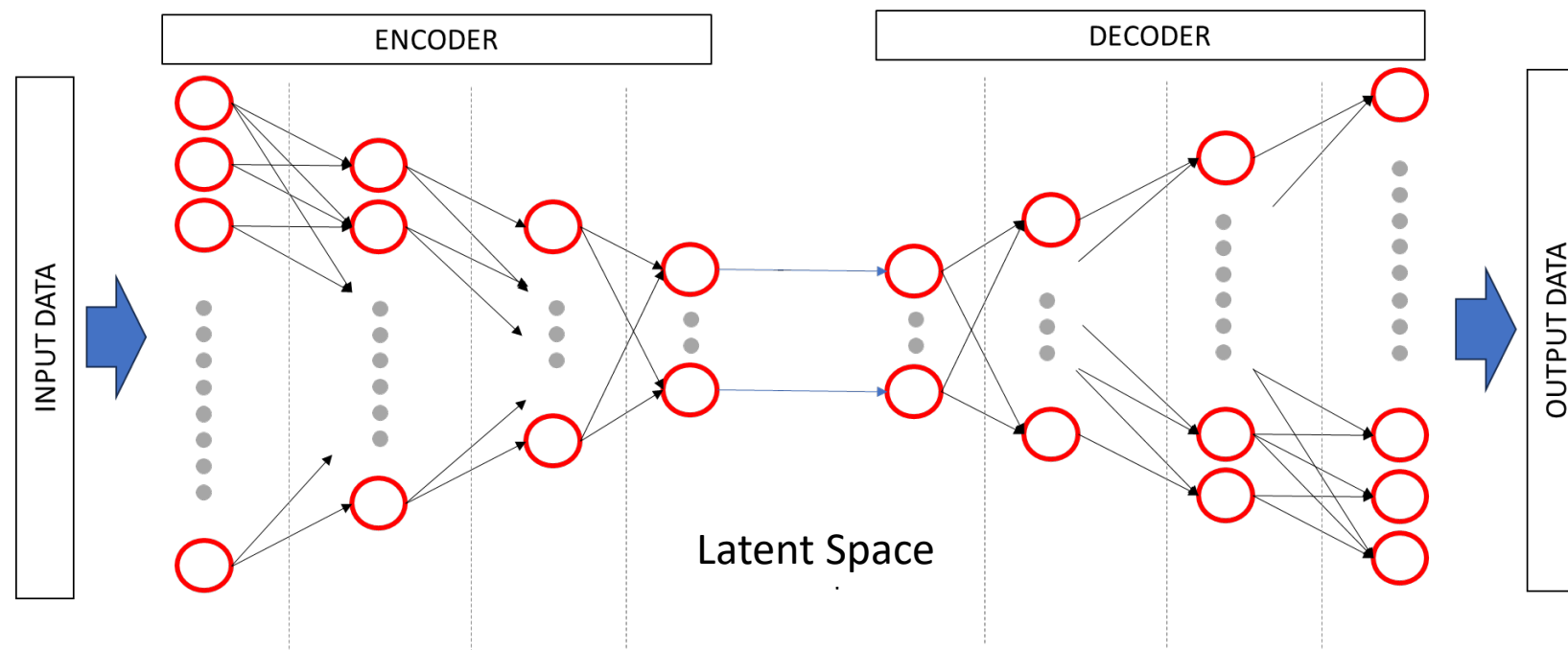
Dimensionality reduction

Unsupervised learning

Artificial Neural Network

Composed of two function:

- encoding
- decoding



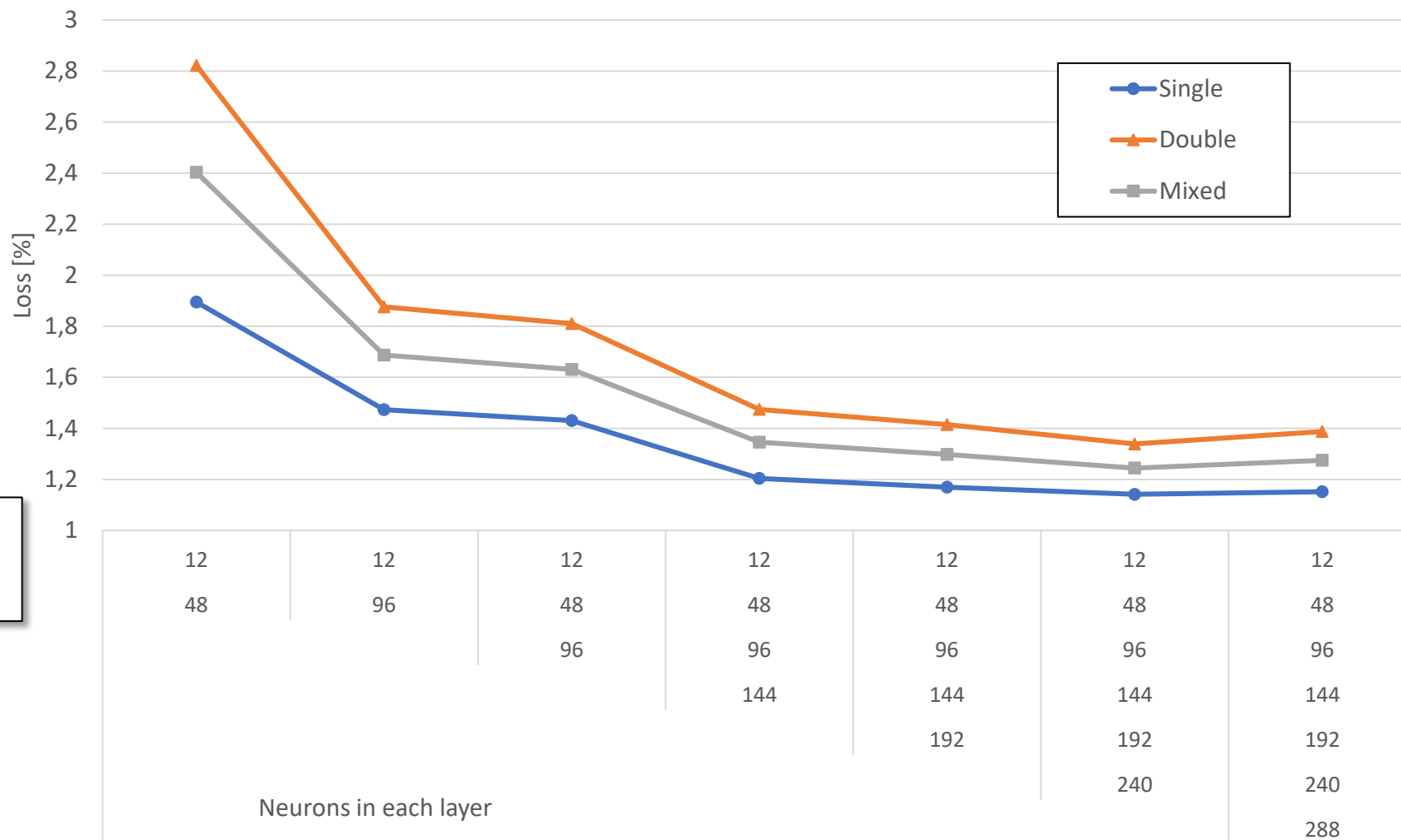
FULLY CONNECTED AUTOENCODER WITH DENSE LAYER

Lossy compression algorithm

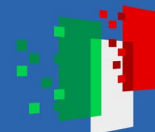
Autoencoder Training: Different configuration

Test increasing:

- Layer number
- Neurons in each layer



Few improvement adding more and more parameters



Autoencoder Training: Different configuration

Chosen configuration

Layers: 3

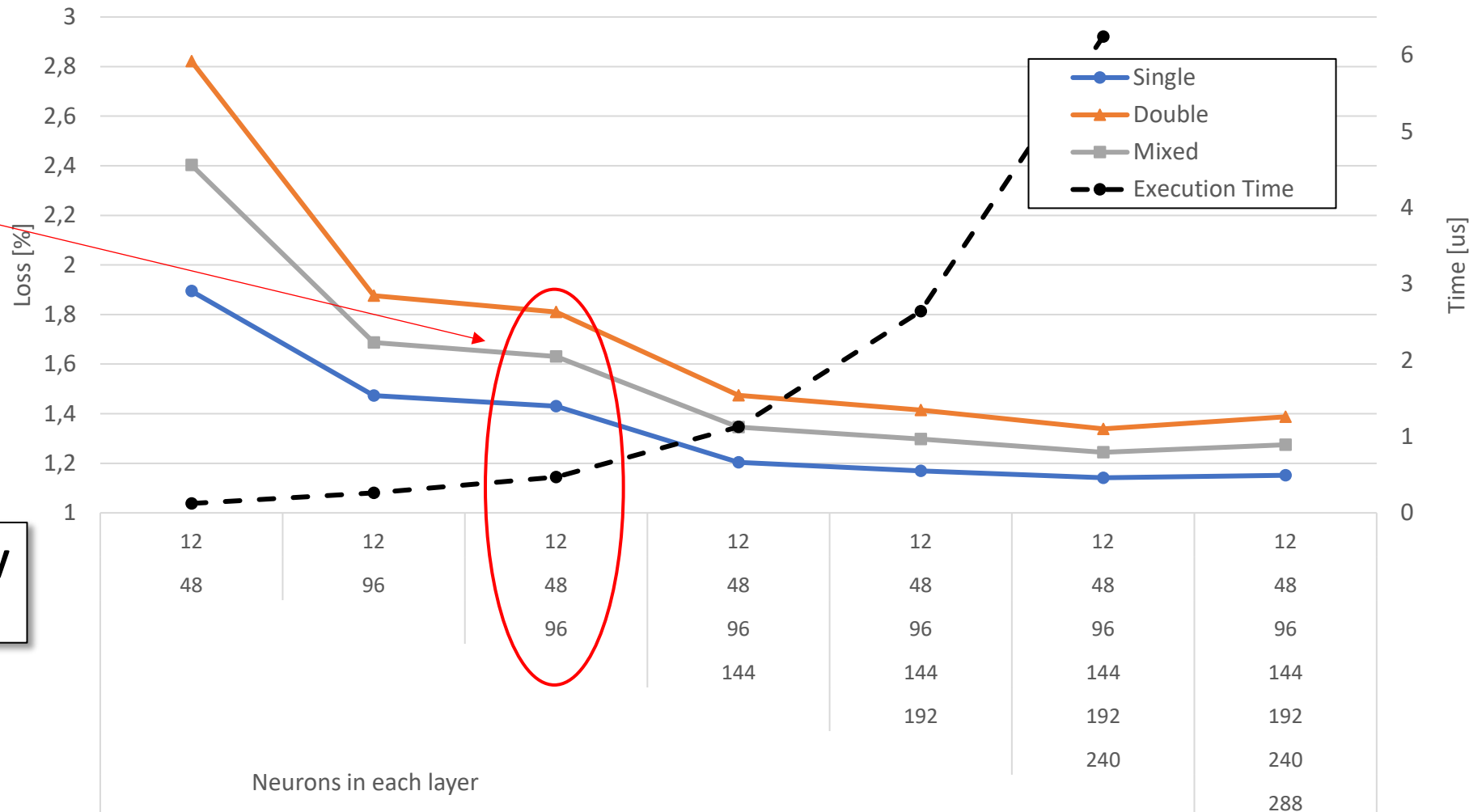
Neurons:

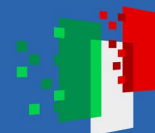
96

48

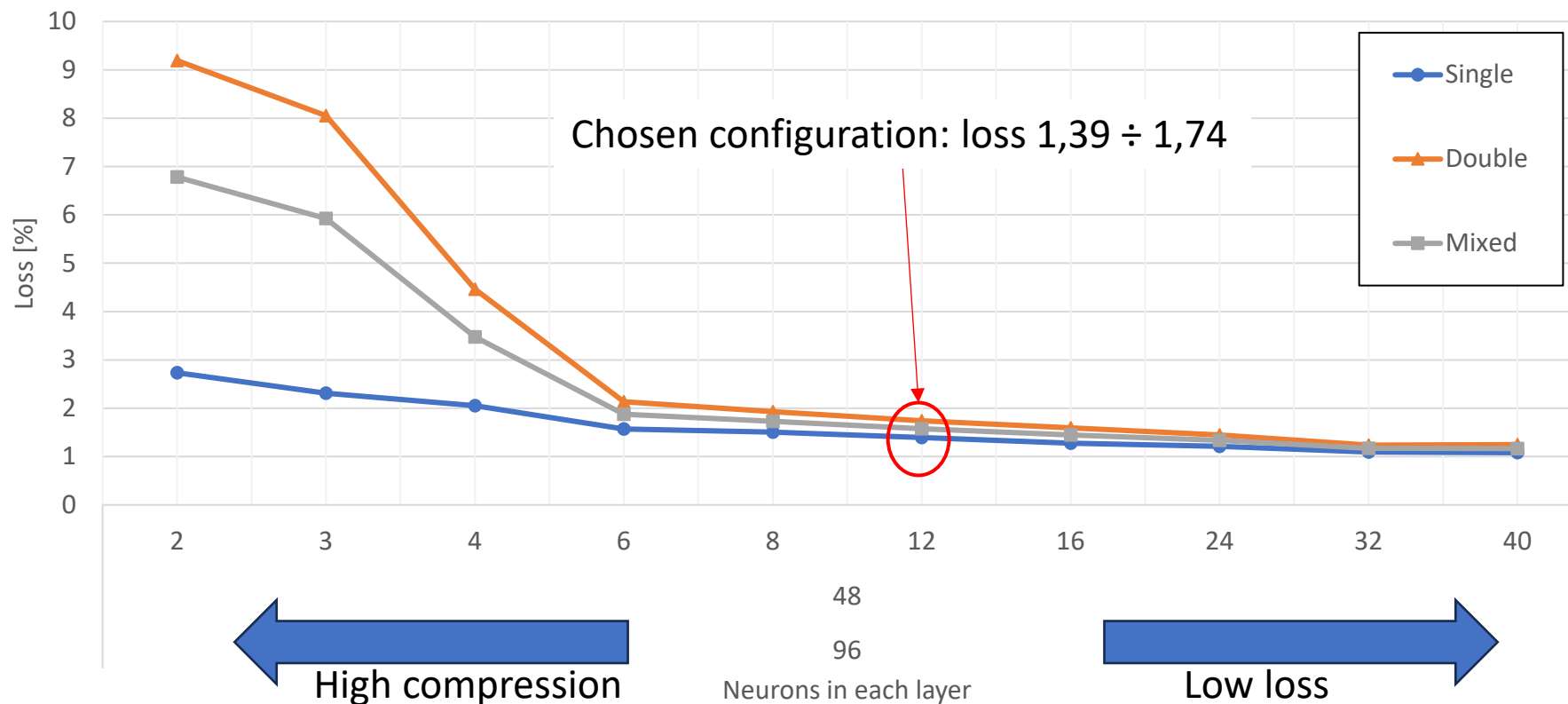
12

Execution time increase very fast with model complexity





Autoencoder Training: Different latent space



Compression ratio is a parameter and could be chosen as loss tradeoff



Finanziato dall'Unione europea
NextGenerationEU



Ministero dell'Università e della Ricerca



Italiadomani
PIANO NAZIONALE DI RIPRESA E RESILIENZA



Future Artificial Intelligence Research

Autoencoder: Training time

GPU

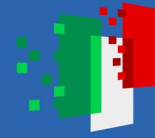


CPU

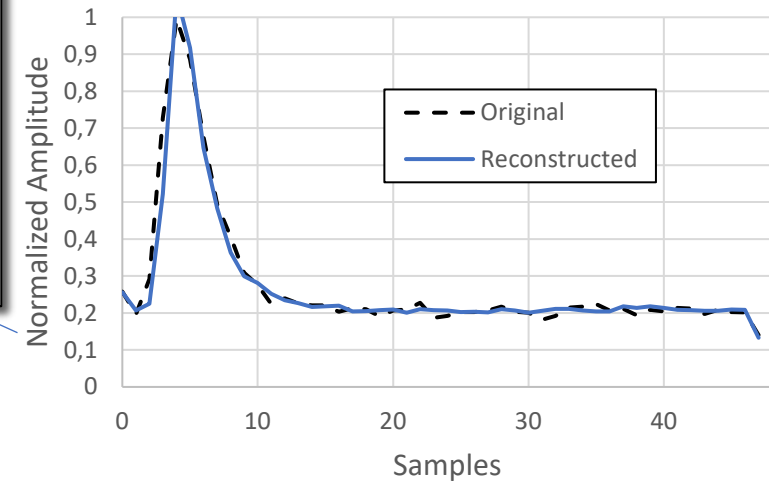
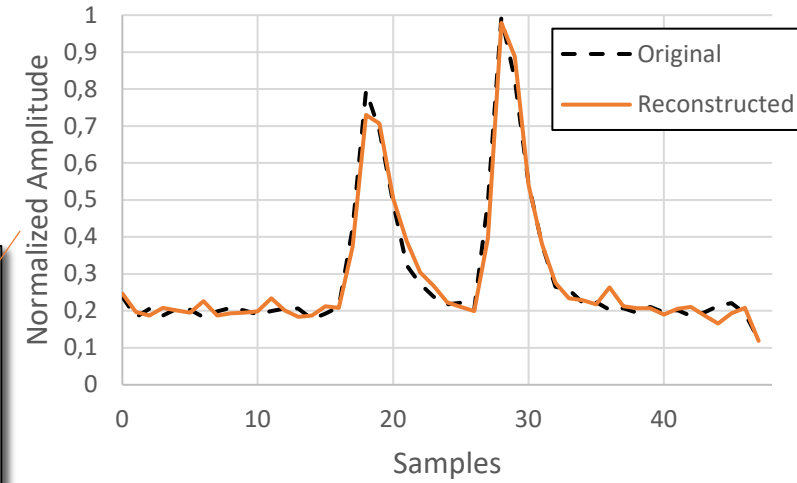
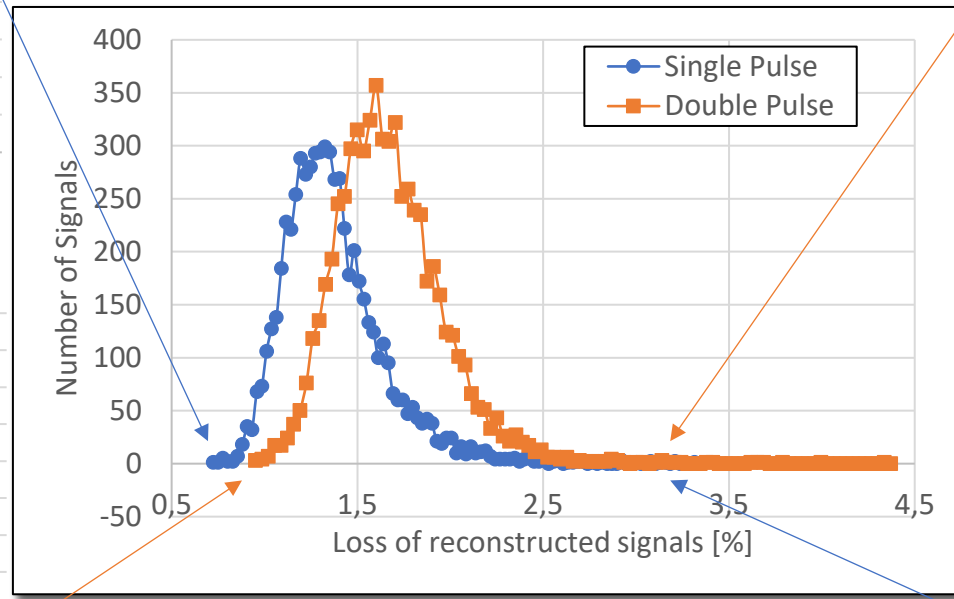
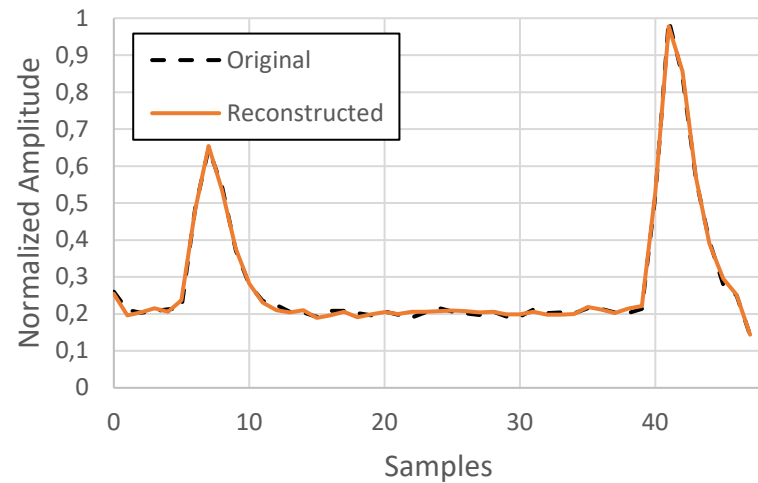
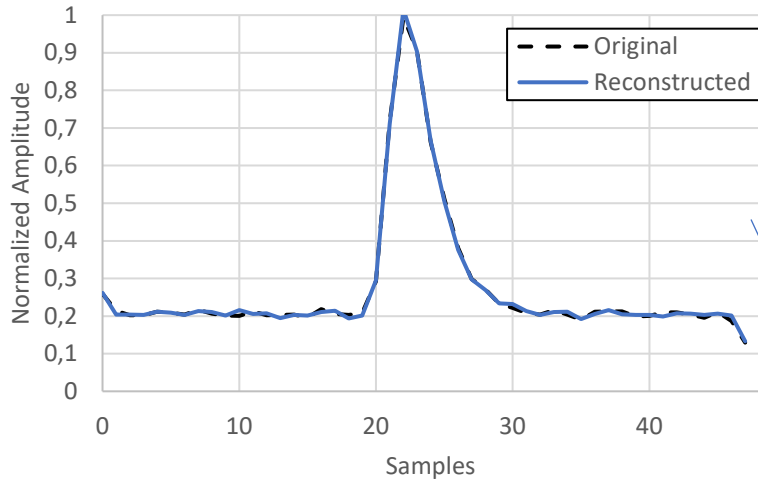
AMD Ryzen 5 5600U
6 core, 12 Logic processor

