# How Much can we Really Compress Scientific Data?

Franck Cappello

Argonne National Laboratory

University of Illinois at Urbana Champaign


Arkaprabha Ganguli (Michigan State University),
David Krasowska(Clemson),
Julie Bessac (ANL), Sheng Di (ANL), Robert Underwood (ANL),
Xiaodong Yu (ANL,
Jon C. Calhoun (Clemson)

# Progress

## Scientific Data Compression:
## From Stone-Age to
## Renaissance

- Background
- Focus on spatial compression
- Best in class lossy compression
- Open questions

This is what we need to compress
(bit map of 128 floating point numbers):

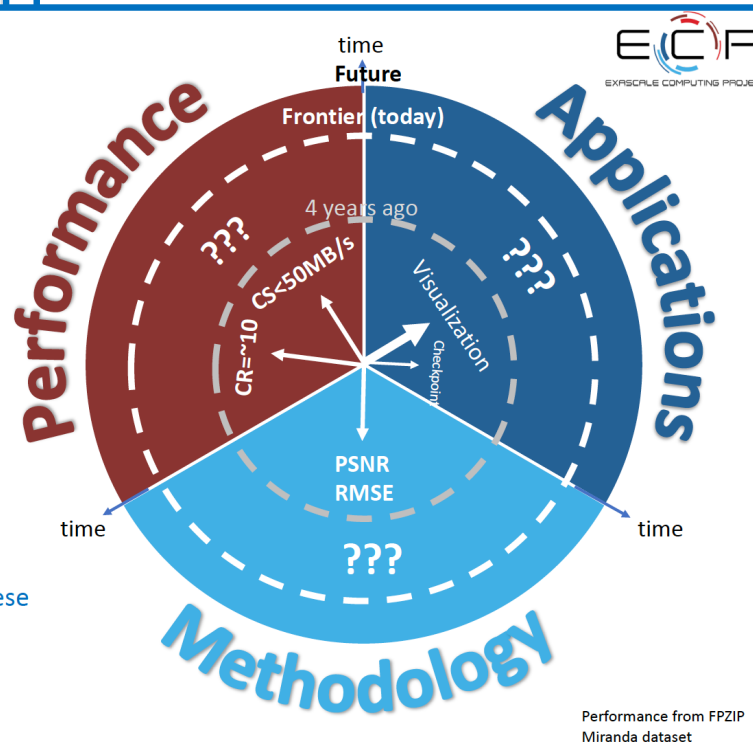## Three Frontiers of Lossy Compression of Scientific Data

## Three Frontiers

**Past (still used today)**
-Lossless compression
-Decimation in space and time + linear or tricubic (or other) interpolation
-Lossy compression mainly for Visualization

Technologies that are at the current frontiers: SZ, ZFP, Z-checker+SDRBench

In this talk, I will illustrate the progresses in these three frontiers based on our explorations and results in the **ECP EZ, CODAR and Exasky Projects and ANL SZ compressor**



