# Unbiased Reconstruction of Lossy-compressed Statistical Data

Boram Yoon (CCS-7)

ISTI Year-End Review
September 13, 2021

LA-UR-21-28904

# Lossy Data Compression using Machine Learning

- Original data: $\left\{ \boldsymbol{X}^{(k)} | \boldsymbol{X}^{(k)} \in \mathbb{R}^D, k = 1,2,3, \dots, N \right\}$

- Compressed data: $\left( \left\{ \boldsymbol{a}^{(k)} | \boldsymbol{a}^{(k)} \in \mathbb{R}^L, k = 1,2,3, \dots, N \right\}, \Phi \right)$
  - $\Phi$: Dictionary containing the information how to reconstruct $\boldsymbol{X}$ from $\boldsymbol{a}$; common for all $k$
  - $\boldsymbol{a}^{(k)}$: Vector in compressed (reduced dimension) space $(L \ll D)$

- Example: principal component analysis (PCA)
  - PCA finds orthogonal directions that maximize the variance as principal components
  - Compression by saving only the coefficients $(\boldsymbol{a}^{(k)})$ of the first few principal components
  - $\boldsymbol{X}^{(k)} \approx \Phi \boldsymbol{a}^{(k)}$

- **Problem:** Reconstruction is not exact, and could be biased.
    How can we quantify the reconstruction error and correct the bias?

# Bias Correction and Error Estimation of Lossy Reconstruction

- We are interested in an expectation value of a function of statistical data
  - $\langle f(\boldsymbol{X}) \rangle = \frac{1}{N}\sum_k f(\boldsymbol{X}^{(k)}) \pm \frac{\sigma_{f(\boldsymbol{X})}}{\sqrt{N}}$

- Lossy reconstruction introduces error $\boldsymbol{X}^{(k)} \neq \boldsymbol{\Phi}\boldsymbol{a}^{(k)} \equiv \boldsymbol{X'}^{(k)}$

  Simple average with the reconstruction is a biased estimator $\langle f(\boldsymbol{X}) \rangle \neq \frac{1}{N}\sum_k f(\boldsymbol{X'}^{(k)})$

- **Key idea:** Unbiased estimator of $\langle f(\boldsymbol{X}) \rangle$ defined using small portion of original data

$$\overline{\boldsymbol{O}} = \frac{1}{N}\sum_{k=1}^{N} f(\boldsymbol{X'}^{(k)}) + \frac{1}{N_{bc}}\sum_{k=1}^{N_{bc}} \left( f(\boldsymbol{X}^{(k)}) - f(\boldsymbol{X'}^{(k)}) \right)$$

  - Bias correction is guaranteed: $\langle \overline{\boldsymbol{O}} \rangle = \langle f(\boldsymbol{X'}) \rangle + \langle f(\boldsymbol{X}) - f(\boldsymbol{X'}) \rangle = f(\boldsymbol{X})$
  - Risk is possibly large statistical error in real problems: The bias correction term (second term) increases the statistical error of $\overline{O}$ accounting for the reconstruction quality

Los Alamos
NATIONAL LABORATORY

# Statistical Error Increase due to Bias Correction

- For simplicity, consider a bias-corrected average of independent observables

$$\overline{\boldsymbol{O}} = \frac{1}{N}\sum_{k=1}^{N}\boldsymbol{X}'^{(k)} + \frac{1}{N_{bc}}\sum_{k=1}^{N_{bc}}\left(\boldsymbol{X}^{(k)} - \boldsymbol{X}'^{(k)}\right)$$

- Variance of the $i$-th component of $\overline{\boldsymbol{O}}$ can be obtained as

$$\sigma_{\overline{O}_i}^2 \approx \frac{1}{N}\sigma_{X_i'}^2 + \frac{1}{N_{bc}}\sigma_{X_i - X_i'}^2 \approx \frac{\sigma_{X_i'}^2}{N}\left(1 + \frac{N}{N_{bc}}\frac{\sigma_{X_i - X_i'}^2}{\sigma_{X_i'}^2}\right)$$

- Quality of lossy-compression on statistical data: $Q^2 \equiv \frac{1}{D}\sum_{i=1}^{D}\frac{\sigma_{X_i - X_i'}^2}{\sigma_{X_i}^2}$

  ➢ Smaller $Q^2$ indicates the better compression

$\left(\begin{array}{cc} N & \sigma^2 \end{array}\right)$
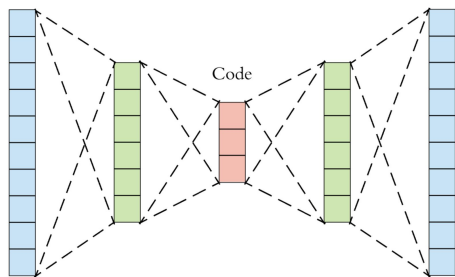
# Numerical Study with ML Compression Algorithms

- **Binary compression using D-Wave**
  - Find a set of vectors ($\Phi$) and their binary coefficients ($\boldsymbol{a}^{(k)}$) reconstructing $\boldsymbol{X}^{(k)}$