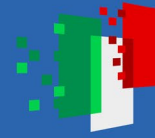


Suitable models for the application



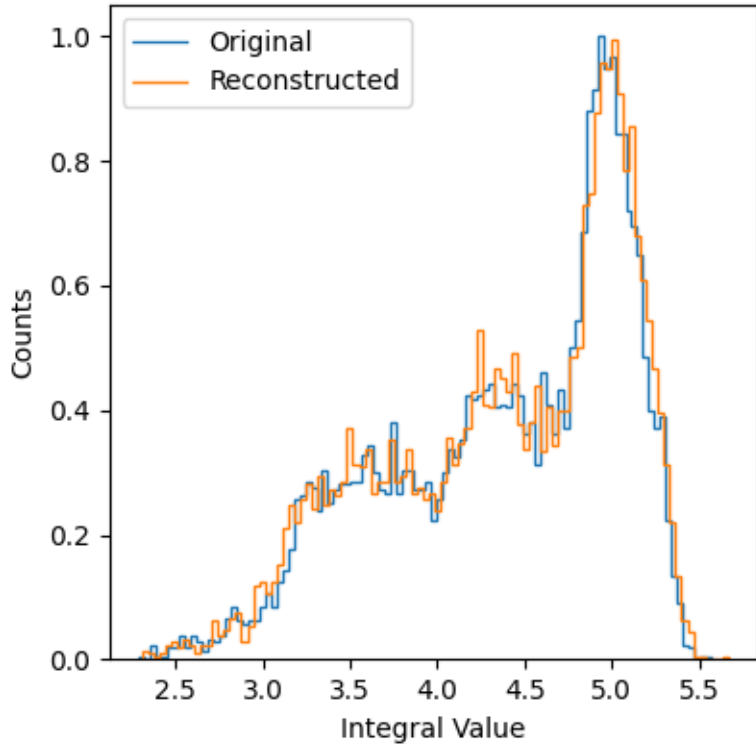
Signals Compression



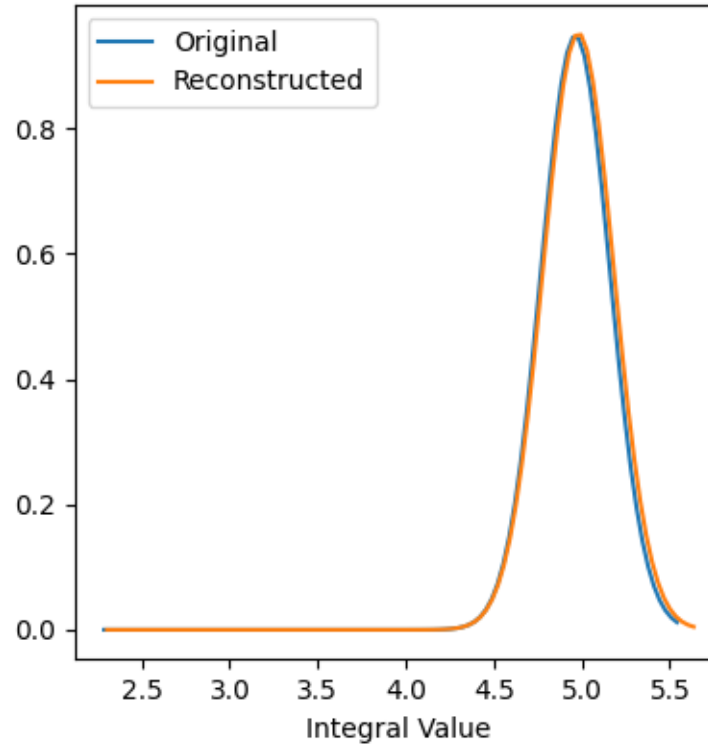


Signals Compression: Integral and spectrum

Integrals histogram



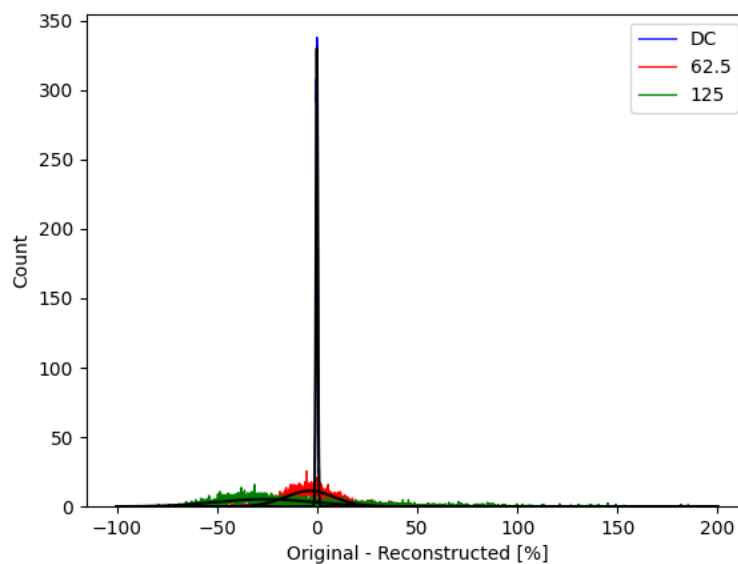
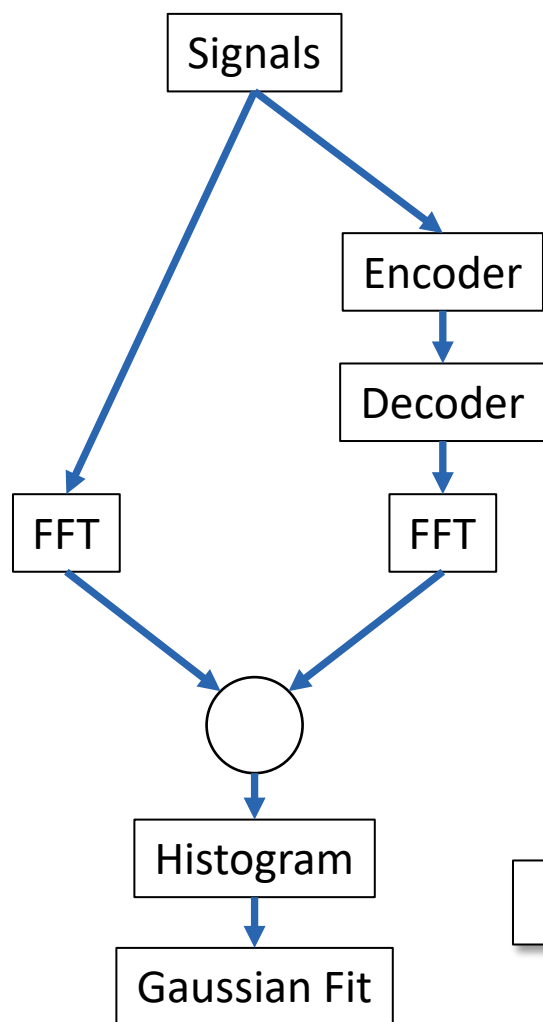
Best Fit