Performance from FPZIP
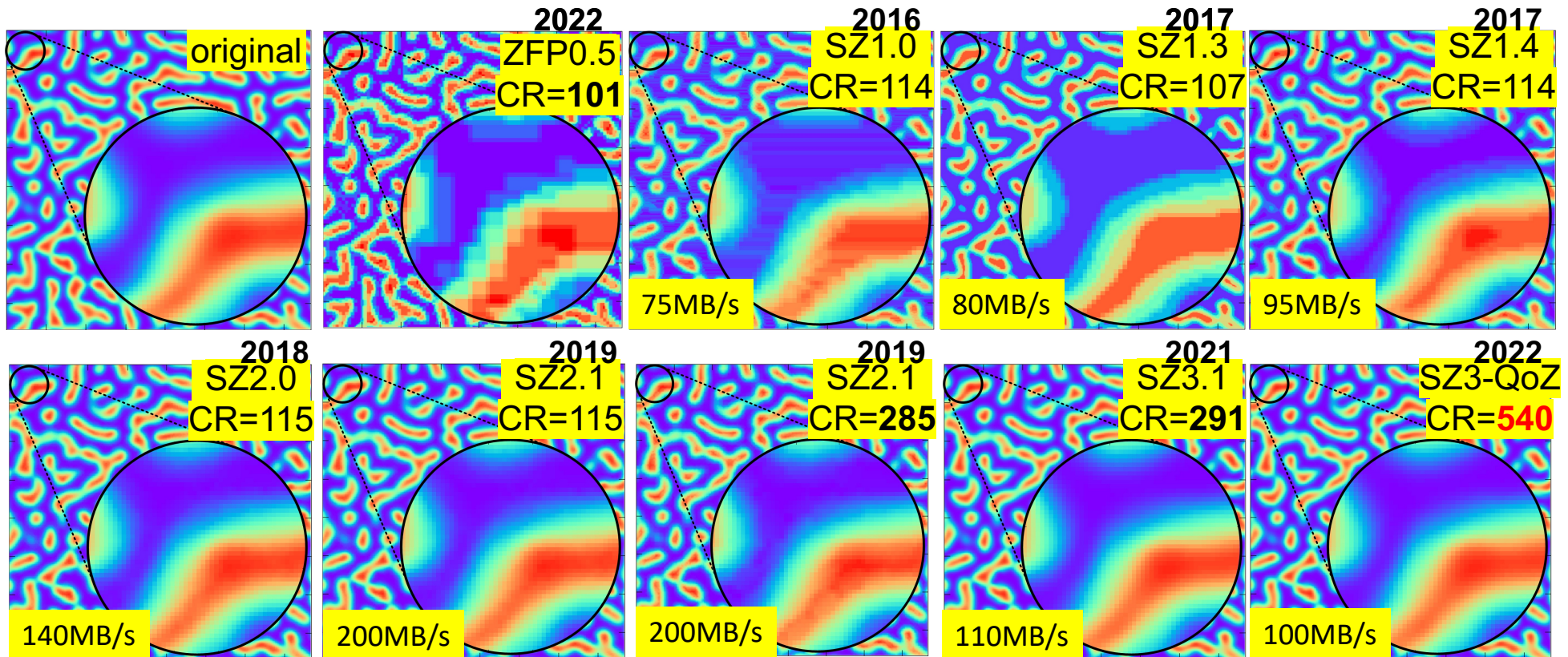Miranda dataset

# Lossy compression of scientific data

- Consist in reducing scientific data volume by leveraging correlations and reducing precision (lossless compression does not reduce scientific data enough)

- Compression ratios (with current compressors) vary depending on use-cases, typically:
  - CR=5 for hard to compress dataset and demanding data/analysis quality preservation
  - CR=10-100 for scientific data presenting high correlation and medium data/analysis quality preservation
  - CR=x100 for visualization (low data/analysis quality preservation)

- Goal: keep the same science (satisfy user's quality requirements WRT QoIs – features)
  - WARNING: You will see images because this is the easiest way to show distortion but compression of scientific data is NOT for images

- Getting significant traction in the scientific community (climate, cosmology, seismic, etc.), IoT community as well (sensors, EKG)

Argonne
NATIONAL LABORATORY

# Progress in Compression Techniques

# Huge Progress in performance in the past 5-6 years

**ECP** EXASCALE COMPUTING PROJECT

Evolution of SZ compression quality and performance using a large-eddy simulation of multicomponent flows with turbulent mixing: Miranda - density field.



Row 1:
- original
- **2022** ZFP0.5 CR=**101**
- **2016** SZ1.0 CR=114 75MB/s
- **2017** SZ1.3 CR=107 80MB/s
- **2017** SZ1.4 CR=114 95MB/s

Row 2:
- **2018** SZ2.0 CR=115 140MB/s
- **2019** SZ2.1 CR=115 200MB/s
- **2019** SZ2.1 CR=**285** 200MB/s
- **2021** SZ3.1 CR=**291** 110MB/s
- **2022** SZ3-QoZ CR=**540** 100MB/s

Visualization of Miranda - density data for SZ's different versions (EB: VRAE 1E-2), Performance on single core CPU (Intel Broadwell)

**SZx compresses at 300GB/s on NVIDIA A100 → Bottleneck is not compression but PCIe**

# More Lossy Compressors

**ZFP** (LLNL): Transform (DCT)
ECP ZFP
Overpreserves data, lower
Compression ratio compared
to SZ, Better speed.

**SPERR** (NCAR): Wavelet
Works well on wave
propagation problem
(Climate, Seismic)

**MGARD** (ORNL)
ECP CODAR
Multigrid adaptive
reduction
MGARD controls the
compression errors in
quantities of interest ($Q$):
Linear expression of the error

