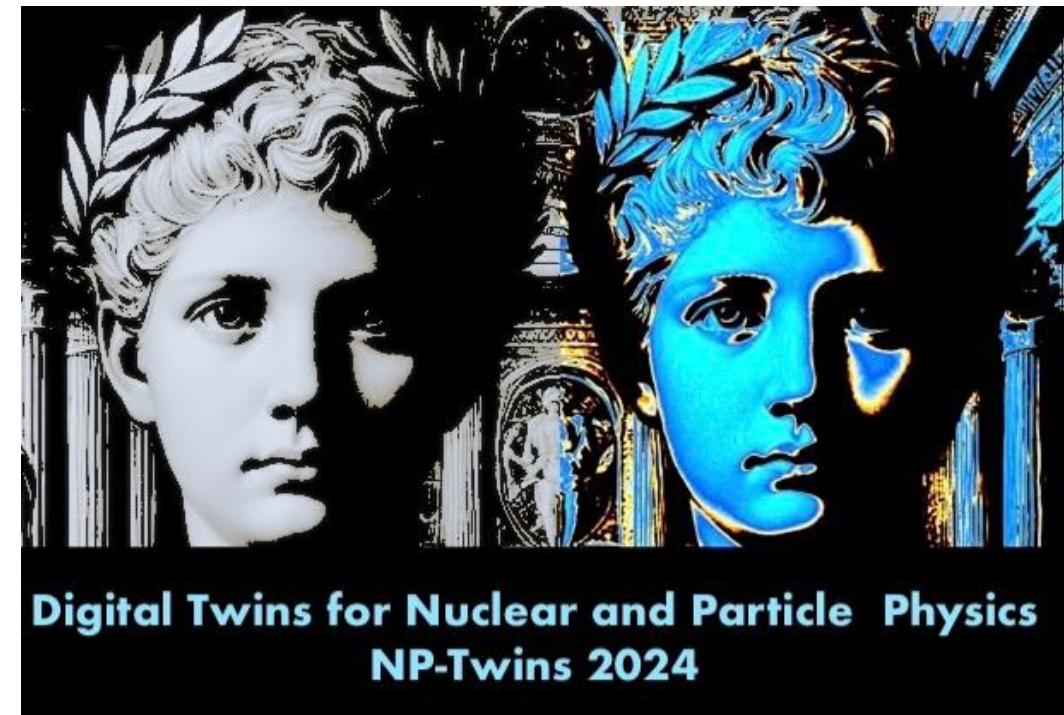# AI for real time data-reduction

## Digital Twins for Nuclear and Particle physics
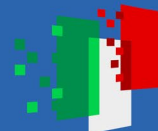
## NPTwins 2024

December 16-18, 2024
Museo Diocesano di Genova

**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**
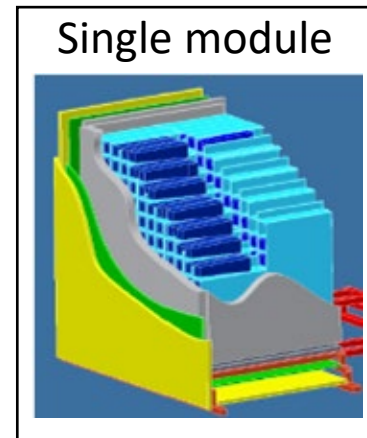Jefferson Lab
Newport News (Virginia)

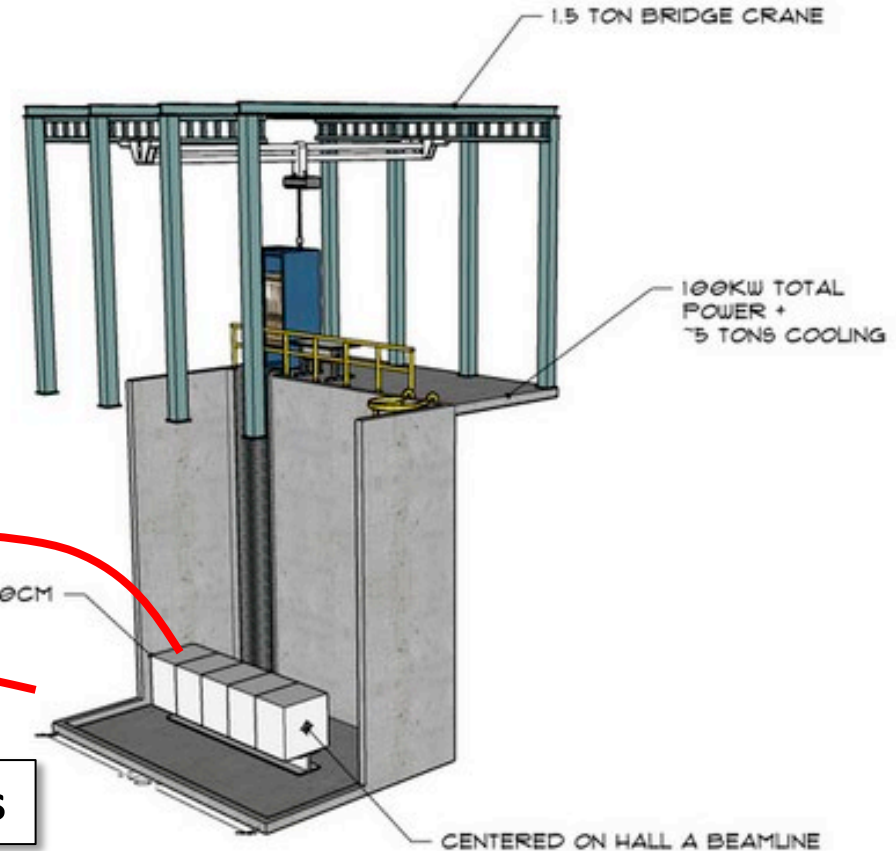# High Energy Physics Experiment: Beam Dump eXperiment (BDX)

