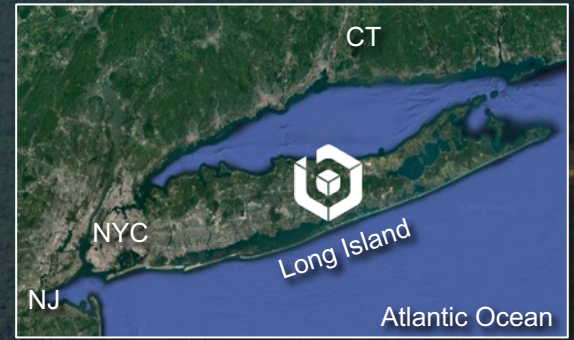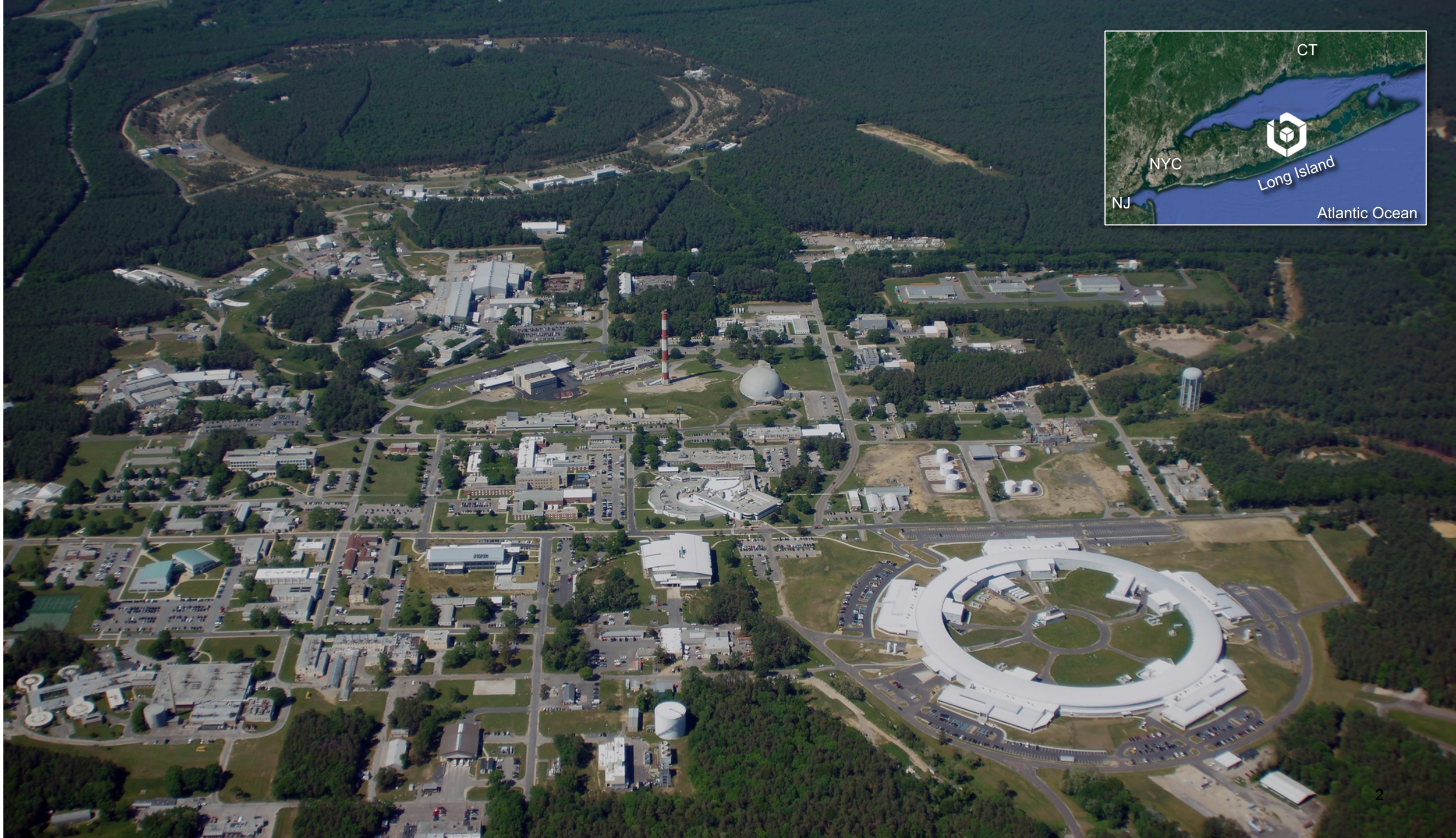# Neural Compression for sPHENIX sparse TPC Data

Yi Huang, Yihui "Ray" Ren, Yeonju Go, Xihaier Luo, Shuhang Li, Thomas Marshall, Joseph D. Osborn, Christopher Pinkenburg, Evgeny Shulga, Shinjae Yoo, Byung-Jun Yoon, Jin Huang (PI)