Largest Compression Ratio For Each Compressor that
Satisfies  Each Pinard et al (2020) Requirements

| Compressor | Pearson R^2 | Spatial Error | KS-test |
|---|---|---|---|
| SZ_Interp | 93 | 93 | **21** |
| SZ (regression) | 14.34 | 14.34 | **14.34** |
| ZFP | 5.45 | 5.45 | **2.36** |
| MGARD | 27.1 | 4.69 | X |
| MGARDx | 14.7 | 6.49 | X |
| TThresh | 16.1 | 16.1 | **2.98** |
| BitGrooming | 1.51 | 1.51 | **1.51** |
| Digit Rounding | 1.86 | 1.86 | **1.86** |
| FPZip | 1.95 | 1.95 | **1.95** |
| NDZip | 1.64 | 1.64 | **1.64** |
| Zstd | 1.35 | 1.35 | |

# More Lossy Compressors

**TTRESH** (LLNL):

HoSVD (Tucker Decomposition)
Quantize the Core tensor
Very high compression ratio
Tendency to blur the overall data
(loose details)
1 or 2 orders of magnitude slower
than SZ or ZFP

**Autoencoders**

Overall architecture of
convolutional autoencoder
(A. Glaws, R. King, and M.
Sprague, "Deep learning for
in situ data compression of
large turbulent flow
simulations," Physical Review
Fluids, vol. 5, no. 11, p.
114602, 2020.)

12 residual blocks
for feature extraction
+ 3 compression layers

Largest Compression Ratio For Each Compressor that
Satisfies Each Pinard et al (2020) Requirements

| Compressor | Pearson R^2 | Spatial Error | KS-test |
|---|---|---|---|
| SZ_Interp | 93 | 93 | **21** |
| SZ (regression) | 14.34 | 14.34 | **14.34** |
| ZFP | 5.45 | 5.45 | **2.36** |
| MGARD | 27.1 | 4.69 | X |
| MGARDx | 14.7 | 6.49 | X |
| TThresh | 16.1 | 16.1 | **2.98** |
| BitGrooming | 1.51 | 1.51 | **1.51** |
| Digit Rounding | 1.86 | 1.86 | **1.86** |
| FPZip | 1.95 | 1.95 | **1.95** |
| NDZip | 1.64 | 1.64 | **1.64** |
| Zstd | 1.35 | 1.35 | **1.35** |

Significant
Smoothing

Argonne
NATIONAL LABORATORY

# What makes SZ3 different: a Highly **Modular/ Customizable** Compression Framework

**ECP**
EXASCALE COMPUTING PROJECT

2021
**R&D 100**
WINNER

**SZ 3 (C++)** library of algorithms for lossy compression and examples of SZ compressors built from the library of algorithms.

To compose and tune a compression pipeline we analyze the data to compress and user requirements in compression speed, ratio and accuracy.



User requirements

**Lossy compressor parameters**

Data Analysis

**Lossy compressor composition**

SZ Library of compression algorithms

Feature detection
Wavelet Transform
Resolution coarsening
Linear regression predictor
Lorenzo Predictor
Spline interpolation predictor
Pattern based predictor
Auto-encoder predictor
Linear quantization
Log transform
Huffman coding
Arithmetic coding
Leading bit coding
Truncation
Zstd

LCLS
SZ-separator
SZ 2.1 (default)
SZ-Pattern
SZ-Interp
SZx

**Argonne**
NATIONAL LABORATORY

# Progress in Applications and Methodology

# Many more Applications (than in 2018)

- Climate
- Combustion
- Cosmology
- Deep Learning
  - Activation data
  - Model coefficients
  - Training data
- Extreme Weather
- Fusion Energy
- Hydrodynamics
- IoT
- Light Sources (Physics Instruments)
- Materials Science
- Molecular Dynamics
- Quantum Chemistry
- Quantum Circuit Simulation
- Seismic Imaging

# Many more Use-cases (than in 2018)

We are seeing an increasing diversity/number of use-cases

"Classic" use-cases:

1) Visualization

2) Reducing storage footprint (offline compression)

3) Reducing I/O time (on-line, in-situ compression)

Recently identified use-cases:

4) **Reducing streaming intensity** (recent for generic floating-point compressors)

5) Lossy checkpoint/restart from lossy state

- reduce checkpoints footprint on storage – adjoint, accelerate checkpointing

6) Re-computation Avoiding by reducing the memory footprint → GAMESS

7) **Running larger simulations by reducing the memory footprint**

8) Accelerating CPU/GPU – memory transfer

9) Reduce DNN model size

10) Accelerate training (I/O read time) of DNNs

Cappello, F., Di, S., Li, S., Liang, X., Gok, A. M., Tao, et Al., Use cases of lossy compression for floating-point data in scientific data sets. *The International Journal of High Performance Computing Applications*, *33*(6), 1201–1220, 2019