**≈1000**
Calorimeter channels
(30MB/s)

**≈ 300**
Veto channels
(500MB/s)

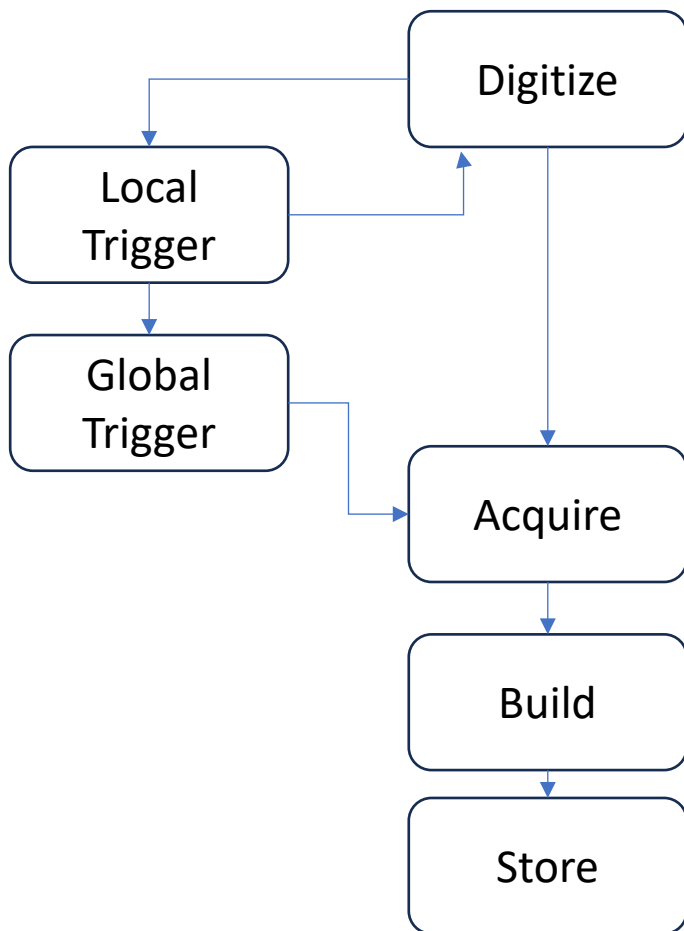Single module

**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).
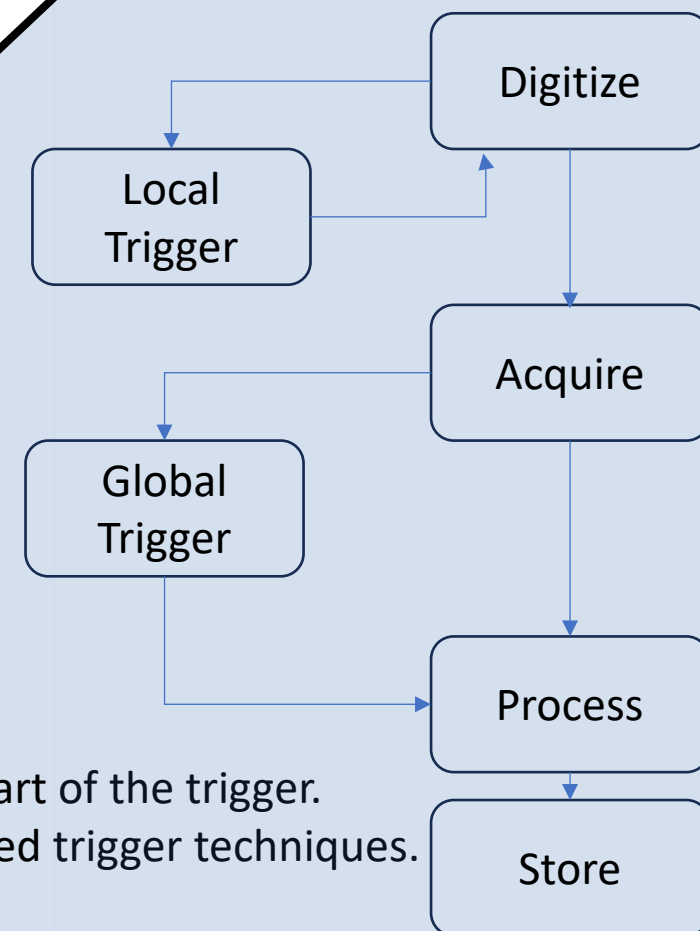
# Traditional triggered DAQ VS Streaming Readout

Digitize

Local Trigger

Global Trigger

Acquire

Build

Store

**Cons:**
Only few information form the trigger.
Trigger logic difficult to implement and debug.
Not easy to adapt to different condition.
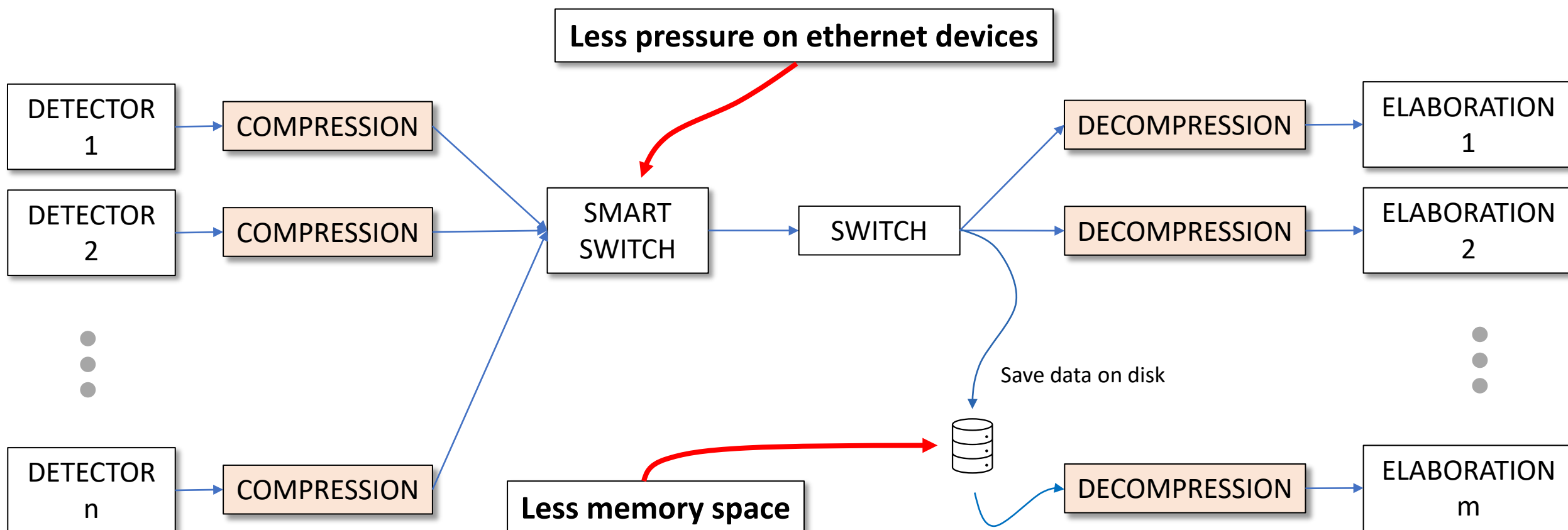
**Pros:**
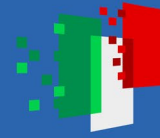It works reliably.

