



# TREASURE

## Future Collaborations

# Data Compression

Marco Montella, with input from:

Akshat Gupta, Sanjiban Sengupta, Caterina Doglioni,  
Antonio Boveia, Thomas Elliott

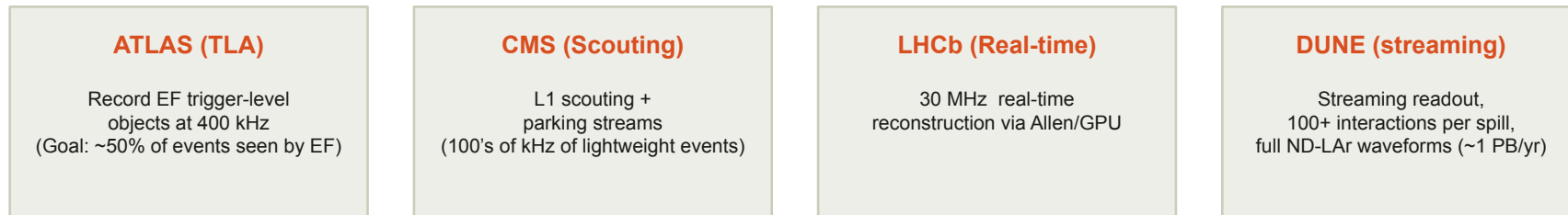
Treasure & Friends Meeting 2026.06.17



# Our Initiative: ML lossless compression

- Explore ML-driven data compression methods to:
  - **Expand Real Time Analysis** on collider data in a fixed-budget model
  - **Reduce data-recording & storage costs** across experiments

## Real Time Analysis dataset “forecast”



- Both Lossy and Lossless compression models under study for ATLAS, CMS, LHCb, and DUNE by a cross-experiment team involving information theory experts
  - As summarised in [Antonio Boveia's talk](#) at April TREASURE Workshop
- **TODAY's GOALS:** → **Outline recent progress**  
→ **Highlight potential areas for collaborations**

# Lossless Compression: Recent Progress

[BOA constrictor: a Mamba-based lossless compressor for scientific data](#)

[Akshat Gupta et al 2026 Mach. Learn.: Sci. Technol. 7 035014](#)

**FRESHLY PUBLISHED!**

[In-depth presentation of work by A. Gupta at CHEP2026](#)

## HIGHLIGHTS

Very competitive compression ratios across data formats / experiments

■ RESULTS - ATLAS

7.23x on ATLAS HDF5, 2.21x on Raw Row-major Jet-only data

<b>HDF5</b> <b>7.23x</b> vs 6.79 LZMA Wins by ~8%. We expect this ratio to be better if we insert a physics informed loss function.	<b>RAW, Jet-only, Row-major</b> <b>2.21x</b> vs 1.70x LZMA Features stored per-jet. ~30% improvement over LZMA.	<b>RAW, Jet-only, Column-major</b> <b>1.82x</b> vs <b>1.92x</b> LZMA Long runs of one variable. Dictionary codes thrive here, BOA struggles.
---	---	--

More work necessary to be competitive with standard libraries

We lag in throughput compared to traditional algorithms.

**3.5-45 MB/s**  
COMPRESSION

**1.5-22 MB/s**  
DECOMPRESSION

■ RESULTS - CMS & GENERALISATION

CMS Raw TBASKET Files – 1500+ Branches

ALGORITHM	EFF. RATIO	COMP.	DECOMP.
LZMA (9)	5.48x	2.8	119.8
BZIP2 (9)	4.87x	25.5	65.8
BROTLI (11)	4.72x	0.9	14.6
ZSTD (19)	4.48x	2.7	1088.3
ZLIB (9)	3.21x	3.4	371.2
<b>BOA (128x2)</b>	<b>5.58x</b>	<b>9.32</b>	<b>6.98</b>



# Short-term directions of work

- **FOUR Summer Students** to work on compression in collaboration with University of Manchester / OSU.

## **Stefan Ionescu (OSU)**

Research: Lossless compression on High-Luminosity ATLAS Trigger-Level data

## **Shanu Sahlot (GSoC / UoM)**

Research: Physics informed bytestream compression implementation

## **Zhengkai Sun (UoM)**

Research: Compression resources vs Disk resources tradeoff

## **- N8CIR**

Research: Sustainable data compression

# Potential Collaboration / Alignment

- **Several areas of potential alignment:**

- Understanding **orthogonality between tokenization and ranged compression**
- Assess the feasibility of a “general map” for BOA with **acceptable cross-experiment performance**
- Compression as a potential **vehicle for open data distribution**
  
- Submitted cross-experiment Phase-I *Genesis* proposal on AI Compression for HEP for broader benchmarking across many datasets, optimization, and entropy floor studies  
**Regardless of review outcome, the core team will pursue this with the aim of joining a Phase-II proposal**

We would be **happy to kickstart a more in-depth discussion** to potentially steer our upcoming work along directions of common interest!