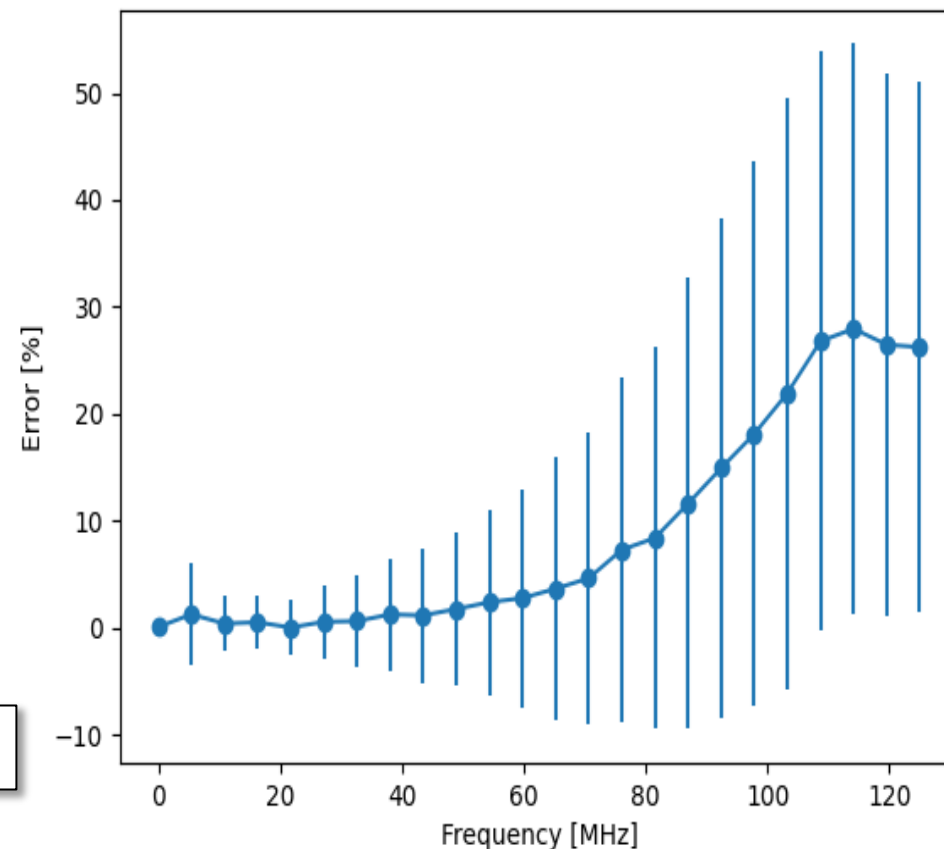
	A	μ	σ
Original	199,124	4,963	0,1952
Recon.	199,992	4,981	0,2006
Diff. [%]	0,44	0,35	2,77

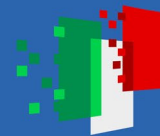
Good performance also on the derived quantities for physical analysis

Signals Compression: FFT analysis

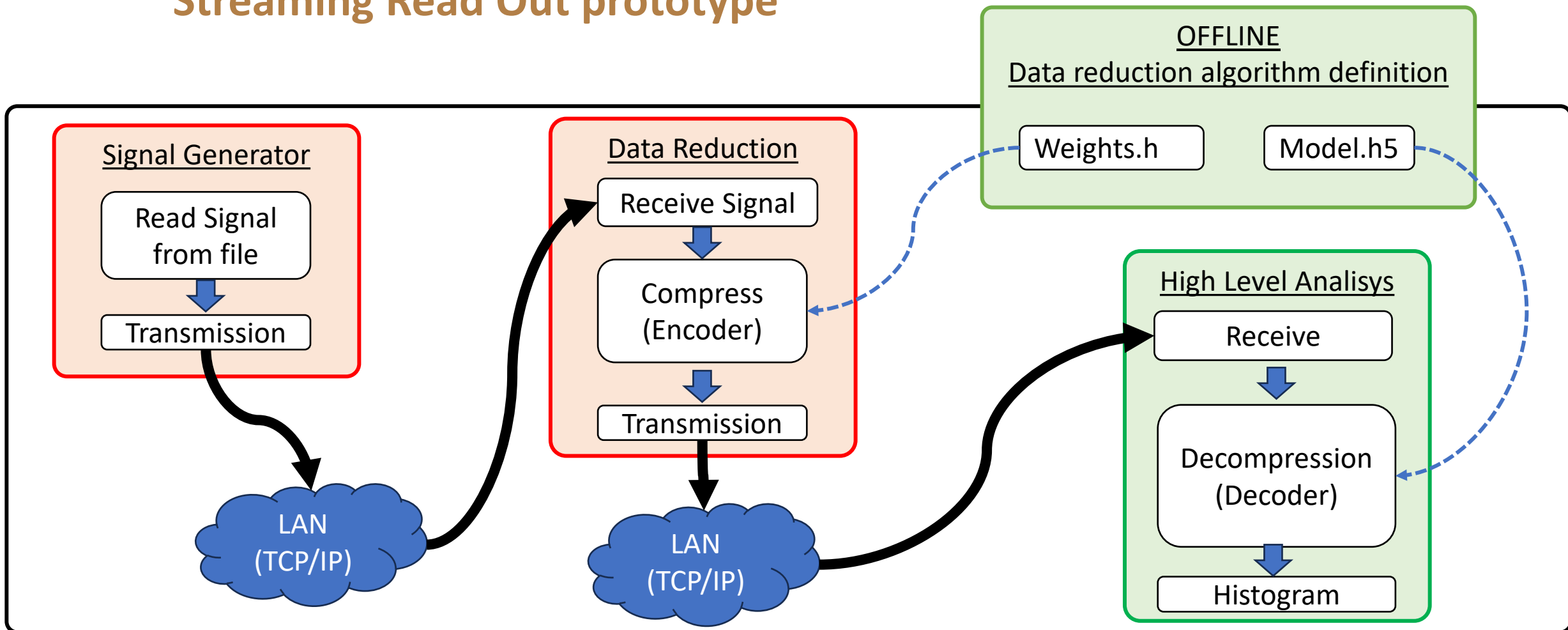


Better reconstruction of low frequency





Streaming Read Out prototype

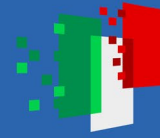




Finanziato dall'Unione europea
NextGenerationEU



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



Data Reduction

Receive Signal

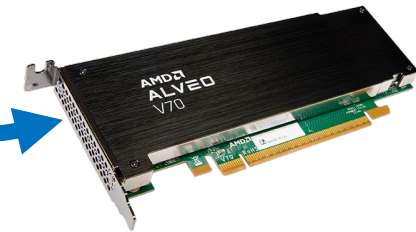
Compress (Encoder)

Transmission

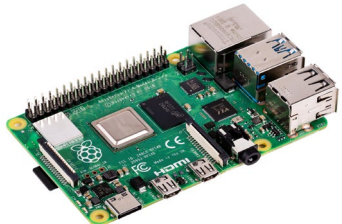
Xilinx XRT



ALVEO V70 FPGA



Raspberry Pi 4 Rev. B

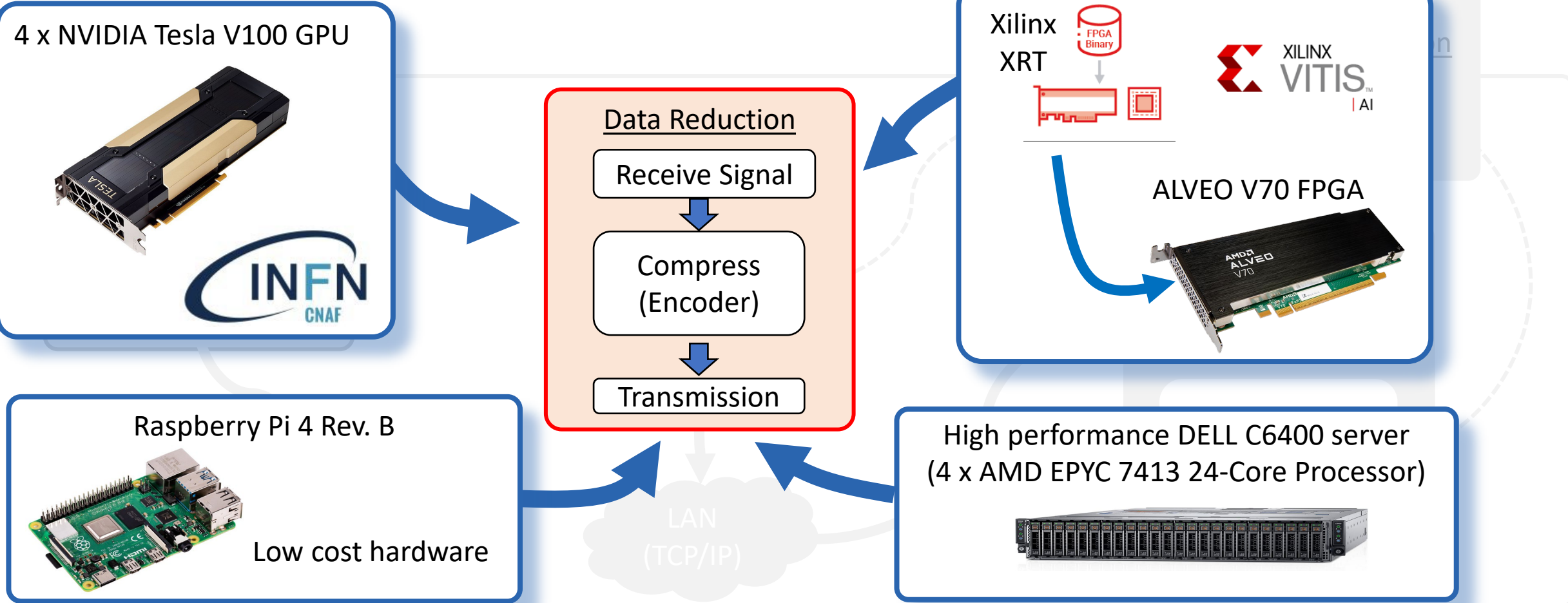


Low cost hardware

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



LAN (TCP/IP)