**Triggered**
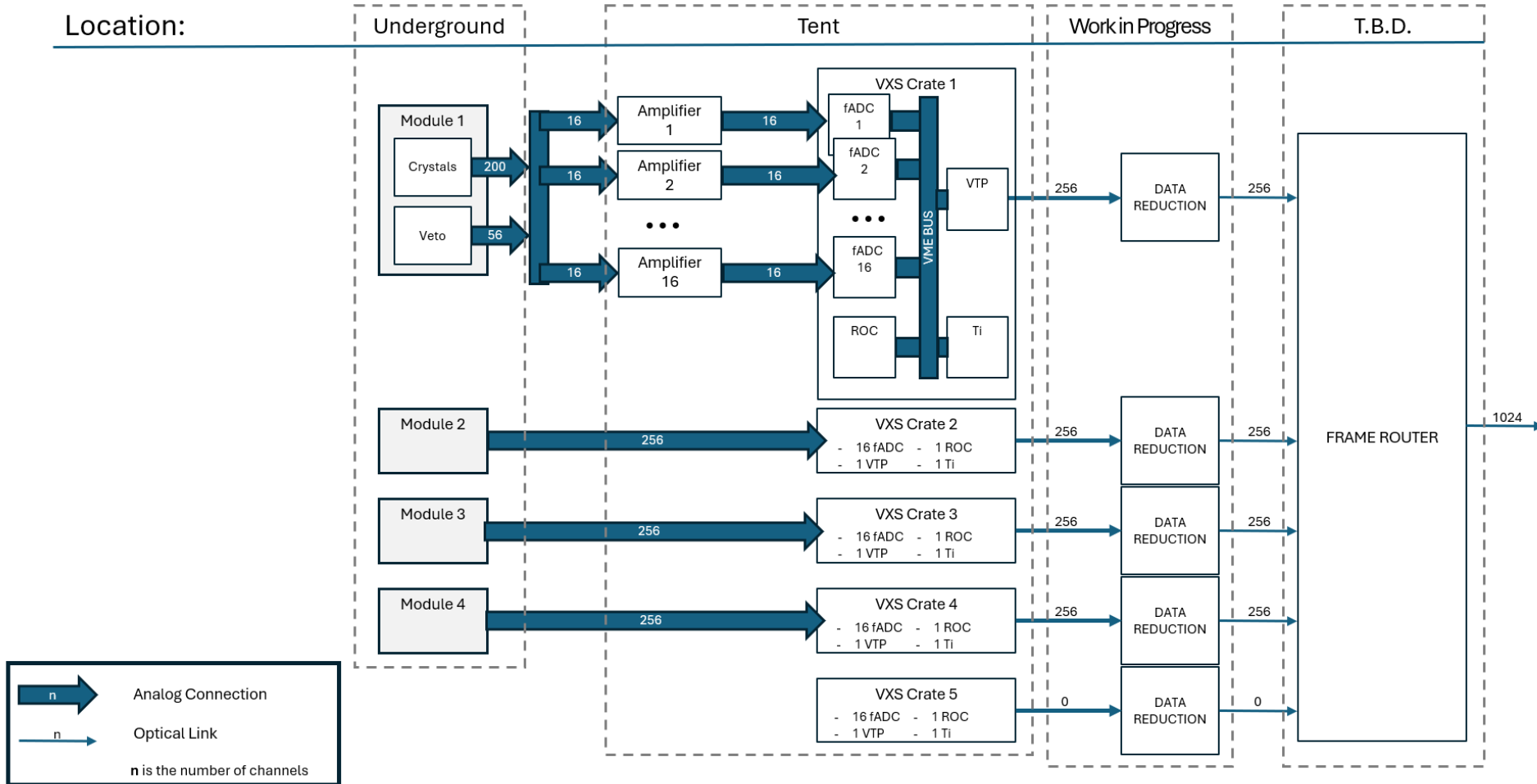
**Streaming**

**Cons:**
**High data rate.**
New design.

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

Digitize

Local Trigger

Acquire

Global Trigger

Process

Store

# Block scheme of data flow

# Detailed BDX data flow scheme

# Data reduction algorithm: Autoencoder

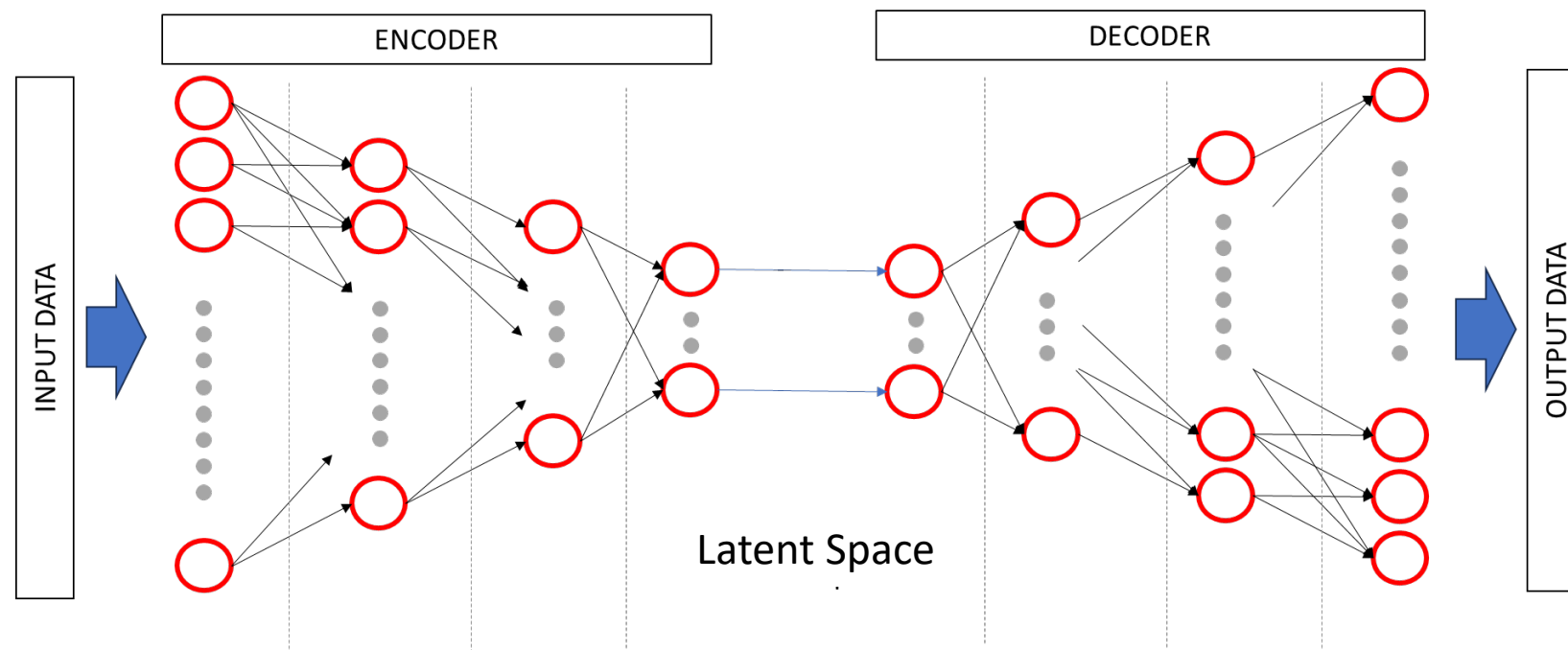Machine Learning Algorithm

Dimensionality reduction

Unsupervised learning
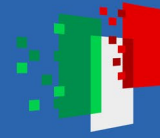
Artificial Neural Network

Composed of two function:
  - encoding
  - decoding

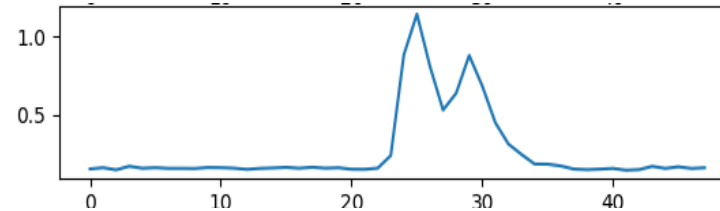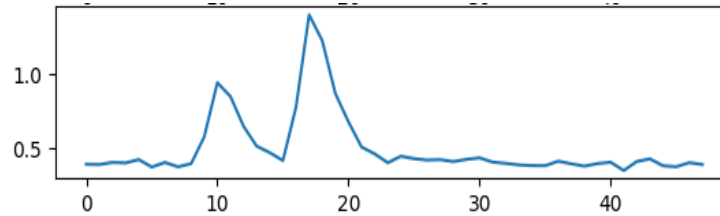FULLY CONNECTED AUTOENCODER WITH DENSE LAYER



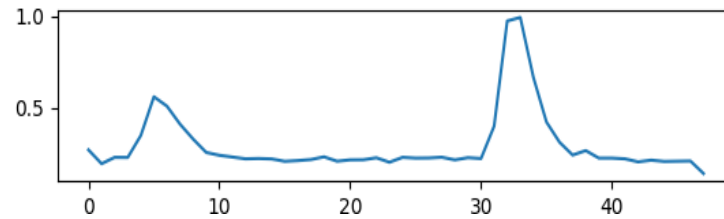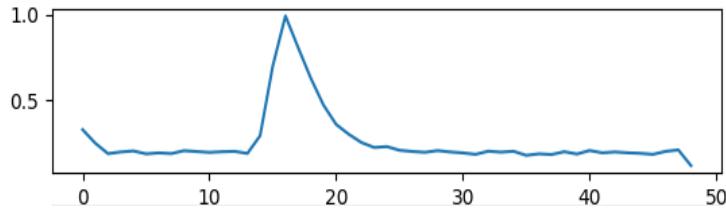ENCODER

DECODER

INPUT DATA

OUTPUT DATA

Latent Space

**Lossy compression algorithm**

# Data from physical Experiment

High

Event Probability

Low

**Very-Low probability signals could be sent uncompressed**

**Data from Experiment 1**

**Data from Experiment 2**

**AE Training**

Weights

**Encoder**

**Encoder***

Weights*

**AE Training**

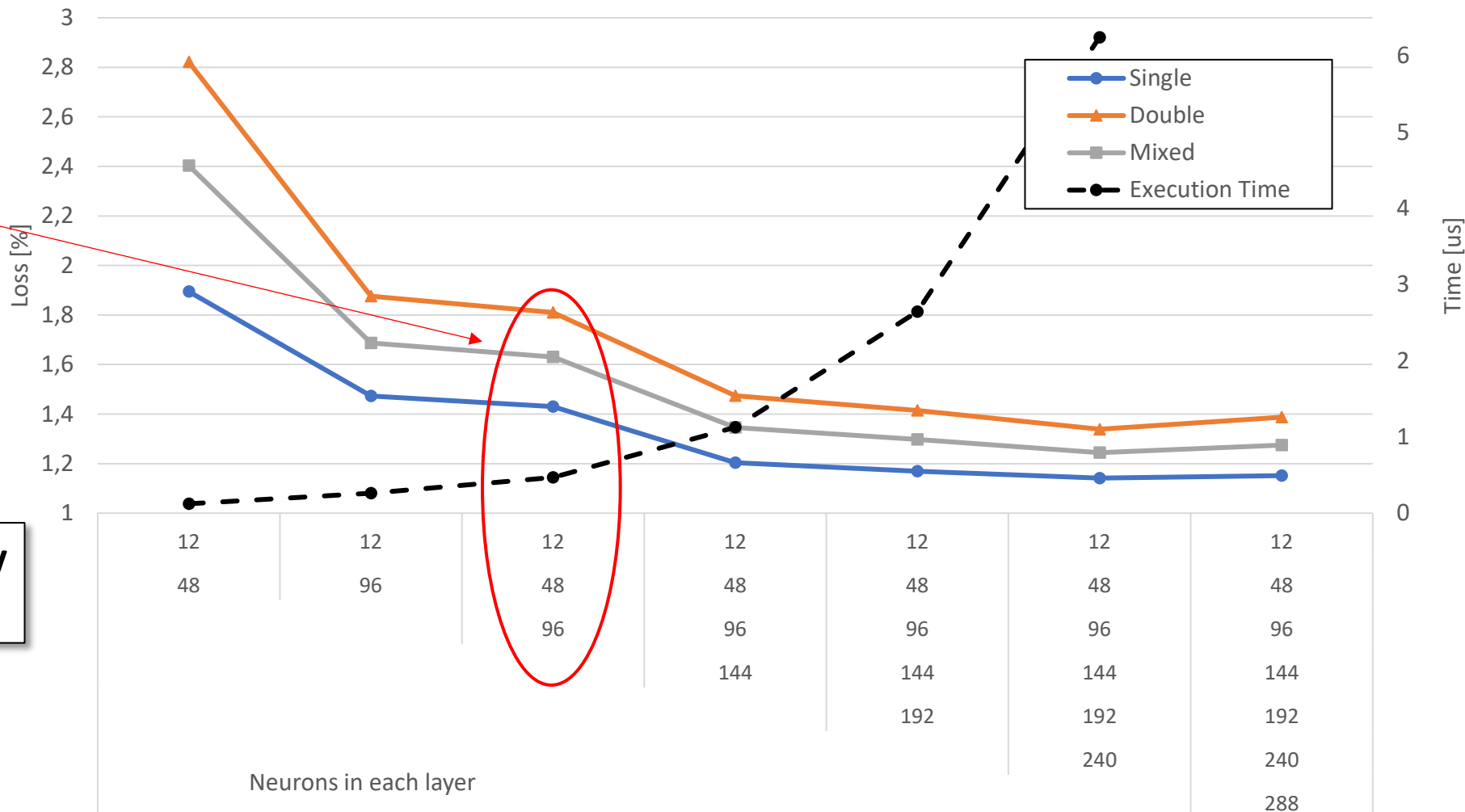Autoencoder Training: Different configuration