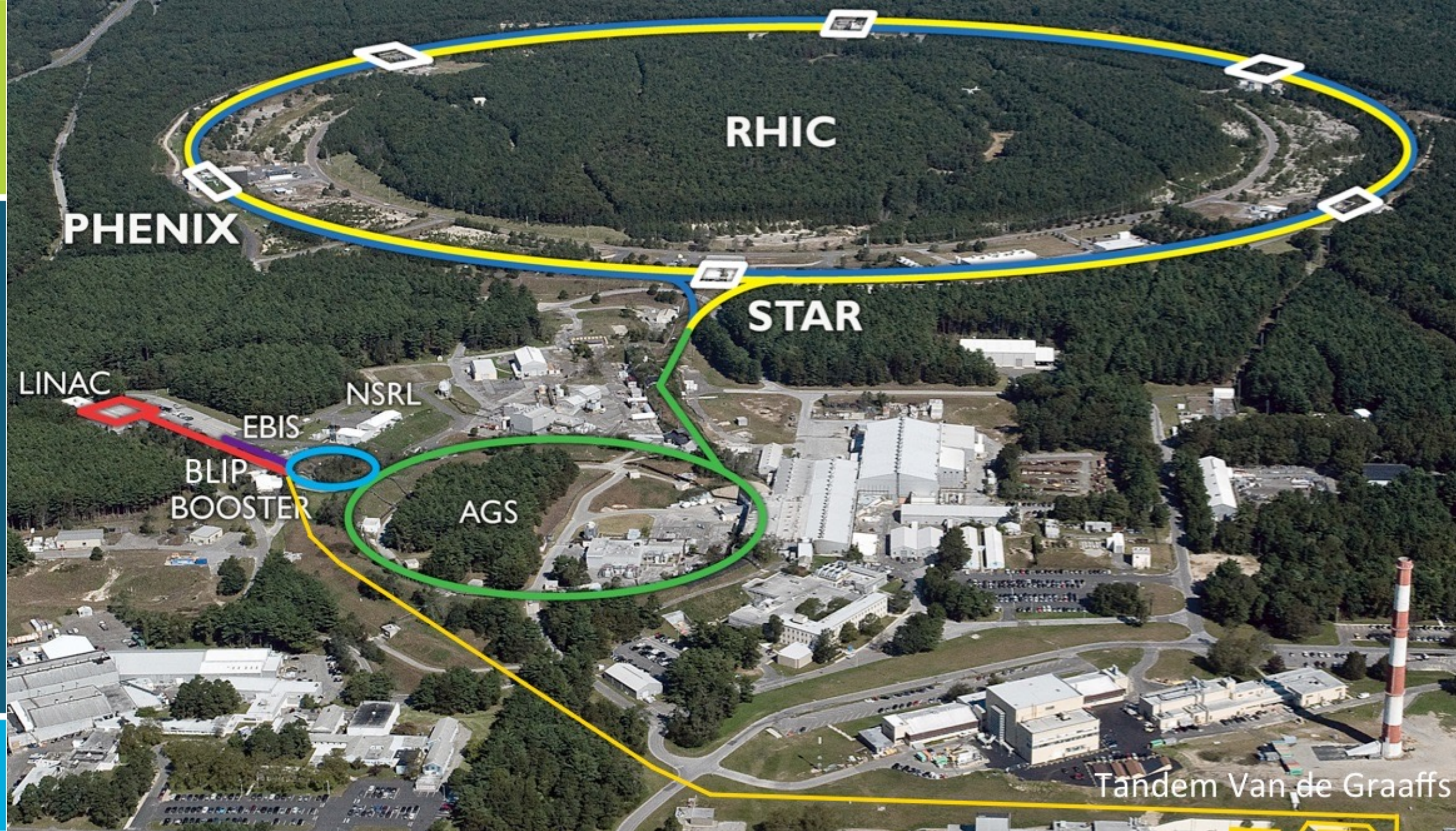Presenter: Yihui "Ray" Ren, AI-CoDesign Group Leader, BNL

@BrookhavenLab

SRO XII, University of Tokyo, Dec. 2-4, 2024

CT

NYC

Long Island

NJ

Atlantic Ocean

2

Relativistic Heavy Ion Collider,
future Electron-Ion Collider
(2.4 miles in circumference)

National Synchrotron
Light Source II

3

Prev: Computational Science Initiative (CSI)
Now: Computing and Data Science (CDS)

# SPHENIX at BNL



Scientific American, 03/01/2023

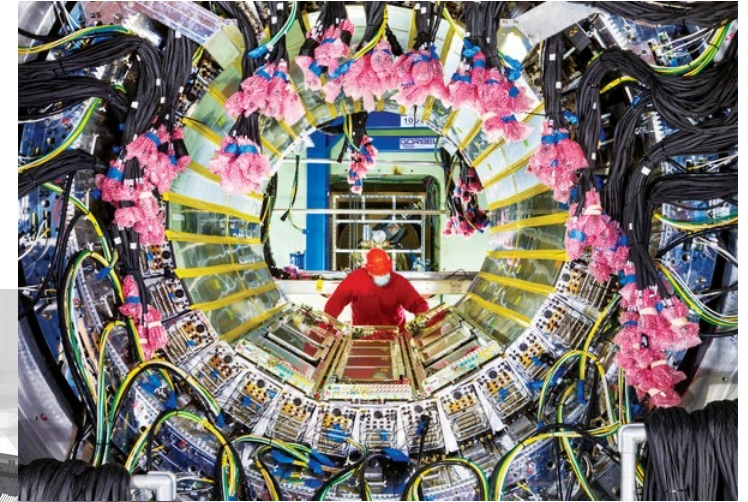**Tracking System**
TPC
INTT
MVTX

**Calorimeters**
Electromagnetic
Inner Hadronic
Outer Hadronic

Data taking began last year!
High-precision tracking system +   Hermetic
Electromagnetic & Hadronic calorimeters

# Time projection Chamber (TPC)

A TPC is composed of 48 layers of rectangular grid of sensor nodes. It acts as a camera capturing 3D particle trajectories.