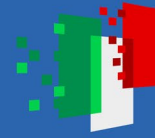




Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



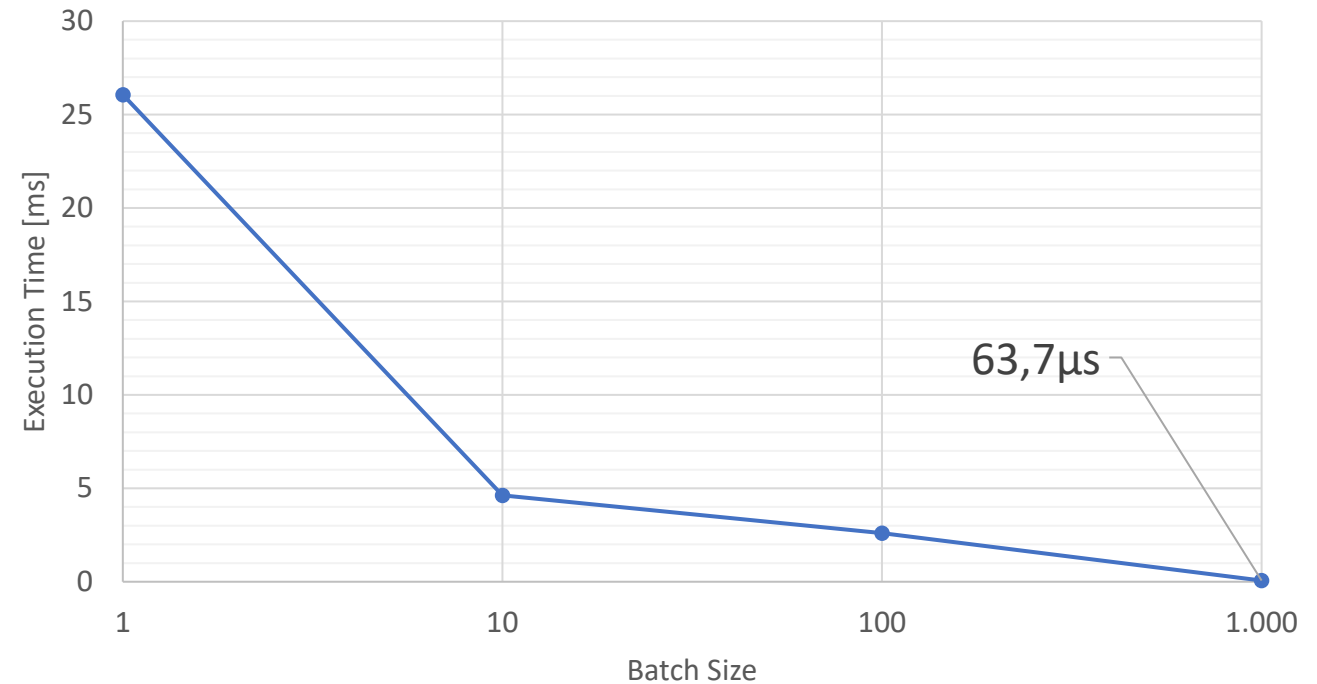
Future
Artificial
Intelligence
Research

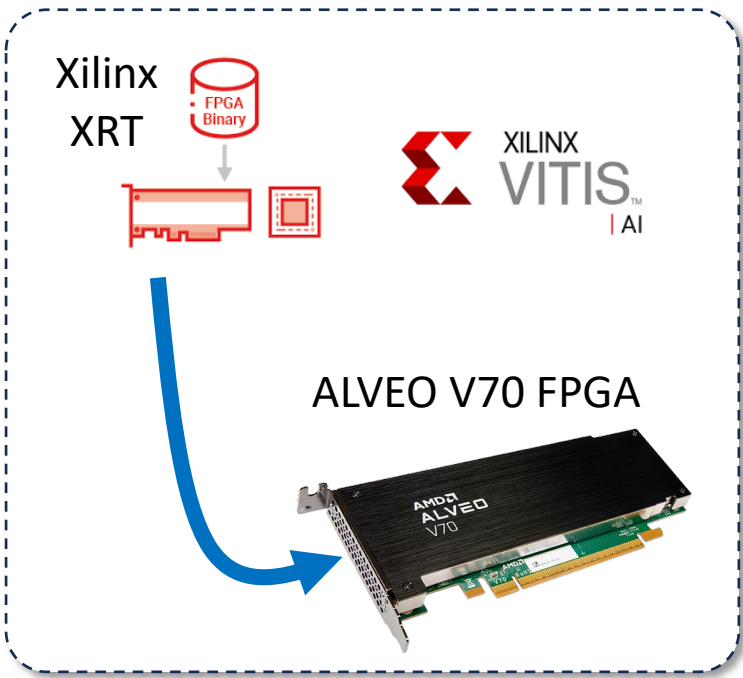
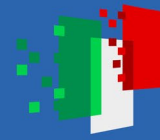
Implementation: GPU

4 x NVIDIA Tesla V100 GPU

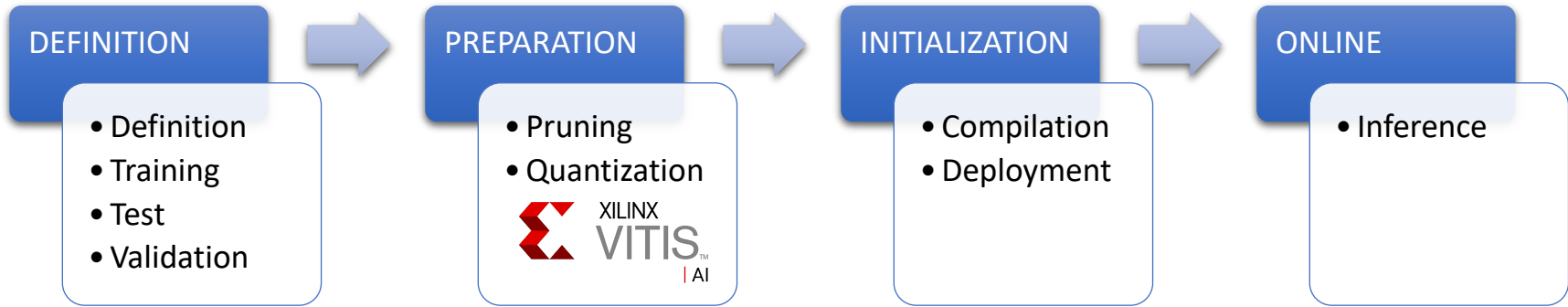


**Execution time not enough
for the application!**



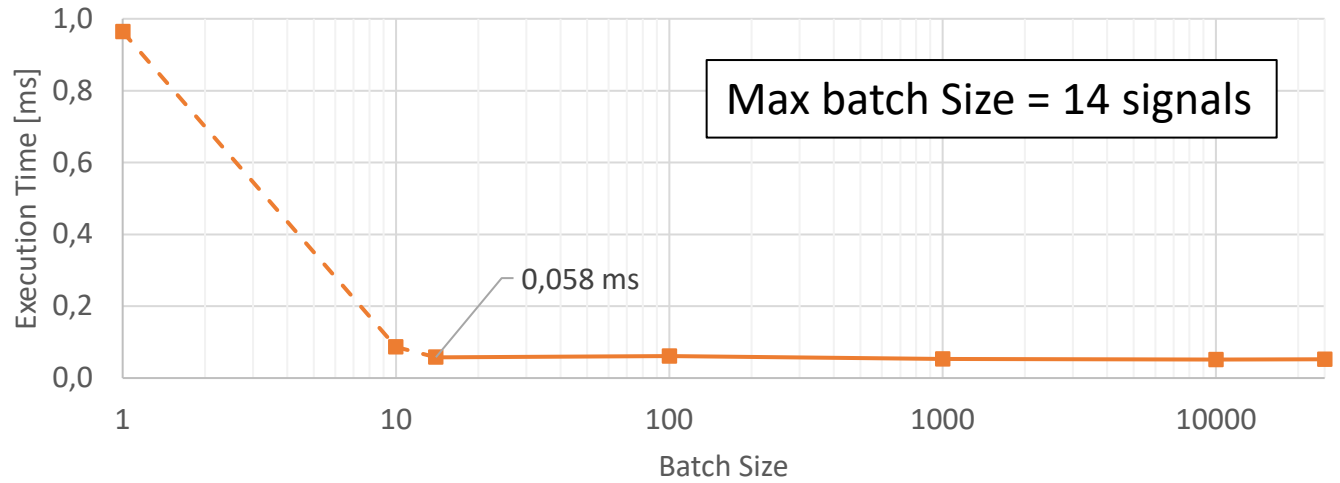


Implementation: FPGA



Execution time still not enough for the application!