Chosen configuration
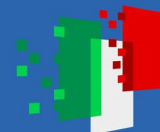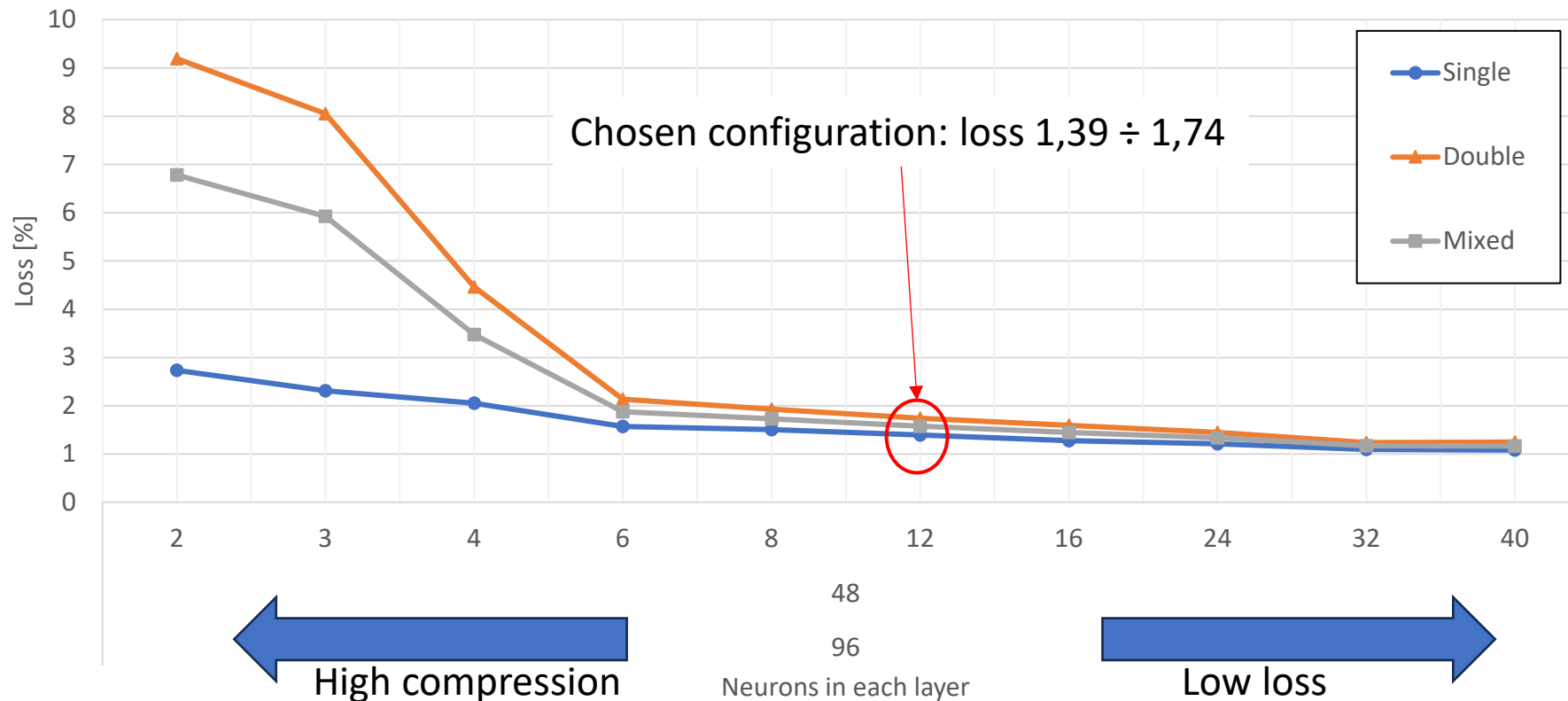
Layers: 3

Neurons:
96
48
12

Execution time increase very fast with model complexity

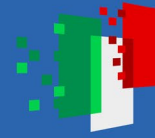Autoencoder Training: Different latent space

Chosen configuration: loss 1,39 ÷ 1,74
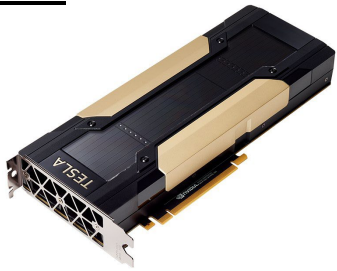
High compression

Low loss

Neurons in each layer

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

# Autoencoder: Training time

# Signals Compression

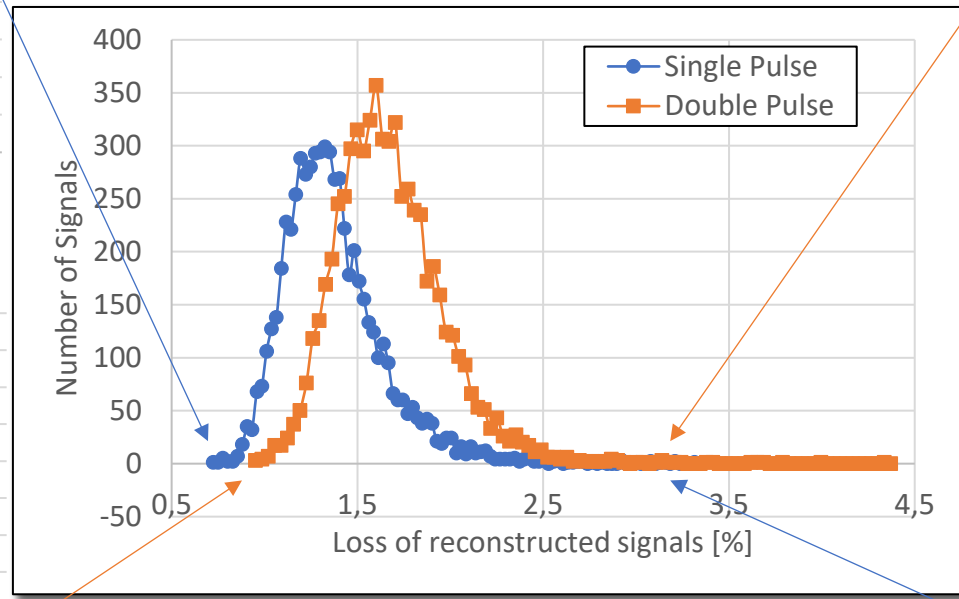# Signals Compression: Integral and spectrum