- three layer groups, 16 layers each
- two sides, divided by the transverse plane passing the collision point
- 12 sectors, 30 degree each

sPHENIX Time Projection Chamber
100 Hz ZDC, MBD Prescale: 2, HV: 4.45 kV GEM, 45 kV CM, X-ing Angle: 2 mrad
2023-06-23, Run 10931 - EBDC03 reference frame 89
Au+Au sqrt(s_{NN})=200 GeV

# TPC Data

**Dataset**: MDC1 AuAu 0-10%C + 170kHz pileup



48 layers, 16 each

249
192
30°

- Number of voxels: (azimuth × z × layer):
  - Outer layer group: 2304 × 498 × 16 ≈ 18M;
  - Middle layer group: 1536 × 498 × 16 ≈ 12M;
  - Inner layer group: 1152 × 498 × 16 ≈ 9M

TPC Drift field

Outer Layer Section:
- 192 × 249 × 16
- Sparser

OUTER HCAL
SC MAGNET
INNER HCAL
EMCAL
TPC
INTT
MAPS
ENDCAP
FLUX RETURN

sPHENIX

Brookhaven National Laboratory

# Neural Auto-Encoder



- A typical Auto-Encoder uses an Encoder network to compress the data into "code"; and a Decoder network to reconstruct the original input.

- The voxel distribution:
  - long-tailed (skewed)
  - sparse (many zero values)
  - zero-suppressed (discontinued)
  - 10-bit integer (saturated)

**Very Challenging for a regular auto-encoder!**

# Bicephalous Convolutional Auto-Encoder (BCAE)

- A dedicated Segmentation decoder to determine whether a voxel has been zero-suppressed.

- Integrate a transformation function $\tau$ into the network:
$\tau = \log(x - 64)/6$ for non-zero values.

input

Segmentation output (mask)

Regression output

Permanent storage

Encoder $E$

$E$ output

$\times m$

code

Seg. Decoder $D_{\text{seg}}$

$D_{\text{seg}}$ output

$\times n$

Reg. Decoder $D_{\text{reg}}$

$D_{\text{reg}}$ output

$\times n$

Regression output × mask

input

# Results

- Compression Ratio:
  **1:27**

- Mean-Squared Error (MSE):
  **218.44**

# Other Lossy Compression Algorithms

We identified three conventional lossy compression algorithms, which were mainly designed for dense data matrices such as in fluid dynamic simulations.

- **MGARD**: MultiGrid adaptive reduction of data.
  https://github.com/CODARcode/MGARD

- **SZ**: Error-bounded lossy compressor.
  https://github.com/szcompressor/SZ

- **ZFP**: Compressor for integer and floating-point data stored in multidimensional arrays.
  https://github.com/LLNL/zfp

# Results

- Conventional methods allow users to change compression ratio.

- Our model has better compression ratio and lower MSE. (good balance)

- Conventional methods do not require "training".



| | Compr. ratio↑ | MSE↓ | log MAE↓ | PSNR↑ |
|---|---|---|---|---|
| MGARD | **27** | 626.28 | 1.213 | 3.223 |
| SZ | 24 | 369.69 | 0.302 | 3.452 |
| ZFP | 19 | 219.48 | 0.267 | 3.678 |
| CAE | 27 | 227.61 | 0.349 | 3.703 |
| BCAEwoT | 27 | 230.59 | 0.193 | 3.706 |
| BCAE | **27** | **218.44** | **0.185** | **3.724** |

**Brookhaven** National Laboratory

# Results

- Conventional methods allow users to change compression ratio.

- Our model has better compression ratio and lower MSE.

- Recover entry value distributions (histogram)



Huang, Y., Ren, Y., Yoo, S., & Huang, J. (2021, December). Efficient data compression for 3d sparse TPC via bicephalous convolutional autoencoder. In *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)* (pp. 1094-1099). IEEE. arxiv:2111.05423

Jin Huang (BNL) "SRO for sPHENIX TPC and Real-time AI" in SRO XI

# 2D Encoder and Decoder

- Smaller, faster encoder
- Bulkier, slower decoder
- Can stronger decoder compensate for a weaker encoder?

# 2D Encoder and Decoder

tunable encoder decoder sizes

# Performance comparison

original BCAE versus BCAE-2D

| model | MAE ↓ | PSNR ↑ | Precision ↑ | Recall ↑ | Encoder size ↓ | Compr. Ratio ↓ |
|---|---|---|---|---|---|---|
| BCAE | 0.198 | 9.923 | 0.878 | 0.861 | 201.7k | 27.041 |
| BCAE-2D | 0.152 | 11.726 | 0.906 | 0.907 | 169.0k | 31.125 |

SRO XII, University of Tokyo, Dec. 2-4, 2024. Presenter: Yihui Ren (BNL)

# Performance comparison
original BCAE versus BCAE-2D

| model | MAE ↓ | PSNR ↑ | Precision ↑ | Recall ↑ | Encoder size ↓ | Compr. Ratio ↓ |
|-------|-------|--------|-------------|----------|----------------|-----------------|
| BCAE | 0.198 | 9.923 | 0.878 | 0.861 | 201.7k | 27.041 |
| BCAE-2D | 0.152 | 11.726 | 0.906 | 0.907 | 169.0k | 31.125 |
| BCAE++ | 0.112 | 14.325 | 0.934 | 0.936 | 226.2k | 31.125 |