# Huge Progress in Methodologies

https://sdrbench.github.io/

https://github.com/robertu94/libpressio

https://github.com/CODARcode/Z-checker

# Putting all Together
# example: LCLS/Crystallography

# Example of Success Story: Crystallography

With Chuck Yoon: Stanford

**1: X-ray Beam**

Detector

Liquid Jet

KB Mirrors

Primary Interaction Point

Undulator (420 m upstream)

Diffraction before destruction
Number of pulses/sec: 120
Millions of diffraction patterns from crystals

**2: Diffraction**

**3: Reduction**

LCLS-II Data System

DAQ (per hutch)

Data Reduction Pipeline (shared - 1 for NEH, partitionable)

FFB (shared - 1 for NEH)

Offline (shared by all)

1 PB

100 PB

Data compression

Event level veto

GPU

Ethernet

IB EDR

Ethernet 10 Gb/s

Timing (input)

PM

Online Monitoring Nodes

Event Builder Nodes

Fast Feedback (SSD)

DAQ Readout Nodes

DRP Nodes

Control

Fast Feedback Analysis Farm

Offline Analysis Farm

Lustre

1 MHz acquisition 250 GB/s

Factor of 10 reduction

25 GB/s Data written in HDF5 format

Context: LCLS II. Goal: Definition of reduction method
Detector produces:
- 2D images @ 250GB/s
- 4M pixel/event unsigned integers, in binary XTC2 format

Compression objectives: **CR of 10 or more** with error bound @ 500 MB/s/core
→ **RoiBinSZ** algorithm (regions of interest extraction + background binning + SZ background compression)

**RoiBinSZ compression pipeline**

Data from detector

Peak finder

Mask regions of interest

Extract background

Low pass filter (binning e.g. 2x2)

Compress background (SZ 2.1)

Rebuild image

Indexing

Merge

Phase

Refine

Image

ECP

Argonne
NATIONAL LABORATORY

# First Level of Analysis Distortion: Indexing

## Roibin SZ on Se-SAD SFX Dataset (Selenium)

selenobiotinyl-streptavidin on a cspad detector

<span style="color:red">Chuck Yoon: Stanford</span>

|  | Original | Riobin SZ |
|---|---|---|
| **Total compression ratio** | 1 | **70.65** |
| **Number of hits** | 744,150 | 744,150 |
| **Number indexed** | 255,065 | 255,918 |
| **Rsplit ↓** | 7.58% | **7.08%** |
| **CC1/2 ↑** | 0.997 | **0.997** |
| **CCano ↑** | 0.087 | **0.104** |
| **Rwork ↓** | 0.206 | **0.199** |
| **Rfree ↓** | 0.231 | **0.223** |
| **Map-model CC ↑** | **0.81** | 0.8 |