Compression time of single signal

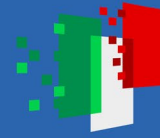




Finanziato dall'Unione europea
NextGenerationEU



Ministero dell'Università e della Ricerca



Italiadomani
PIANO NAZIONALE DI RIPRESA E RESILIENZA



Future Artificial Intelligence Research

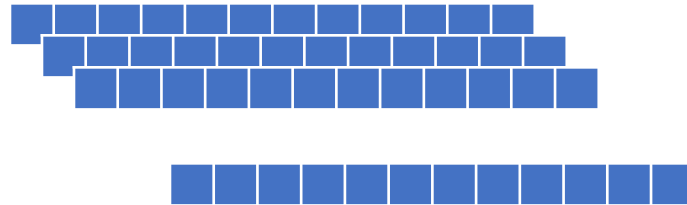
Implementation: High performance server

High performance
DELL C6400 server

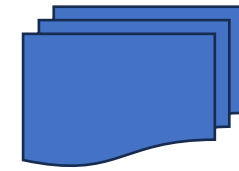


4 x AMD EPYC 7413
24 Core Processor

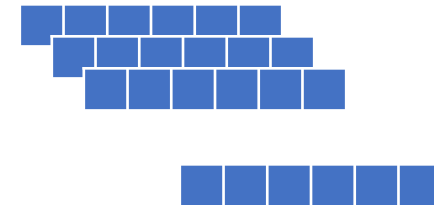
Input Batch



Parallel Execution
(openmp)



Compressed Batch



Single process

Weights

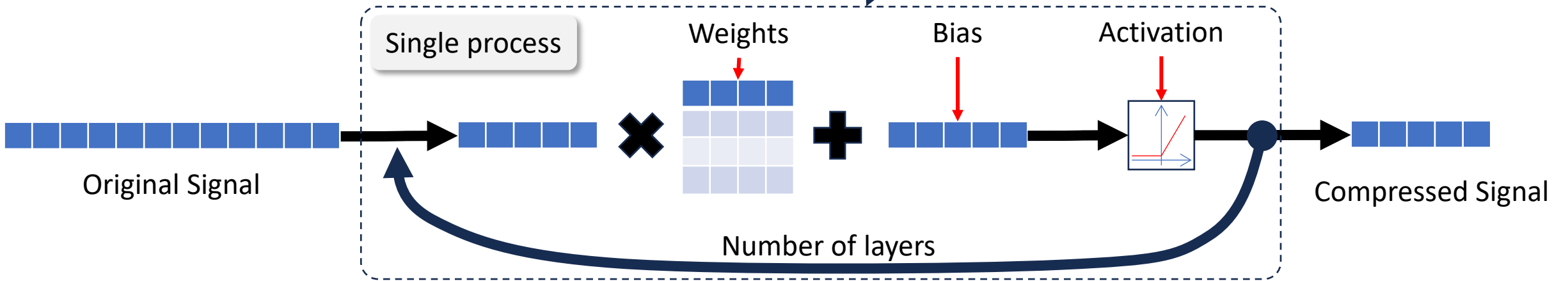
Bias

Activation

Original Signal

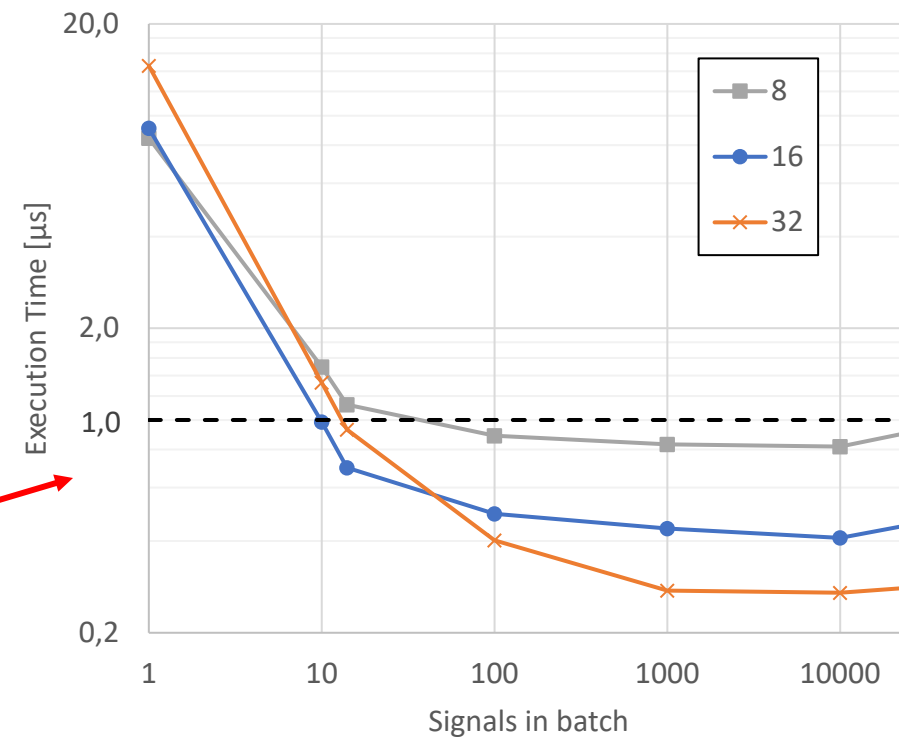
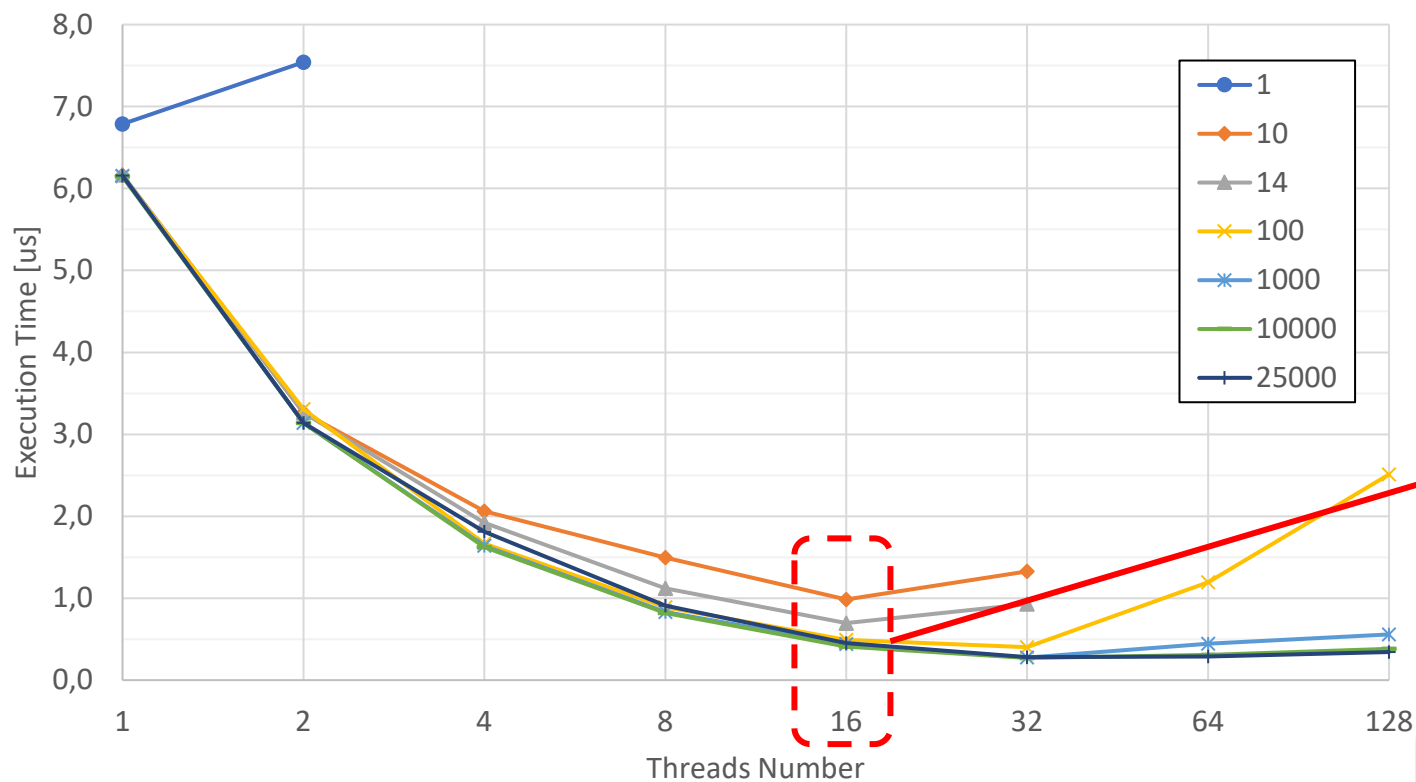
Compressed Signal