**From BCAE to BCAE++**
1. 3D convolution
2. Pad (16, 192, 249) to (16, 192, 256) for easy halving and an increased compression ratio
3. Remove normalization

- Better reconstruction performance
- Still slow

# Performance comparison
original BCAE versus BCAE-2D

| model | MAE ↓ | PSNR ↑ | Precision ↑ | Recall ↑ | Encoder size ↓ | Compr. Ratio ↓ |
|---|---|---|---|---|---|---|
| BCAE | 0.198 | 9.923 | 0.878 | 0.861 | 201.7k | 27.041 |
| BCAE-2D | 0.152 | 11.726 | 0.906 | 0.907 | 169.0k | 31.125 |
| BCAE++ | 0.112 | 14.325 | 0.934 | 0.936 | 226.2k | 31.125 |
| BCAE-HT | 0.138 | 12.376 | 0.916 | 0.915 | 9.8k | 31.125 |

**From BCAE to BCAE-HT**
1. 3D convolution
2. Pad (16, 192, 249) to (16, 192, 256) for easy halving and an increased compression ratio
3. Remove normalization
4. Much smaller intermediate output channels for higher throughput

- Slightly better reconstruction performance
- Super small model size
- Higher throughput

SRO XII, University of Tokyo, Dec. 2-4, 2024. Presenter: Yihui Ren (BNL)

# Performance comparison

original BCAE versus BCAE-2D

SRO XII, University of Tokyo, Dec. 2-4, 2024. Presenter: Yihui Ren (BNL)

# Throughput comparison



- **Full:** encode with float 32, save code as float 16, decode with float 32
- **Half:** encode with float 16, save code as float 16, decode with float 32

| model | mode | MAE | precision | recall |
|---|---|---|---|---|
| BCAE-2D | Full | 0.151937 | 0.905469 | 0.906916 |
| | Half | 0.151965 | 0.905326 | 0.907050 |
| BCAE++ | Full | 0.112347 | 0.933817 | 0.935779 |
| | Half | 0.112342 | 0.933852 | 0.935741 |
| BCAE-HT | Full | 0.138443 | 0.915891 | 0.914562 |
| | Half | 0.138441 | 0.915780 | 0.914701 |

# Throughput comparison



Measured on A6000

SRO XII, University of Tokyo, Dec. 2-4, 2024. Presenter: Yihui Ren (BNL)

BCAE-2D  BCAE-HT  BCAE++

half-precision
full-precision

Higher throughput

ground truth

reconstruction

BCAE-2D  BCAE-HT  BCAE++

difference

BCAE-2D  BCAE-HT  BCAE++

Better recon. performance

Huang, Y., Ren, Y., Yoo, S., & Huang, J. (2023, November). Fast 2D Bicephalous Convolutional Autoencoder for Compressing 3D Time Projection C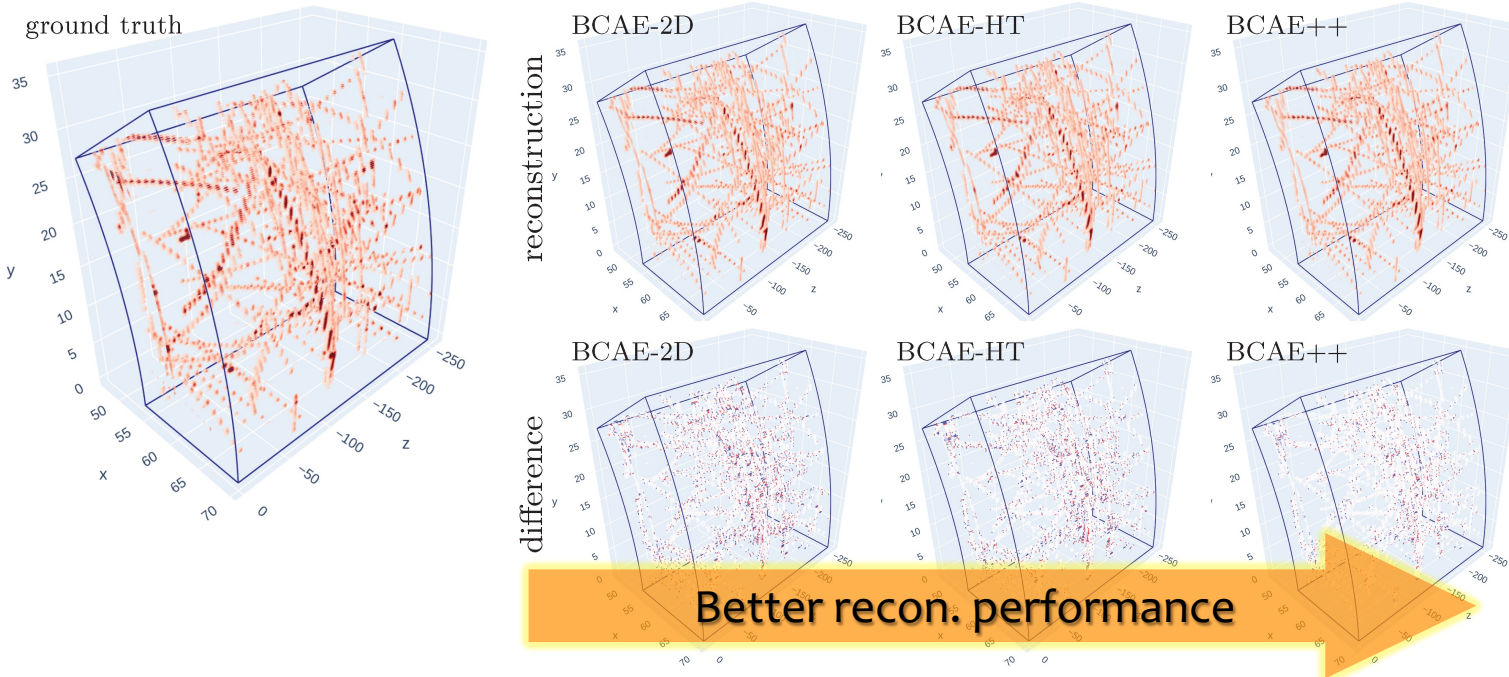hamber Data. In *Proceedings of the SC'23 Workshops of The International Conference on High Performance Computing, Network, Storage, and Analysis* arxiv:2310.15026