↑: higher the better          ↓: lower the better

- Number of hits: An image with at least 15 peaks is considered a hit
- Number indexed: Number of crystals extracted from hits
- Rsplit: measure precision of averaged intensities/amplitudes
- CCano: The correlation coefficient of the Bijvoet differences of acentric reflections
- CC1/2: Pearson correlation coefficient.
- Rwork: measure of the agreement between the crystallographic model and the experimental X-ray diffraction data
- Rfree: Rwork computed on a small, random sample of data
- Map-model CC: cross-correlation between electron density map and model.
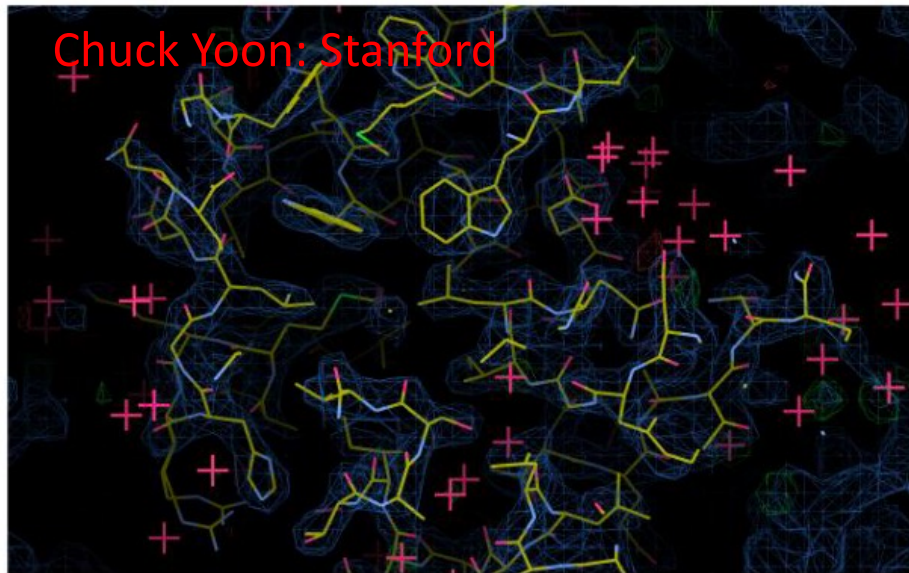
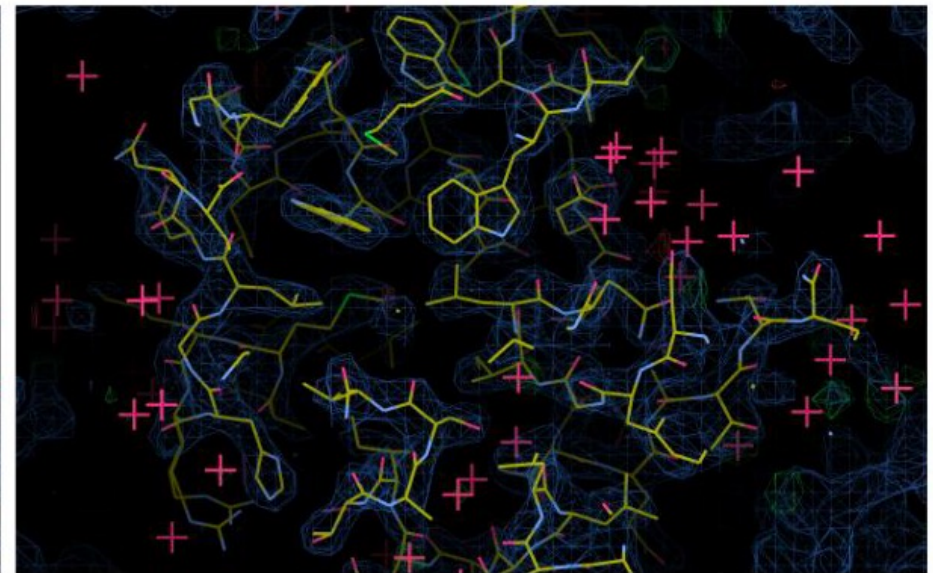# Final Level of Analysis Distortion: Protein Reconstruction

Lysozyme on a jungfrau4m detector

## Reconstruction of Electron Densities Lysozyme

**Very important role in our immune system: breaks up (digests) components of the cell walls of bacteria.**



(a) original          (b) roibin-sz

The data on the right is 196x smaller (or 631× if also using Non-Hit Rejection)

# What's Next

- ## Compressing Memory/ Communications
  - Panda's paper at IPDPS21: MVAPICH+ZFP. Collaboration with UTK on SZ for "mix precision" computation. Difficult problem: how for formulate the impact of lossy compression error on execution results?

- ## Feature Preservation
  - Preservation of derivatives, Structures (Blobs, Halos, Critical points, etc.), Reduction of artifacts, coupling with feature detections. Difficult problem: how to formulate/design error control from quality requirement on features

- ## Automation of Compression Pipeline Construction
  - Many possibilities of compression stages association (pre-processing, decorrelation, quantization, encoding). Difficult problem: How to automatically identify the best pipeline responding to user defined constraints in compression ratio, accuracy, speed