Number of layers

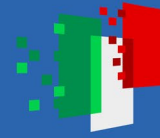


Implementation: High performance server

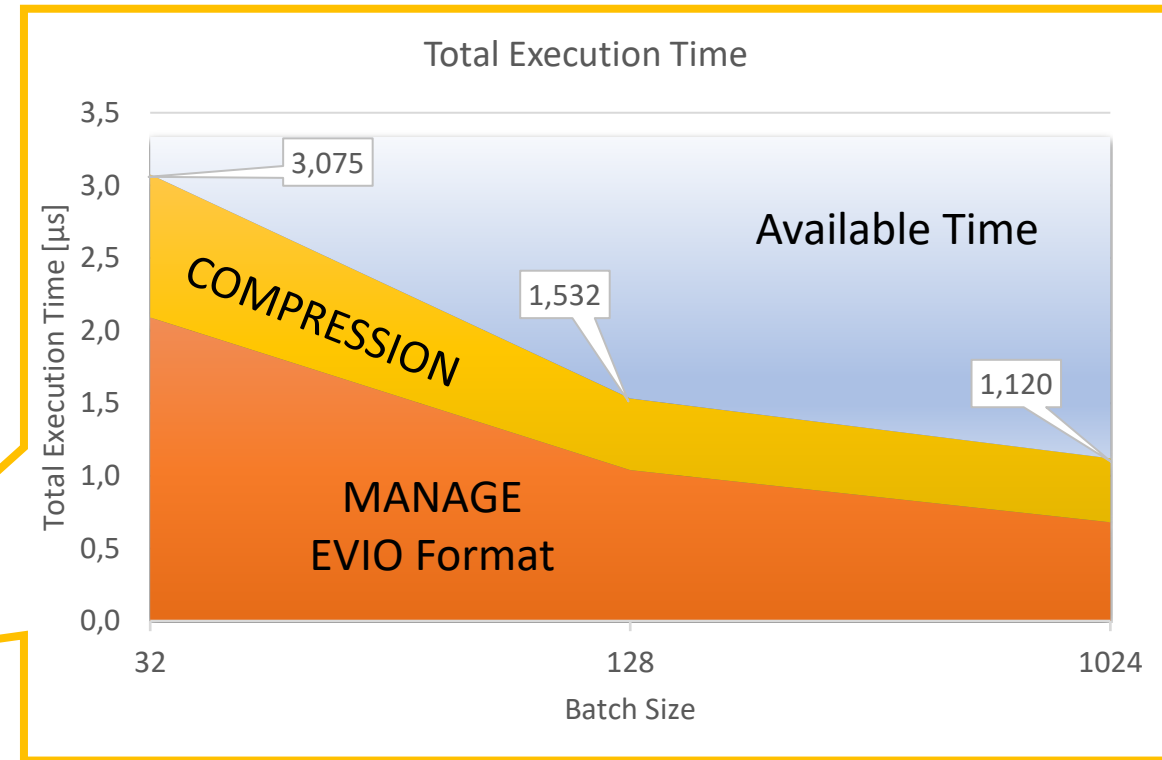
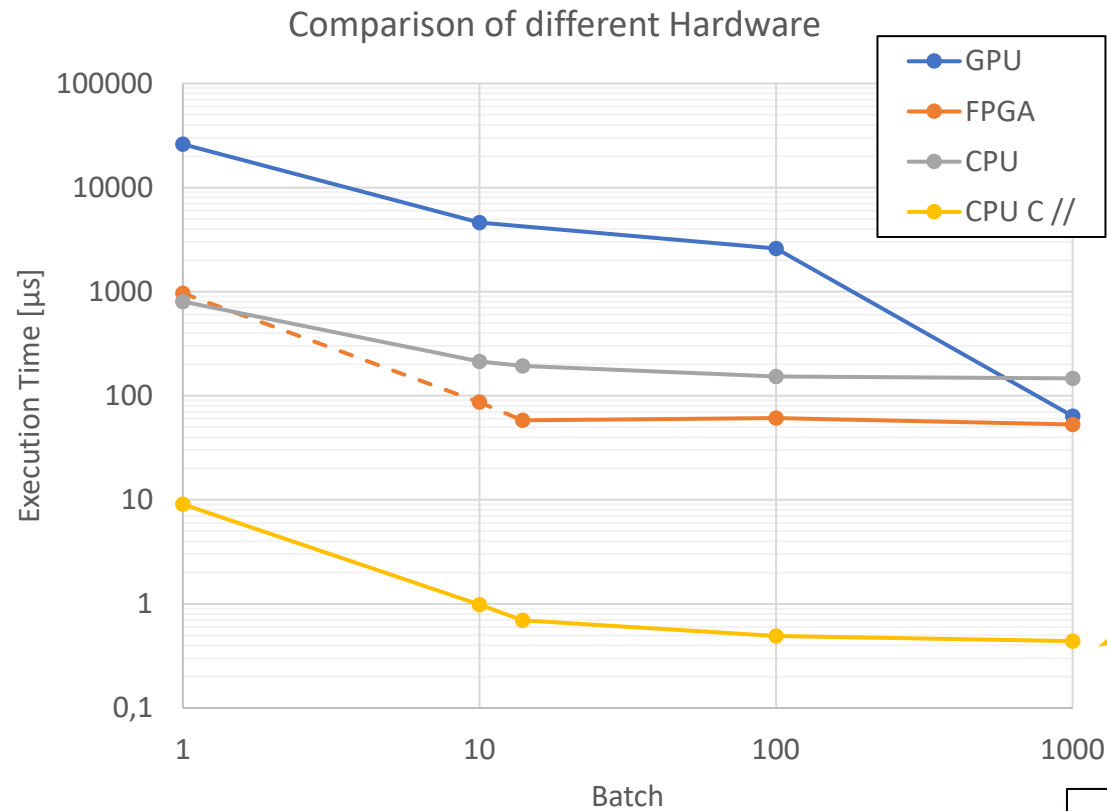
Execution time of different batches and threads number



Chosen 16 Threads
Reasonable execution time for the application



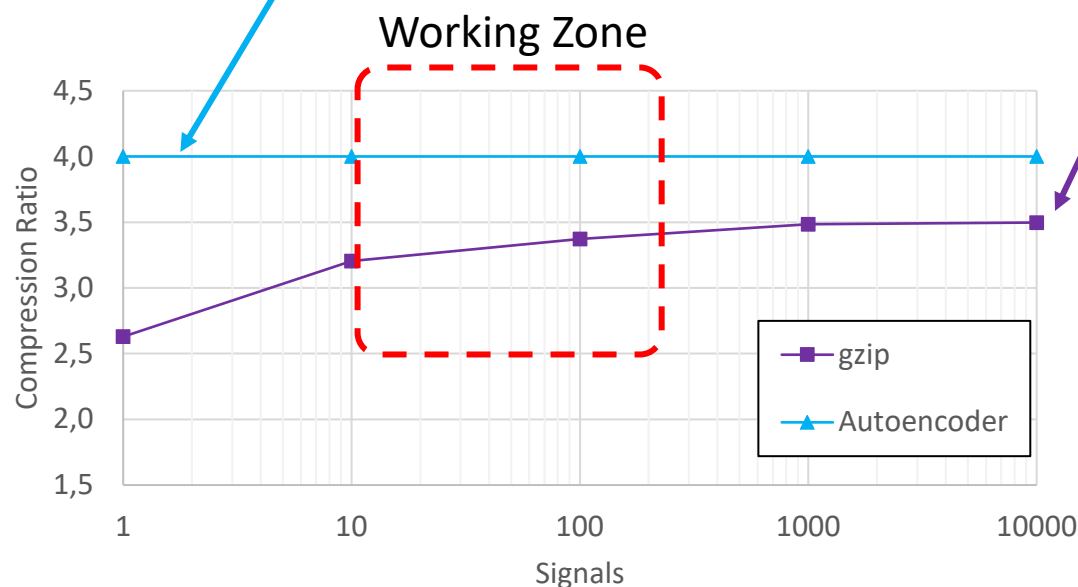
Conclusion



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

Comparison with standard lossless compression

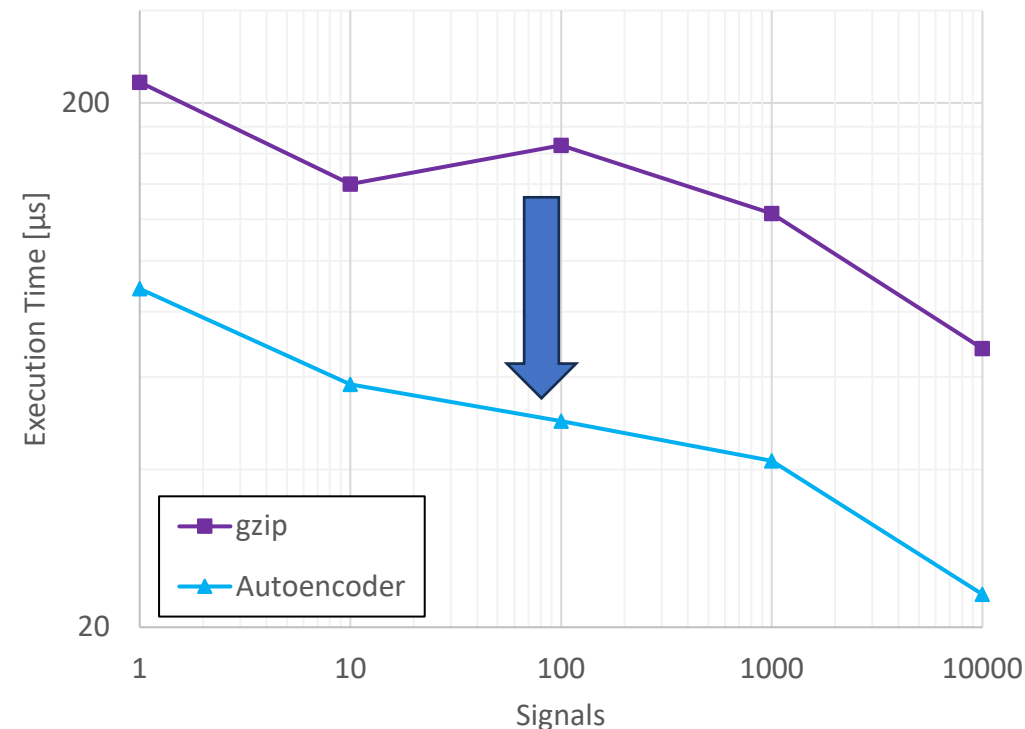
Autoencoder compression ratio is a Parameter