# Question 1:
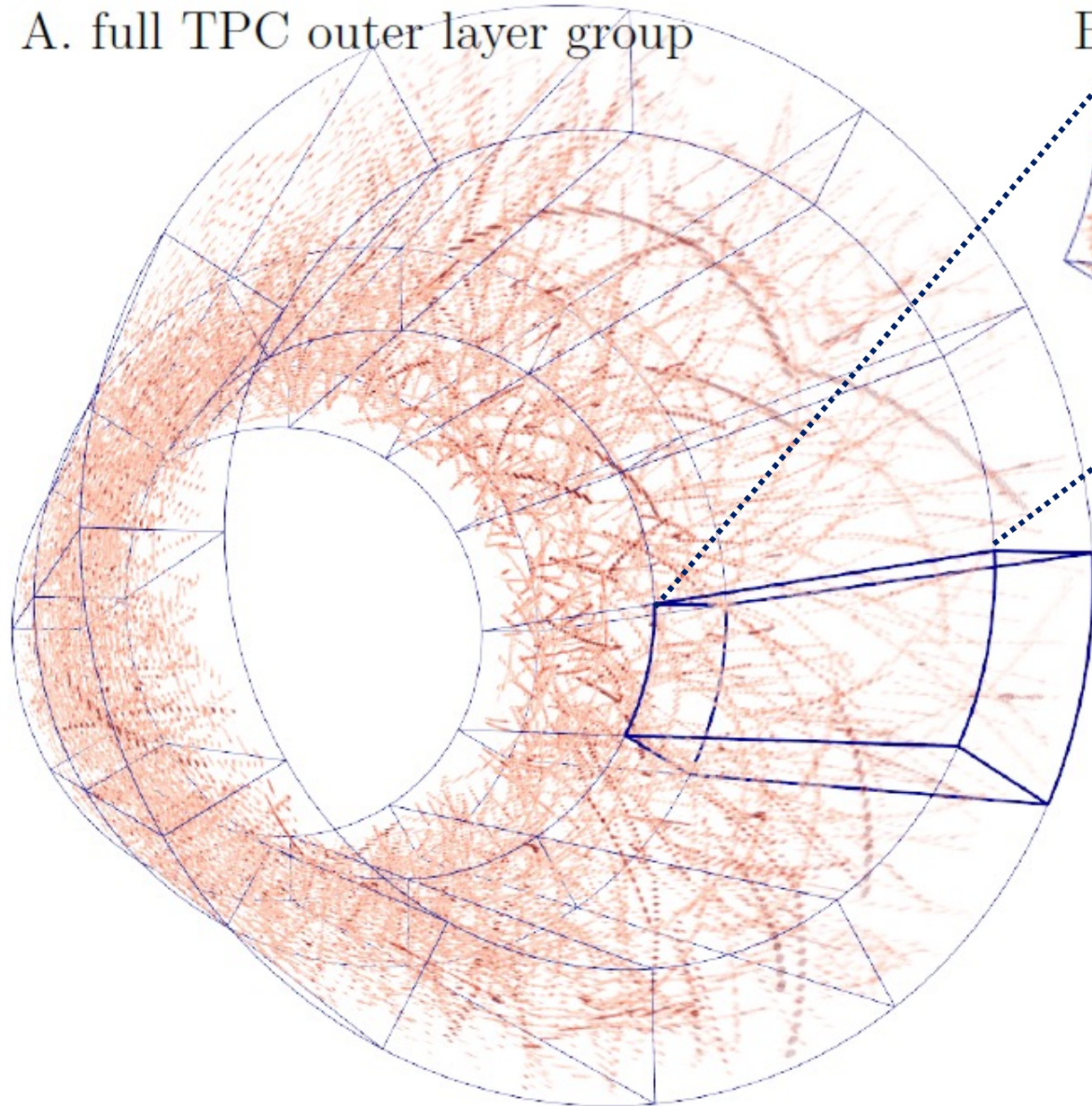
Can we have better performance and better throughput?

# Question 2:

Can we have variable compression ratio depending on occupancy?