- ## Compressibility Bounds

Argonne
NATIONAL LABORATORY

# Lossy Compressibility Bounds

Estimate "absolute" compression bounds →
compression ratio that no compressor can exceed given a user defined
quality constraints: e.g. local max absolute/relative error
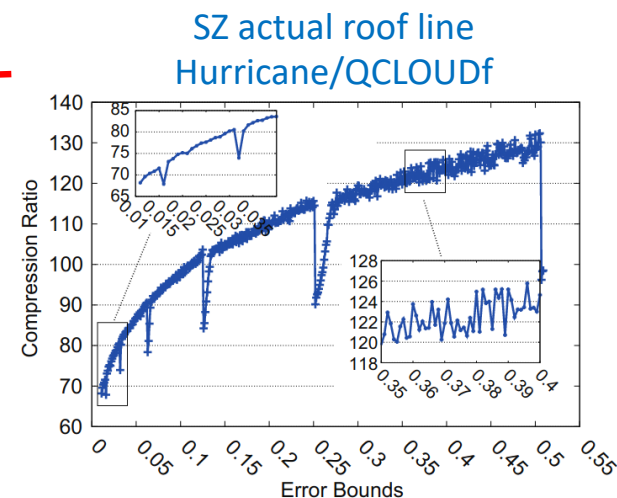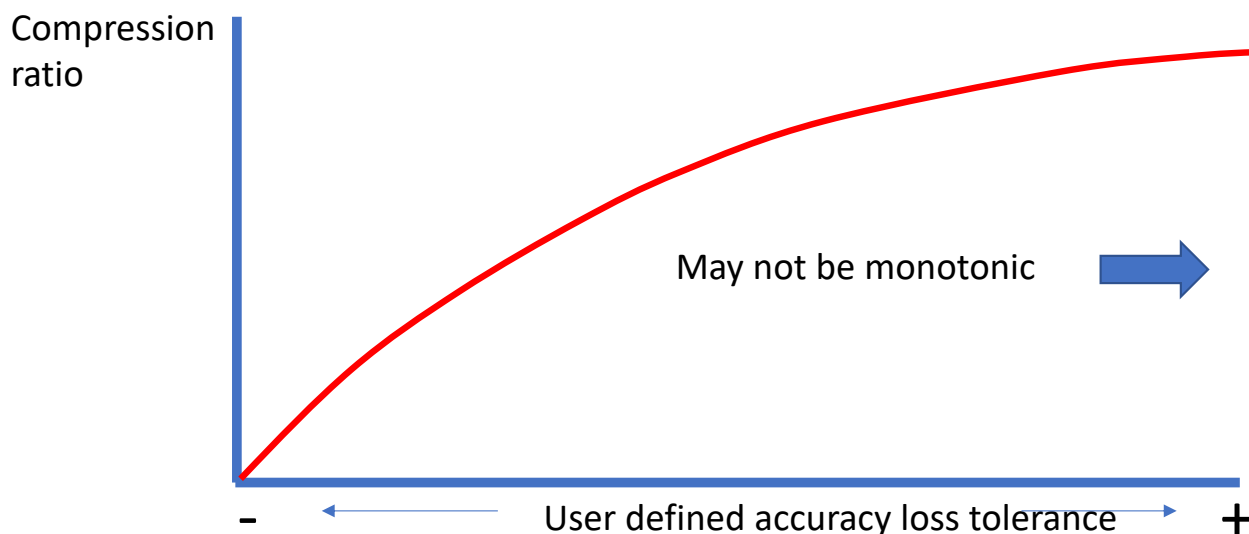
Why does this matter:
- We cannot know if current compressors are really good or not

Important: Scientific datasets in general can be considered as:
Non-Gaussian (distribution), with Memory (correlations), non-Stationary,
non-Ergodic (do not visit all possible state) random processes

# Objective → Roofline of lossy compressibility
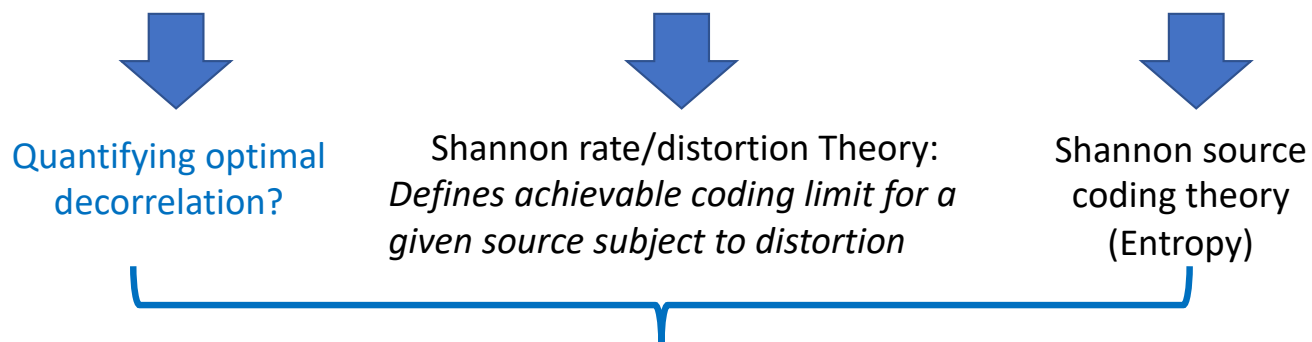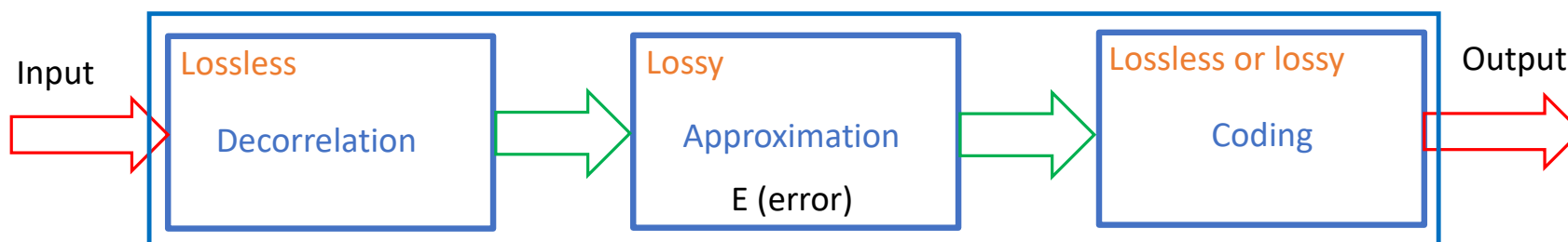
- Lossy compressibility depends on user defined accuracy loss tolerance
- We want a compressibility estimation that capture correlations