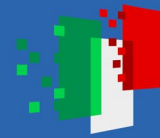
Better compression ratio

Gzip compression ratio depends on signals number

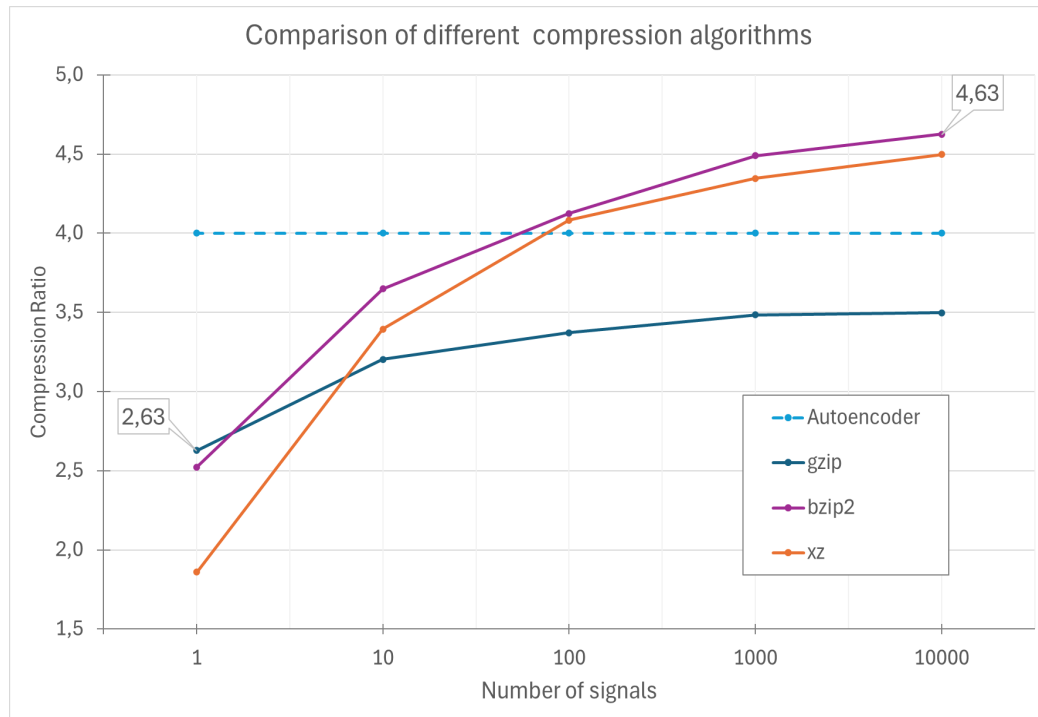
gzip



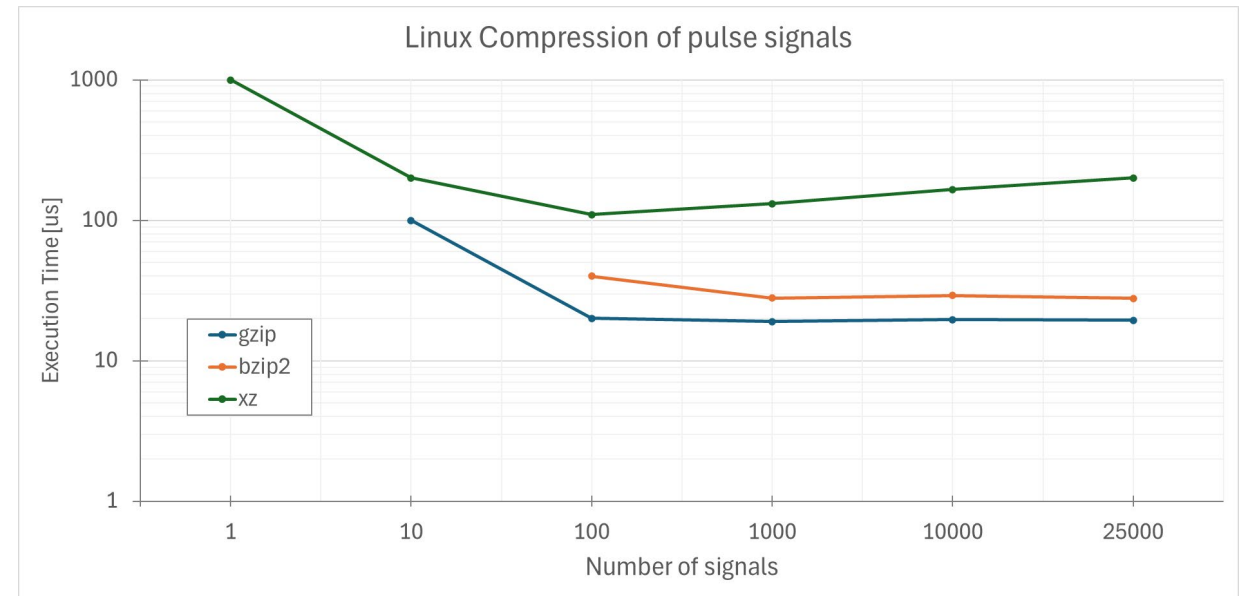
Better also on execution time



Comparison with standard lossless compression



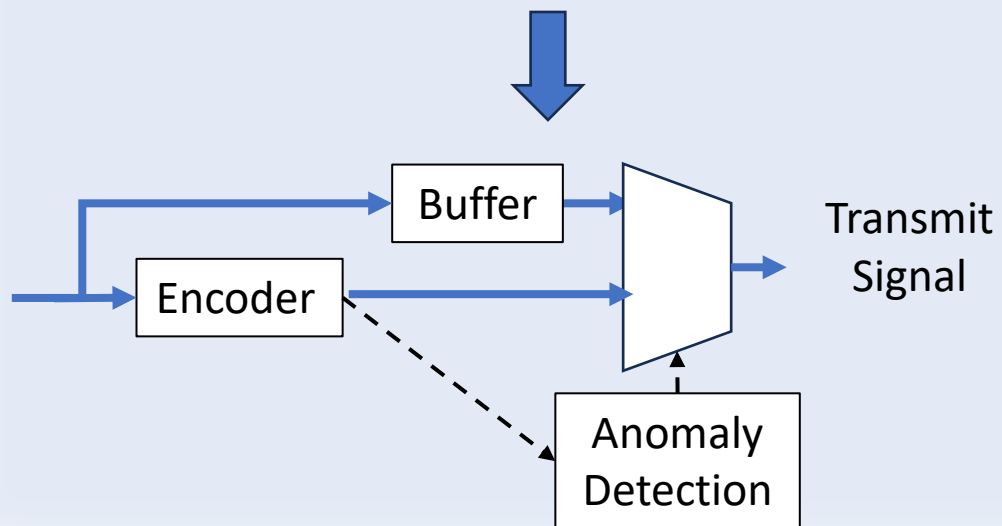
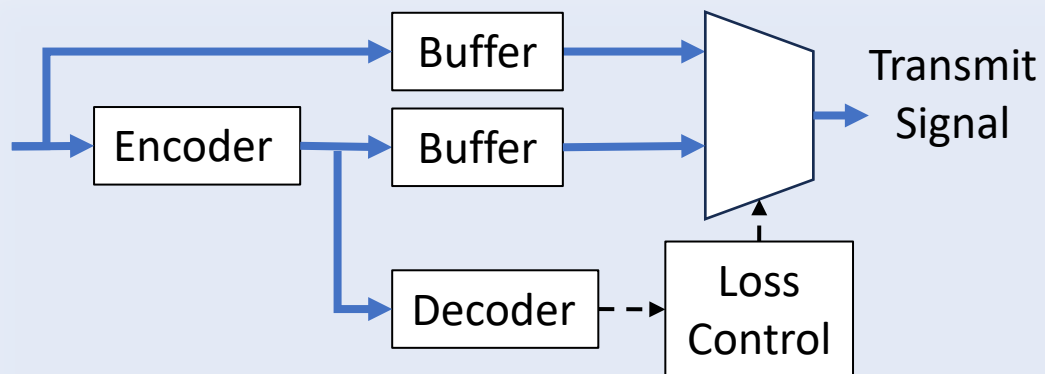
High compression ratio with Bzip2 and X



bzip2
xz

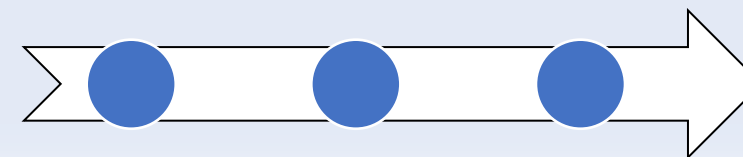
Very poor execution times:
bzip2 50% slower
XZ 10 times slower

Further Studies



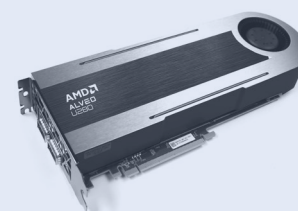
Reduce execution time removing the decoder

Statistical analysis of signals in each EVIO packets



Estimate performance on real acquisition

Low level FPGA implementation



Dedicated connectivity
(2xQSFP28 @ 100GbE)

...or very Low level



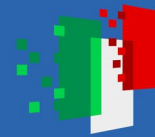
Reduce execution time and maybe save money



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



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.**