Integrals histogram
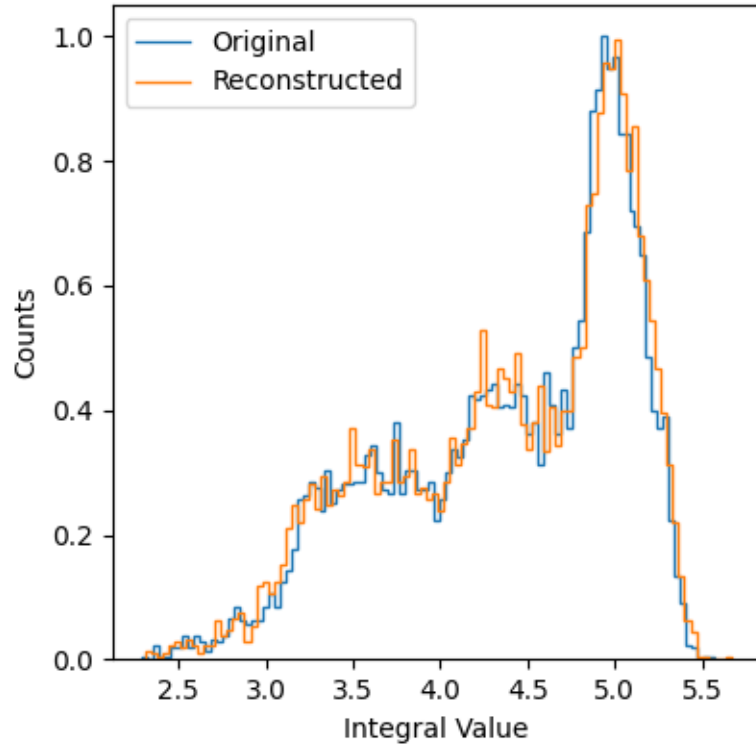


Best Fit

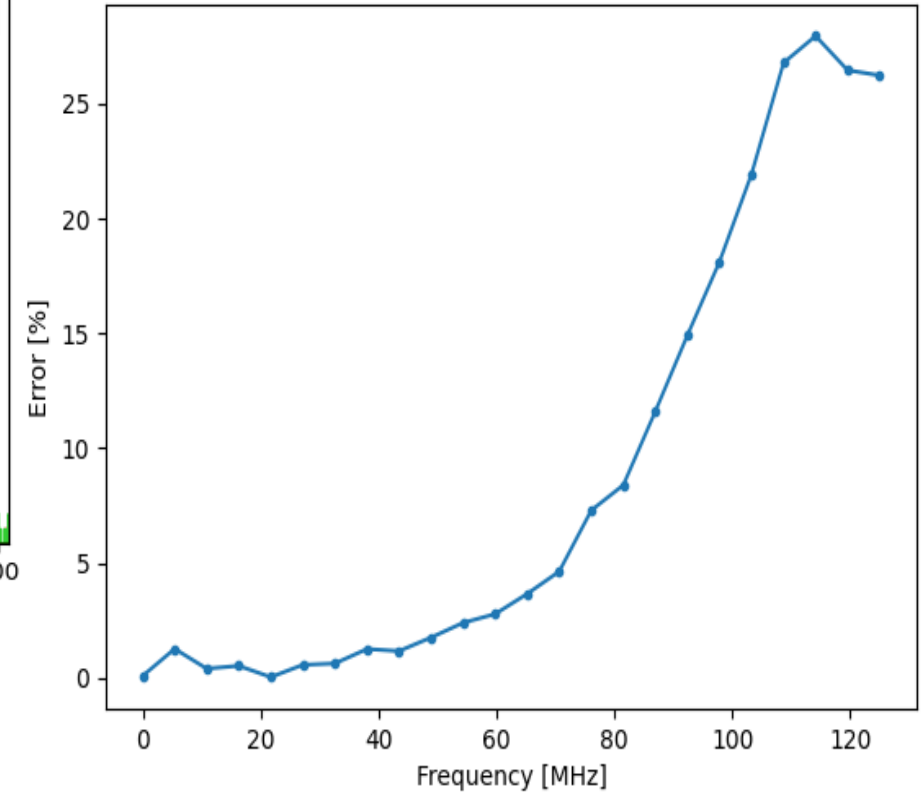|  | A | μ | σ |
|---|---|---|---|
| **Original** | 199,124 | 4,963 | 0,1952 |
| **Recon.** | 199,992 | 4,981 | 0,2006 |
|  |  |  |  |
| **Diff. [%]** | 0,44 | 0,35 | 2,77 |

**Good performance also on the derived quantities for physical analysis**
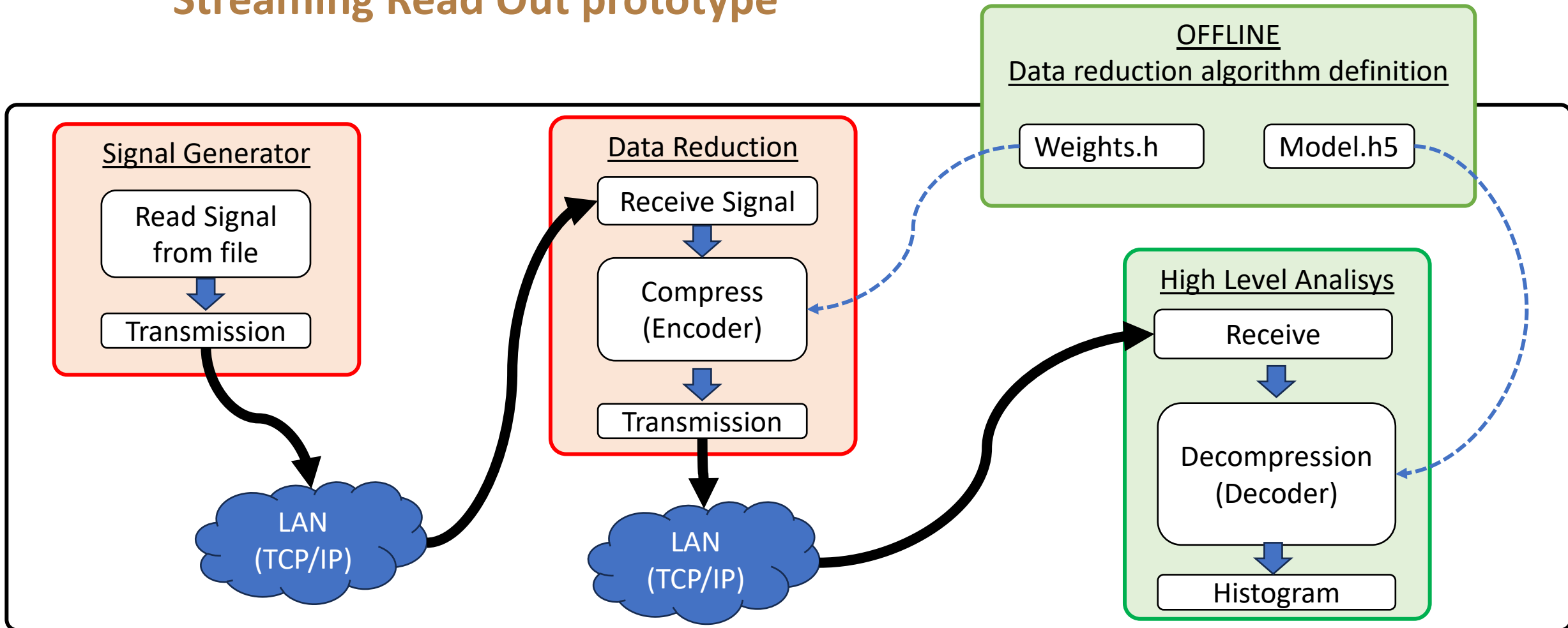
# Signals Compression: FFT analysis

Signals → Encoder → Decoder

Signals → FFT

Decoder → FFT

FFT − FFT → Histogram → Gaussian Fit

**Better reconstruction of low frequency**

# Streaming Read Out prototype



OFFLINE
Data reduction algorithm definition

Weights.h    Model.h5

**Signal Generator**
- Read Signal from file
- Transmission

**Data Reduction**
- Receive Signal
- Compress (Encoder)
- Transmission

**High Level Analisys**
- Receive
- Decompression (Decoder)
- Histogram

LAN (TCP/IP)

LAN (TCP/IP)

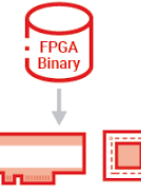# Implementation of Data Reduction Node

4 x NVIDIA Tesla V100 GPU

**Data Reduction**
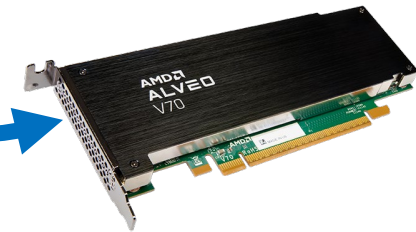
Receive Signal

↓

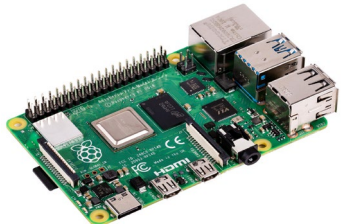Compress
(Encoder)

↓

Transmission

Xilinx XRT

ALVEO V70 FPGA

Raspberry Pi 4 Rev. B

Low cost hardware

LAN
(TCP/IP)