# Question 3:

Can we have variable throughput? Sparser the data, the less compute

A. full TPC outer layer group

B. a TPC wedge

C. zoom-in of a trajectory

**Brookhaven**
National Laboratory

C. zoom-in of a trajectory

layer
time
azimuthal

# BCAE-VS: Bicephalous Convolutional Autoencoder with Variable ratio Compression for Sparse input



layer
time
azimuthal

Locate the most valuable signals, and compress by down-selecting the signals



Keypoint Detection

Brookhaven
National Laboratory

# BCAE-VS
## Bicephalous Convolutional Autoencoder with Variable ratio Compression for Sparse input



Implemented by **MinkovskiEngine** from NVIDIA

# Reconstruction Accuracy

# Variable Compression Ratio and Throughput as Function of Occupancy

throughput = number of TPC wedges processed by one GPU per second. GPU we used is NVIDIA 6000 ADA.

fraction of signals saved to the persistent storage (~1/compression_ratio for ZS data)

# Performance comparison

| model | comp. ratio ↑ | reconstruction performance | | | | | efficiency | |
|---|---|---|---|---|---|---|---|---|
| | | $L_1$ ↓ | $L_2$ ↓ | PSNR ↑ | recall ↑ | precision ↑ | encoder size | throughput ↑ |
| BCAE-2D | 31 | .152 | .862 | 20.6 | .907 | .906 | 169k | **9.6k** |
| BCAE-HT (3D) | 31 | .138 | .781 | 20.8 | .916 | .915 | 9.8k | **9.6k** |
| BCAE++ (3D) | 31 | .112 | .617 | 21.4 | .936 | .934 | 226k | 3.2k |
| BCAE-VS | **34** | **.028** | **.089** | **26.0** | **.988** | **.996** | **382** | 5.6k |

Huang, Y., Go, Y., Huang, J., Li, S., Luo, X., Marshall, T., ... & Yoon, B. J. (2024). Variable Rate Neural Compression for Sparse Detector Data. *arXiv preprint arXiv:2411.11942.*

**Brookhaven** National Laboratory

# Summary

- Bicephalous Convolutional Auto-encoder (BCAE) for sparse TPC data
- Faster BCAE-2D with Encoder-Decoder tradeoff
- Variable compression rate and computation with BCAE-VS.
- (Future) Improve BCAE-VS in low occupancy region
- (Future) Improve throughput of BCAE-VS on GPU and other hardware
- (Future) noise rejection, tracking efficiency, etc.

Source Code:
- BCAE https://github.com/BNL-DAQ-LDRD/NeuralCompression
- BCAE-2D https://github.com/BNL-DAQ-LDRD/NeuralCompression_v2
- BCAE-VS https://github.com/BNL-DAQ-LDRD/NeuralCompression_v3

# Summary

- Bicephalous Convolutional Auto-encoder (BCAE) for sparse TPC data
- Faster BCAE-2D with Encoder-Decoder tradeoff
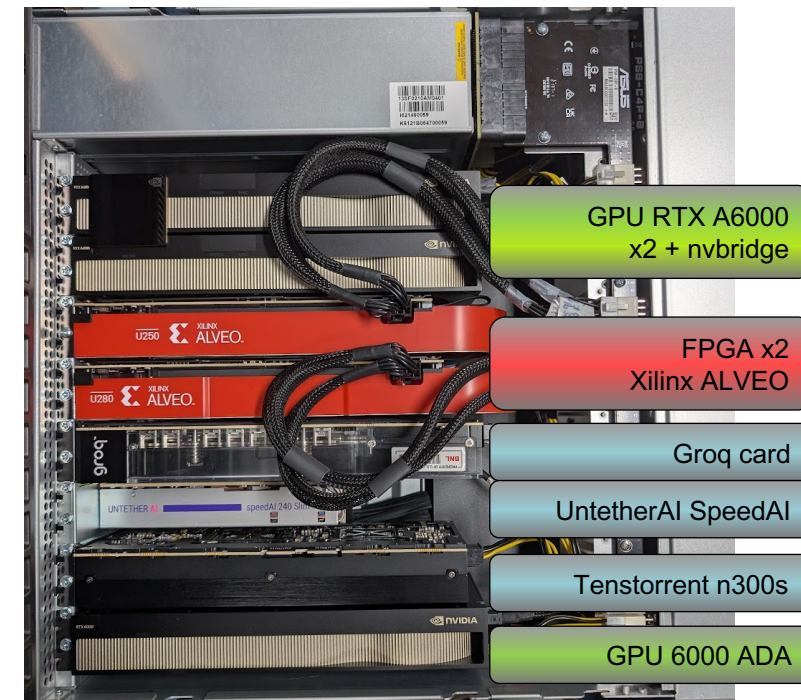- Variable compression rate and computation with BCAE-VS.
- (Future) Improve BCAE-VS in low occupancy region
- (Future) Improve throughput of BCAE-VS on GPU and other hardware
- (Future) noise rejection, tracking efficiency, etc.