SZ actual roof line Hurricane/QCLOUDf

- Should be specific to each dataset and user defined accuracy loss tolerance criteria
- Should be independent of existing compressors

# Lossy Compression Pipeline

## Typical design of a lossy compressor for scientific data



Input → **Lossless** | **Decorrelation** → **Lossy** | **Approximation** | E (error) → **Lossless or lossy** | **Coding** → Output

Quantifying optimal decorrelation?

Shannon rate/distortion Theory: *Defines achievable coding limit for a given source subject to distortion*

Shannon source coding theory (Entropy)

Formulation of Compressibility Bound

Argonne
NATIONAL LABORATORY

# Bound on Lossy Compressibility (Absolute: even if not practical)?

Many researches/results in information theory extending the Shannon rate/distortion theory.
→ But assumptions do not correspond to Scientific data compression: e.g. Gaussian source with memory or Non-gaussian memoryless, non-stationary Gaussian, etc.
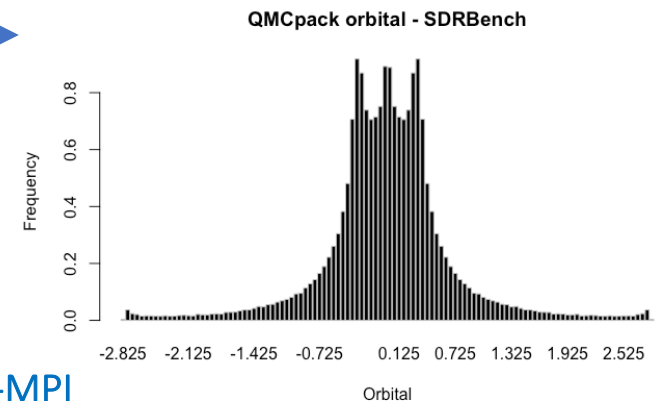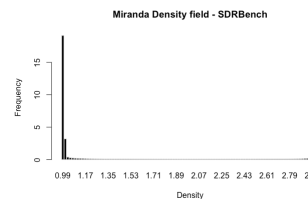
Let's look specifically at the decorrelation stage:

Both produce optimal decorrelation WRT coding gain, in Gaussian case. However, assumptions and data overhead limit their applicability.

KLT (Karhunen–Loève transform) is know to be optimal WRT for decorrelation efficiency (given a bit budget in the compressed format, KTL minimizes the distortion)

But KLT is not optimal (transform) for scientific datasets
→ Because the data distribution is not Gaussian
→ Only for transform based compressors
→ Data needs to be homoscedastic (same variance everywhere on the field)
→ Intrinsically needs a notion of blocks to compute the covariance

QMCpack orbital - SDRBench

Miranda Density field - SDRBench

SVD (and HoSVD: Tucker decomposition): TTRESH and Tucker-MPI
→ Not recommended for Images (DCT performs better), but works well for Vis. ≥3D
→ There is not yet proper consensus on a mathematical framework to select the core tensor elements for optimal truncation
→ Needs to identify an elimination strategy on coefficient to coefficient basis (truncation strategy)

Argonne
NATIONAL LABORATORY

# Lossy Compressibility Bounds

A less difficult problem:

Estimate compression/quality for several compressors. →
Compression bound WRT existing compressors

Why does this matter:
- **Enable fast, automated configuration of a single compressor** to have max quality that will fit in available storage [7]. Cosmology simulations [8], climate simulations [9] and X-ray crystallography [10]
- **Enable quickly choosing between several compressors** with highest CR at runtime in order to minimize data size.
- **Accurately pre-allocate memory** when using compression to expand the amount of on-node logical memory to run applications that utilize in-line compression [11] and quantum chemistry simulations [12].
- **Accurately foreseeing the data transfer time** of I/O or on networks across different devices or sites when compression is used to optimize resource utilization across network links when streaming data [13].

Argonne
NATIONAL LABORATORY

# Bound on Lossy Compressibility (compressor)?

Existing prediction models use knowledge of compressor designs (white box).
→ Not accurate enough (prediction err can be >100%)

- **Formulation of a generic (compressor free) statistical prediction model.**
- **Use training from observed compression ratios (black box) for its specialization to a compressor.**

2 steps:
Step 1: Use Statistical predictors on dataset:
a) SVD to exploit spatial correlation, b) standard deviation to account for variability, c) quantized entropy to represent lossyness and codding
Step 2: Train compression models from regressions (linear and spline-based) to fit predictors to actual observations
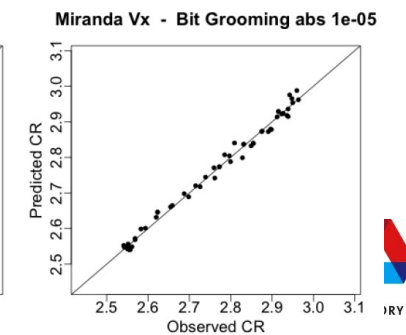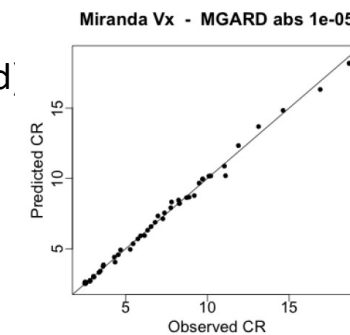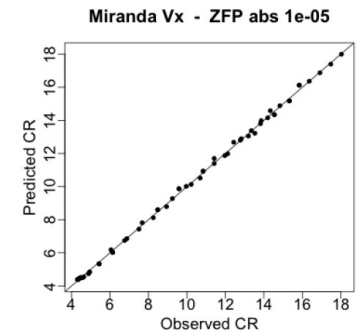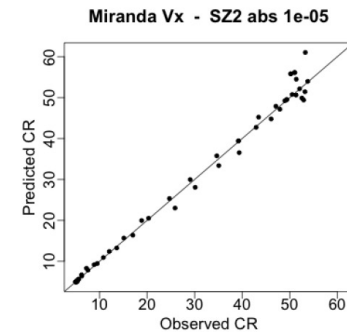
Can be applied and specialized to any lossy compressor before compression

*Uses log because values may appear from dividing by the standard deviation*

Linear model:

$$\overset{*}{\log}(CR) = a + b \times \log(\text{q-ent}) + c \times \log\left(\frac{\text{SVD-trunc}}{\sigma}\right)$$
$$+ d \times \log(\text{q-ent}) \times \log\left(\frac{\text{SVD-trunc}}{\sigma}\right) + \epsilon, \quad (1)$$

where $\epsilon$ is a Gaussian random variable with mean $0$ and standard deviation $\sigma_{eps}$. Coefficients $a$, $b$, $c$, $d$ and $\sigma_{eps}$ are estimated by least-square estimation with the R-function lm.

Tested on 4 different Compressors:
(out-of-sample prediction)

- SZ (prediction)
- ZFP (Transform)
- MGARD (Multi-grid)
- Bit Grooming (truncation and bit operations)

# Conclusion

Lossy Compression for scientific data:

- A very active research topic
- Many teams working on the topic
- Excellent progress in the past 5-6 years
- Significant interest/adoption by apps
- Still many interesting open questions

# Thanks

**Argonne**
NATIONAL LABORATORY



RESEARCH SPONSORED BY

The Exascale Computing Project

A Collaborative effort of the U.S. Department of Energy Office of Science And the National Nuclear Security Administration

(17-SC-20-SC)

**U.S. DEPARTMENT OF ENERGY** | Office of Science

**NNSA** National Nuclear Security Administration