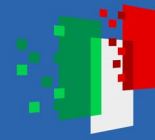High performance DELL C6400 server
(4 x AMD EPYC 7413 24-Core Processor)
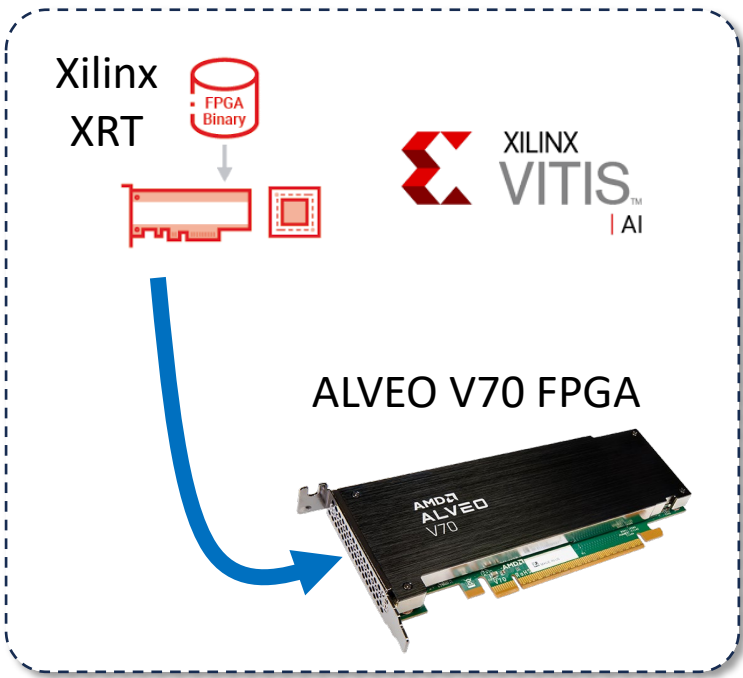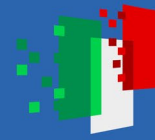
# Implementation: GPU

4 x NVIDIA Tesla V100 GPU

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

63,7μs

# Implementation: FPGA

Xilinx XRT

FPGA Binary

XILINX VITIS | AI

ALVEO V70 FPGA

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

**DEFINITION**
- Definition
- Training
- Test
- Validation

**PREPARATION**
- Pruning
- Quantization

XILINX VITIS | AI

**INITIALIZATION**
- Compilation
- Deployment

**ONLINE**
- Inference

Compression time of single signal

Max batch Size = 14 signals

0,058 ms

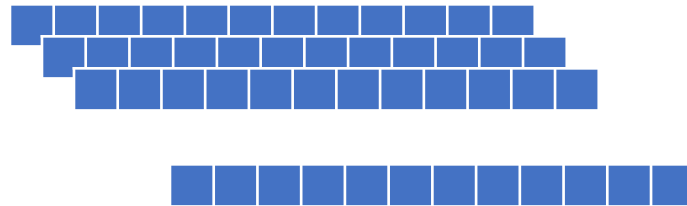Execution Time [ms]

Batch Size

# Implementation: High performance server



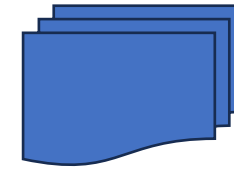High performance
DELL C6400 server

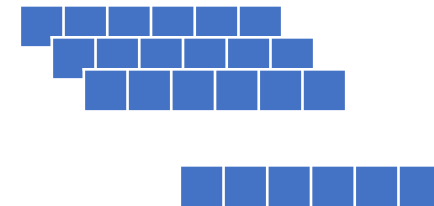4 x AMD EPYC 7413
24 Core Processor

Input Batch

Parallel Execution
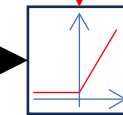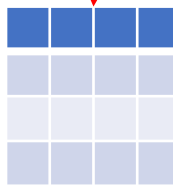(openmp)

Compressed Batch

Single process

Weights

Bias
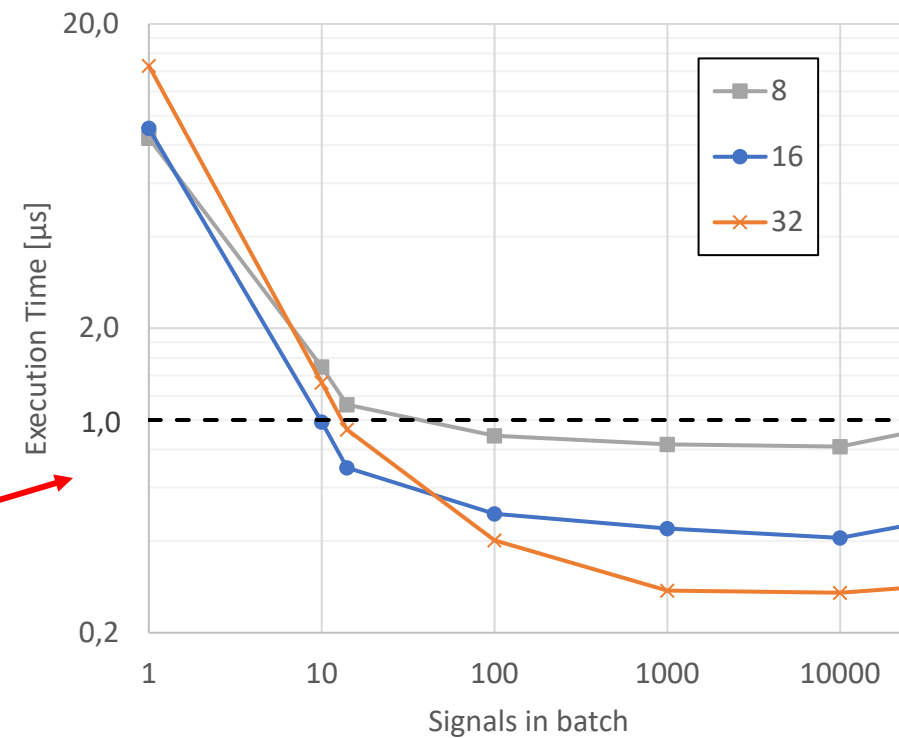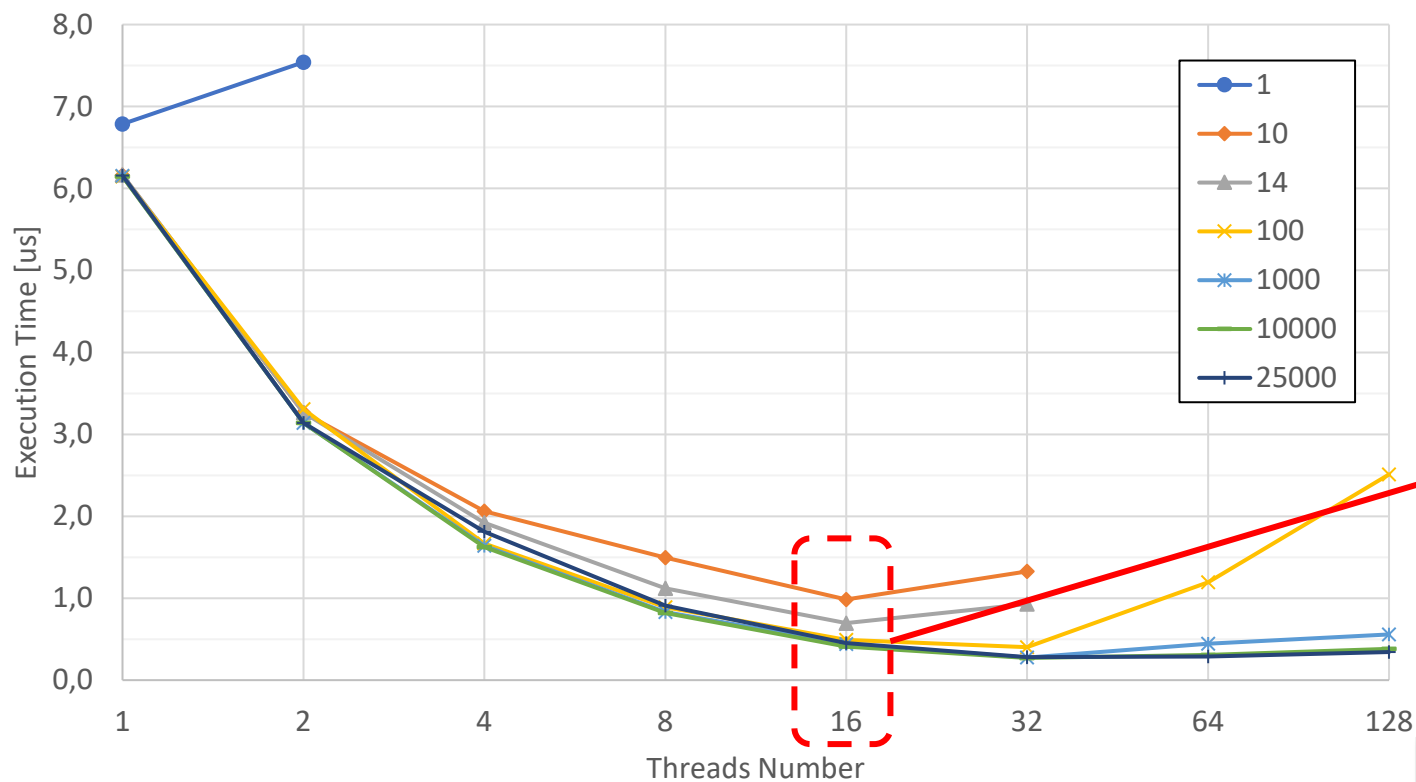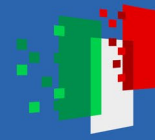
Activation