**Feel free to try on your streaming detector data!** Source Code:
- BCAE https://github.com/BNL-DAQ-LDRD/NeuralCompression
- BCAE-2D https://github.com/BNL-DAQ-LDRD/NeuralCompression_v2
- BCAE-VS https://github.com/BNL-DAQ-LDRD/NeuralCompression_v3



GPU RTX A6000 x2 + nvbridge
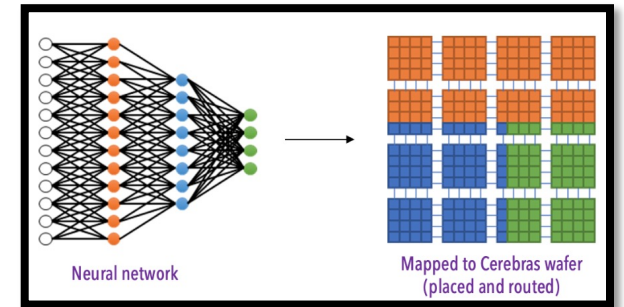
FPGA x2 Xilinx ALVEO

Groq card

UntetherAI SpeedAI

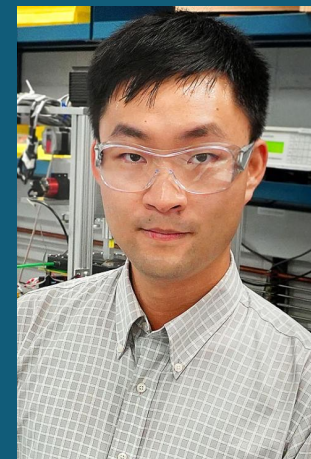Tenstorrent n300s

GPU 6000 ADA
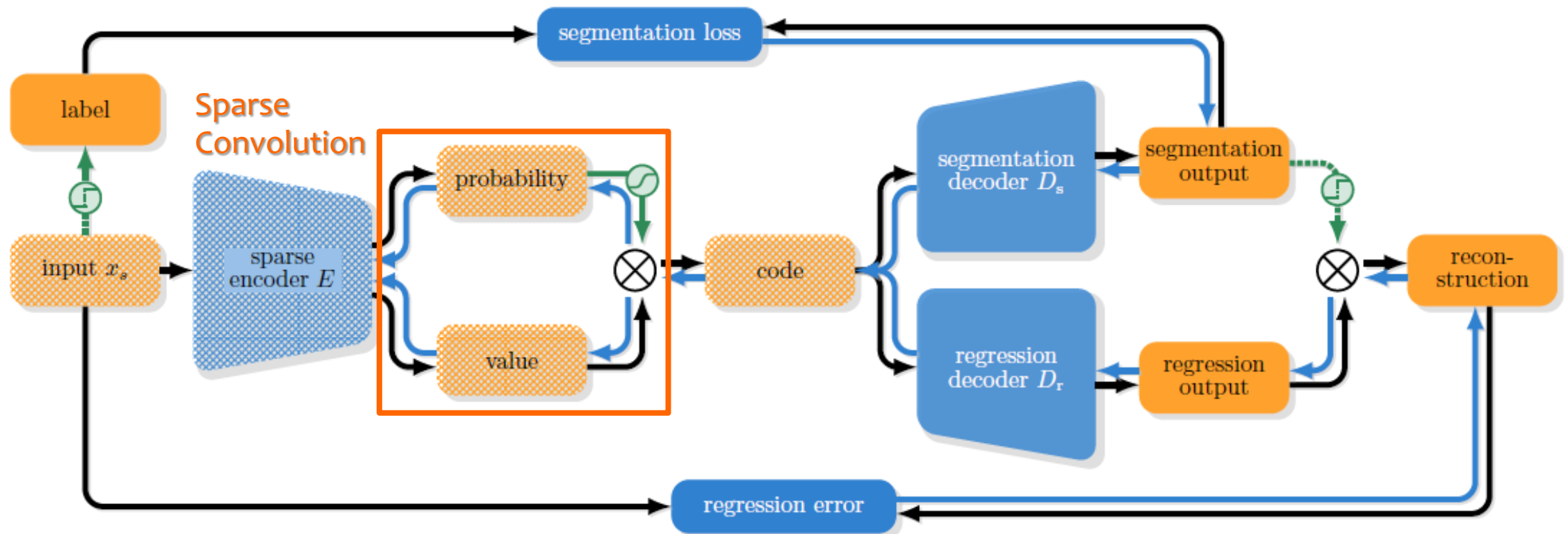
# Thank you!

# ありがとうございます

Yihui "Ray" Ren (yren@bnl.gov)

**Yi Huang**, Yihui "Ray" Ren, Yeonju Go, Xihaier Luo, Shuhang Li, Thomas Marshall, Joseph D. Osborn, Christopher Pinkenburg, Evgeny Shulga, Shinjae Yoo, Byung-Jun Yoon, **Jin Huang** (PI)

# BCAE-VS

Variable ratio Compression for Sparse input

# Sanity Check with lossless compressions

- Compare to lossless compression tools. More than twice compression ratio.

- Compressed code of BCAE can not be compressed further.