Original Signal

Compressed Signal

Number of layers

# Implementation: High performance server

Execution time of different batches and threads number



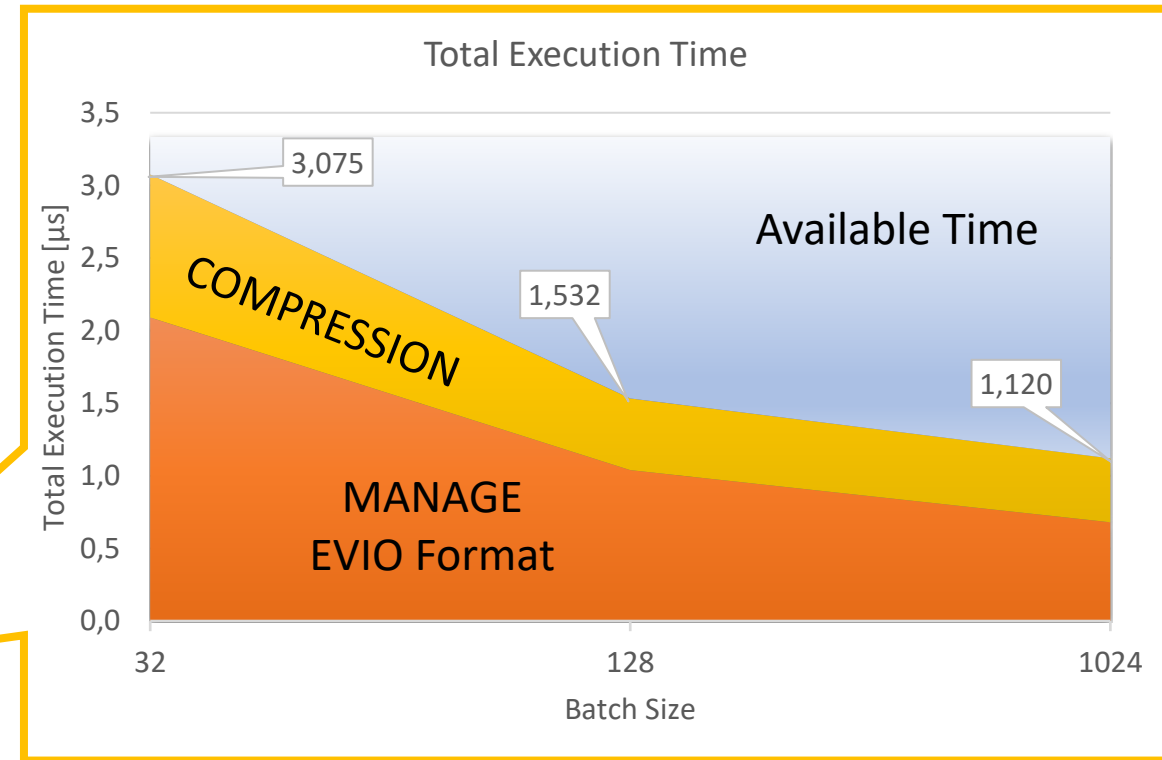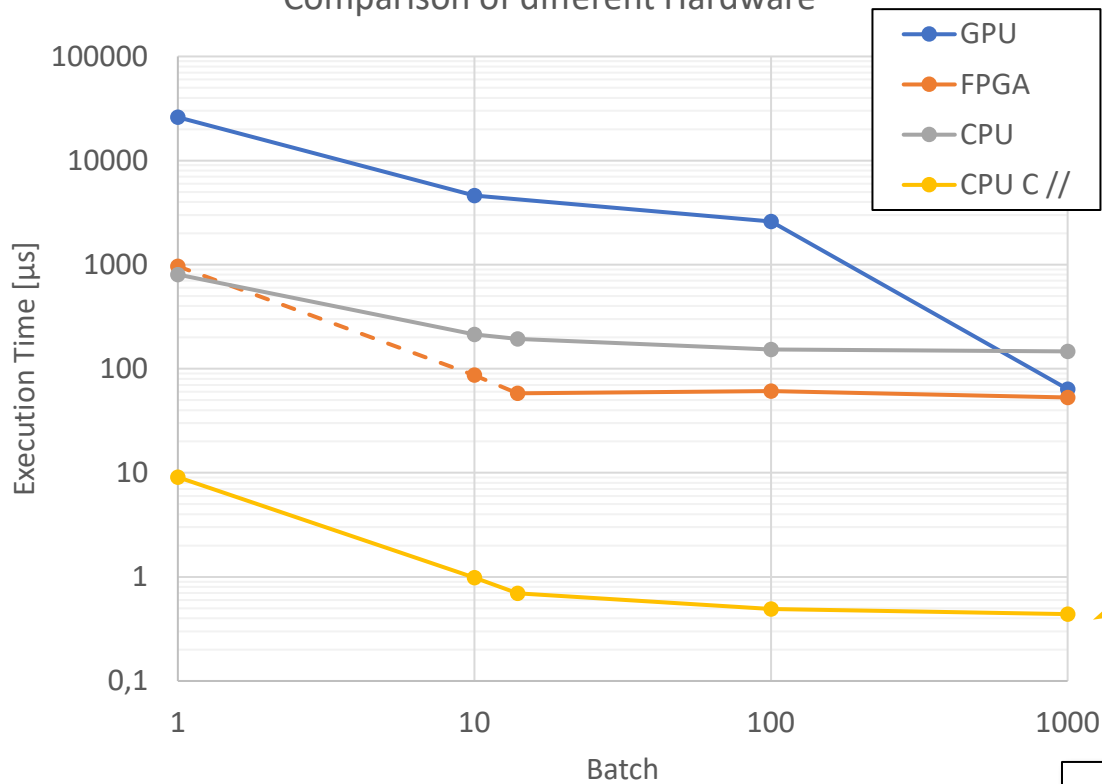**Chosen 16 Threads**
**Reasonable execution time for the application**

Finanziato dall'Unione europea
NextGenerationEU

Ministero dell'Università e della Ricerca

Italiadomani
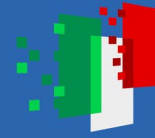PIANO NAZIONALE DI RIPRESA E RESILIENZA

FAIR
Future Artificial Intelligence Research
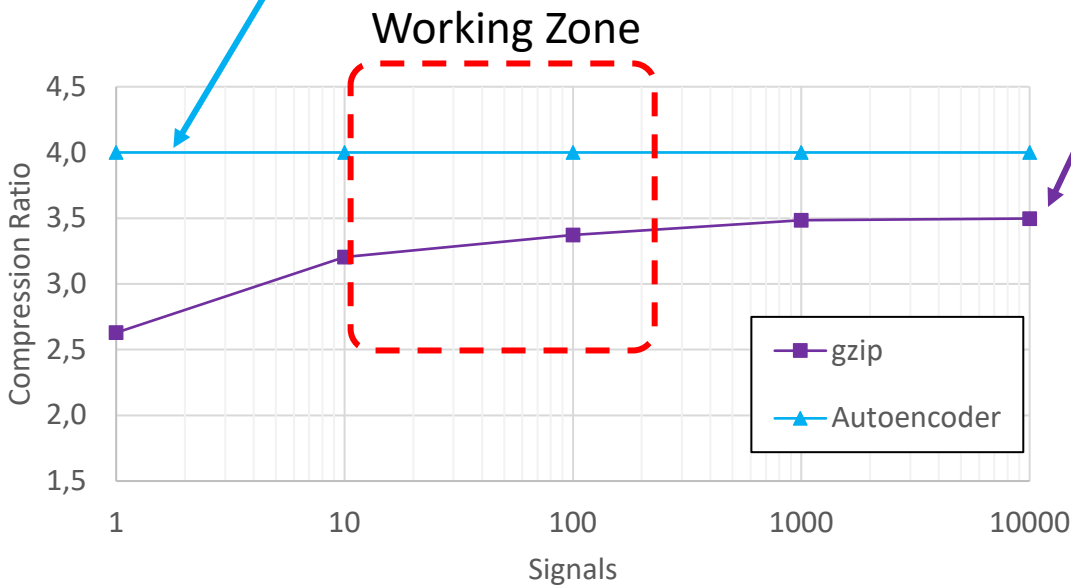
# Conclusion