**How is our model comparing to other lossy compression algorithms?**



Zip Ratios of Raw

| | |
|---|---|
| gzip mean 10.10 | |
| bzip mean 16.20 | |
| 7z mean 12.44 | |
| BCAE ratio 26.98 | |



Zip Ratios of BCAE-compressed

| | |
|---|---|
| gzip mean 1.08 | |
| bzip mean 1.09 | |
| 7z mean 1.11 | |

Brookhaven National Laboratory

# Real-time AI accelerator

Gpuserver0 upgrade yesterday

L40S removed for RMA

Three AI chips:
- Groq: TPU based
- UntetherAI:  INT8 optimizes
- Tenstorrent: RISC-V cores

More photos:
https://photos.app.goo.gl/z3An
Nhfrd4bqZTeN6

GPU RTX A6000
x2 + nvbridge

FPGA x2
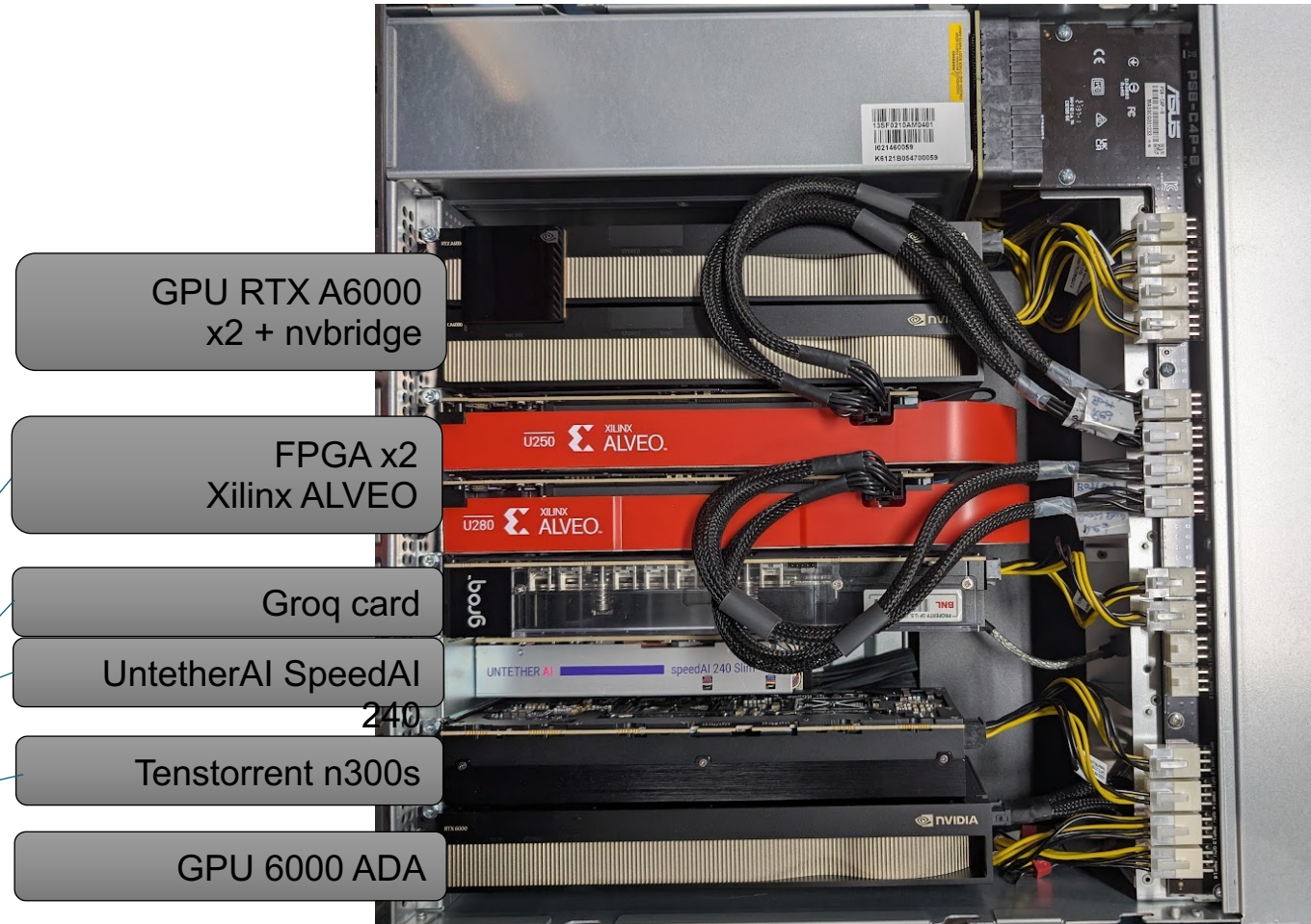Xilinx ALVEO

Groq card

UntetherAI SpeedAI
240

Tenstorrent n300s

GPU 6000 ADA

```
jinhuang@gpuserver0:~$ sudo lspci | grep cc
41:00.0 Processing accelerators: Xilinx Corporation Device d004
61:00.0 Processing accelerators: Xilinx Corporation Device d00c
81:00.0 Processing accelerators: Device 1e67:0004 (rev 01)
a1:00.0 Processing accelerators: Groq TSP100 Tensor Streaming Processor
c1:00.0 Processing accelerators: Device 1e52:401e (rev ff)
```

# Summary

- Bicephalous Convolutional Auto-encoder (BCAE) for sparse TPC data
- Faster BCAE-2D with Encoder-Decoder tradeoff
- Variable compression rate and computation with BCAE-VS.
- (Future) Improve BCAE-VS in low occupancy region
- (Future) Improve throughput of BCAE-VS on GPU and other hardware
- (Future) noise rejection, tracking efficiency, etc.

Source Code:
- BCAE https://github.com/BNL-DAQ-LDRD/NeuralCompression
- BCAE-2D https://github.com/BNL-DAQ-LDRD/NeuralCompression_v2
- BCAE-VS https://github.com/BNL-DAQ-LDRD/NeuralCompression_v3



Neural network → Mapped to Cerebras wafer (placed and routed)

Brookhaven National Laboratory