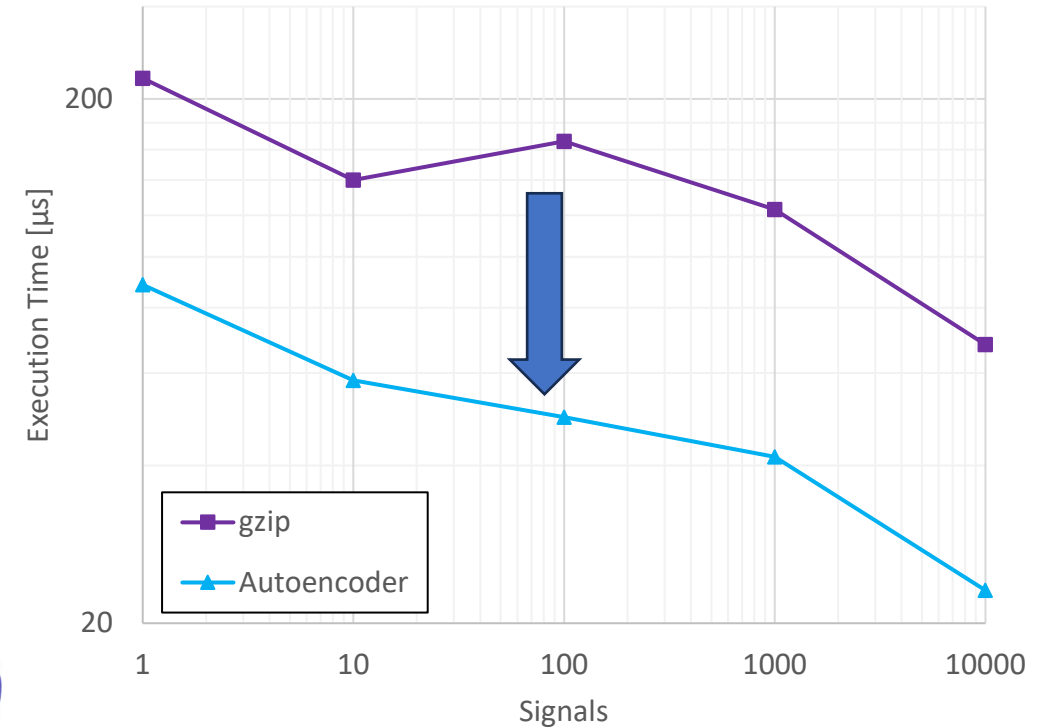**Rate can be managed for EVIO packet with at least 32 signals**

# Comparison with standard lossless compression
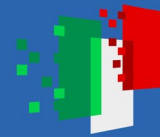
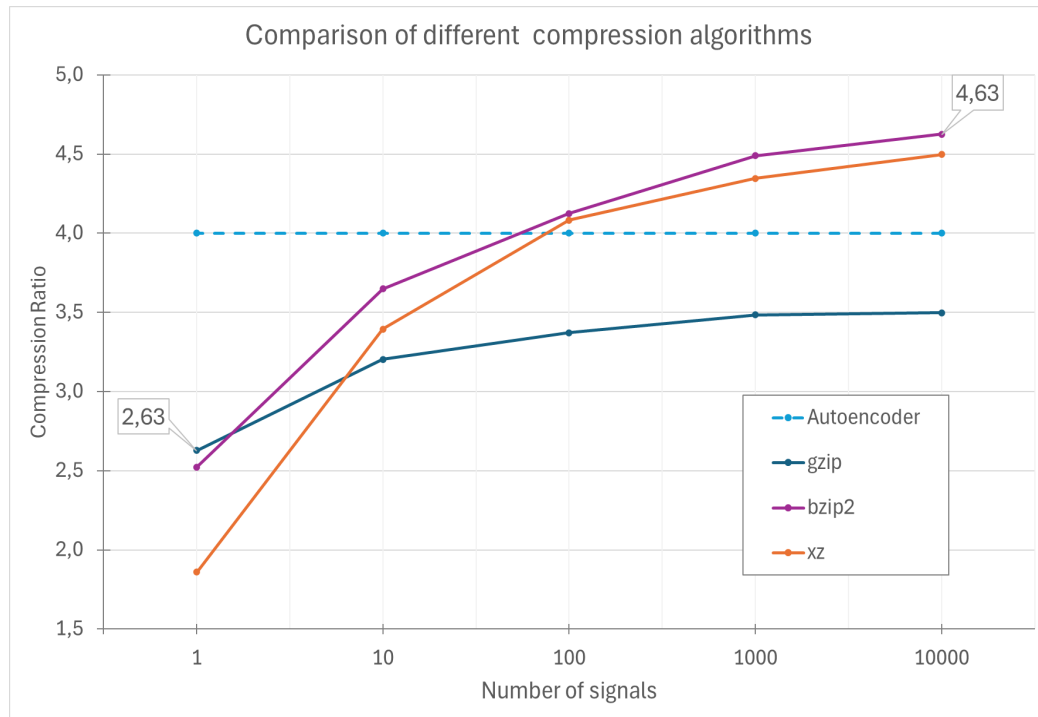Autoencoder compression ratio is a Parameter

Gzip compression ratio depends on signals number

Working Zone



**Better compression ratio**

**Better also on execution time**

# Comparison with standard lossless compression



Comparison of different compression algorithms

Autoencoder
gzip
bzip2
xz

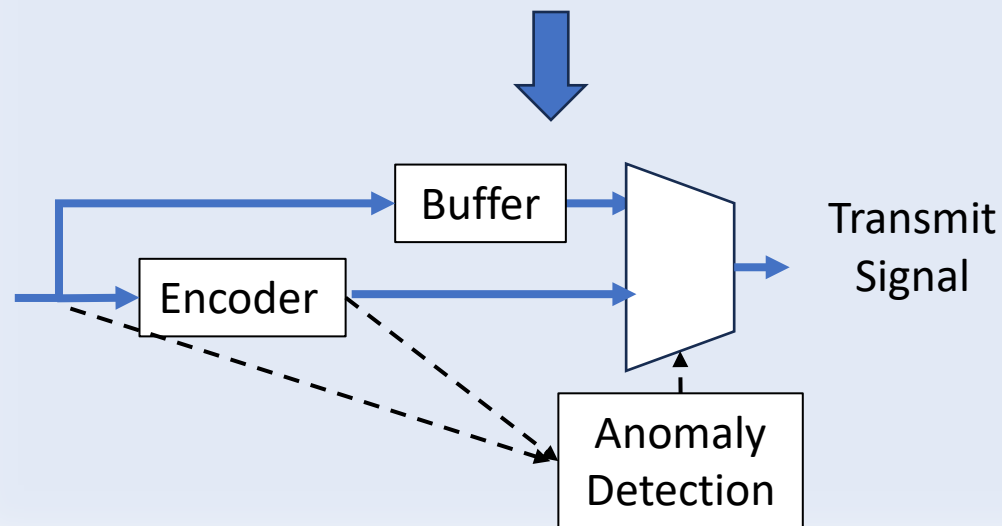Linux Compression of pulse signals
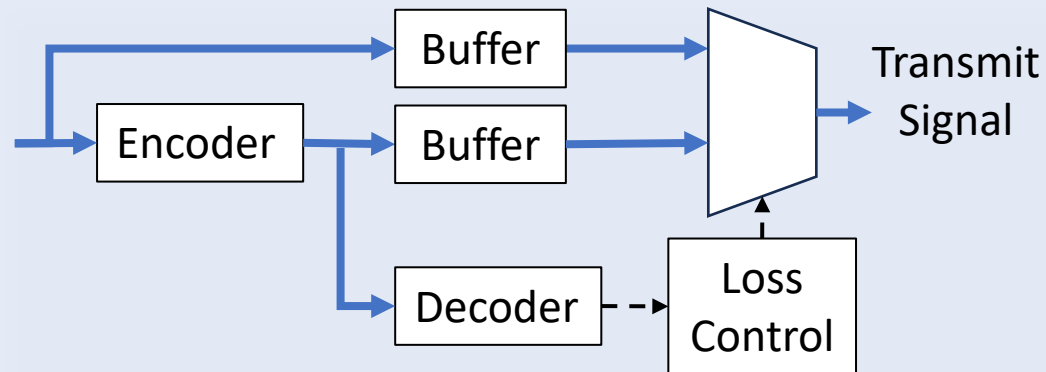
gzip
bzip2
xz

**High compression ratio with Bzip2 and X**
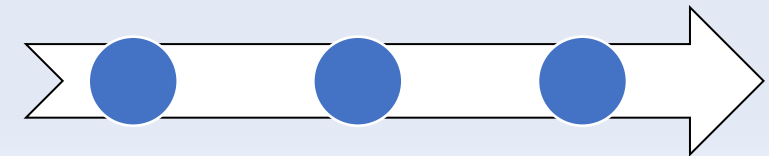
**Very poor execution times:**
bzip2      50% slower
XZ          10 times slower

# 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**

# Thank you for your attention

https://fondazione-fair.it

https://www.jlab.org

https://www.ge.infn.it

https://sealab.unige.it