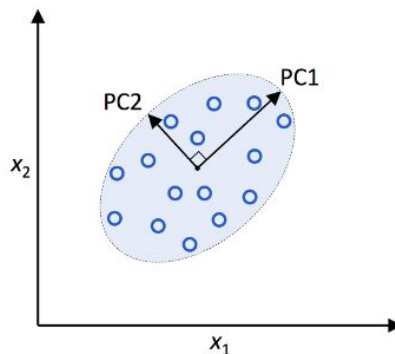$$\min_{\Phi} \sum_k \min_{\boldsymbol{a}^{(k)}} \left[ \sum_i \left( X_i^{(k)} - \left[ \Phi \boldsymbol{a}^{(k)} \right]_i \right)^2 \right]$$

- **Bottle-neck Autoencoder (AE)**
  - Fully connected NN with ReLU
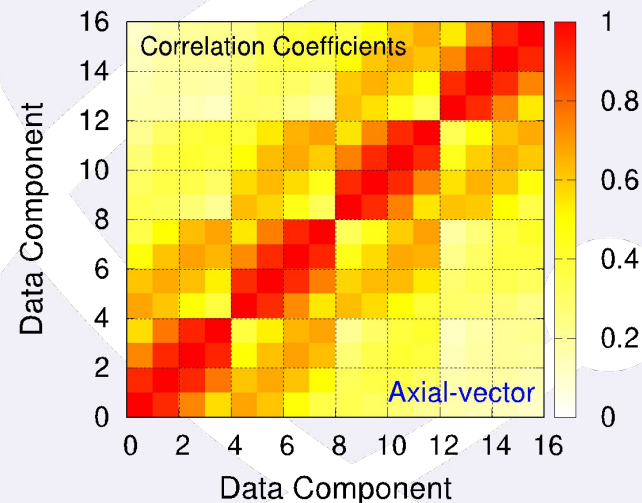  - Encoder: (16, 128, 64, 32, $N_z$)
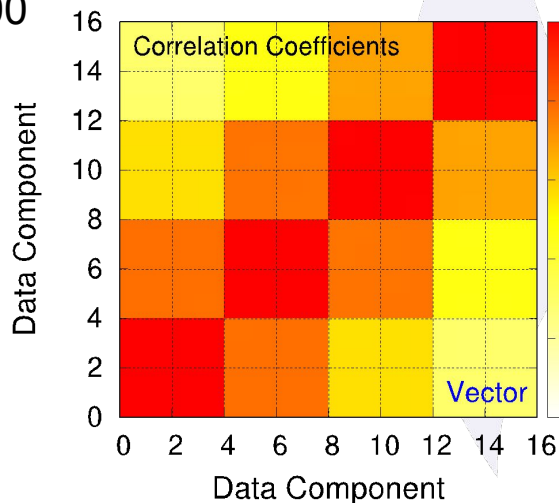  - Decoder: ($N_z$, 32, 64, 128, 16)

- **Principal Component Analysis (PCA)**
  - Compression by saving the first $N_z$ coefficients of the principal components
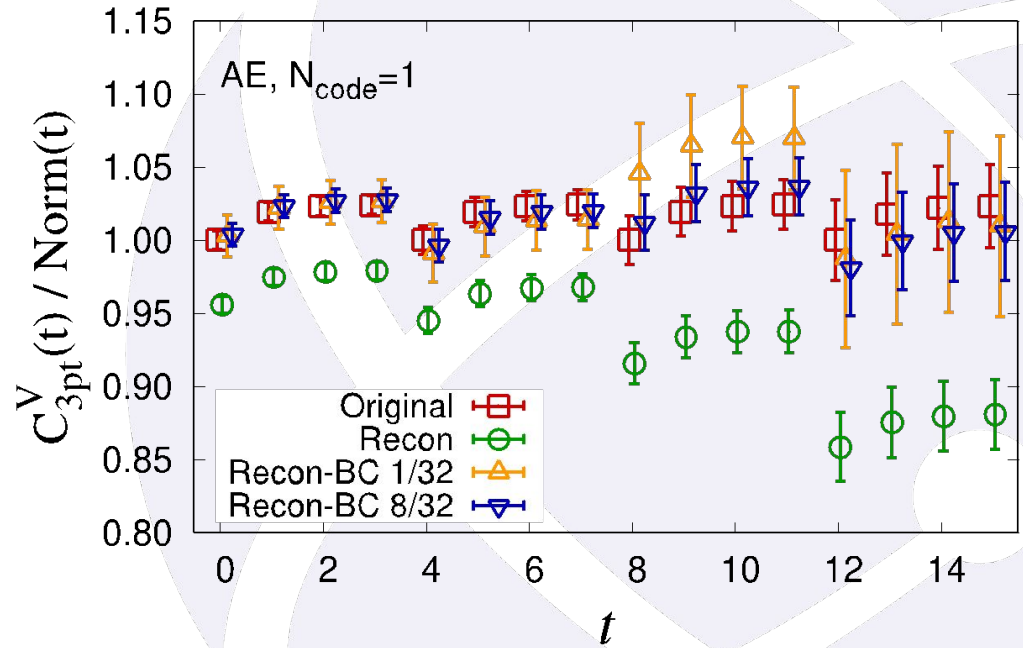
# Numerical Study with ML Compression Algorithms

- Data
  - Two different sets of simulation results in lattice quantum chromodynamics (QCD)
  - Describing the interactions between a nucleon and external currents
  - Vector length: D = 16
  - Data size: N = 12800
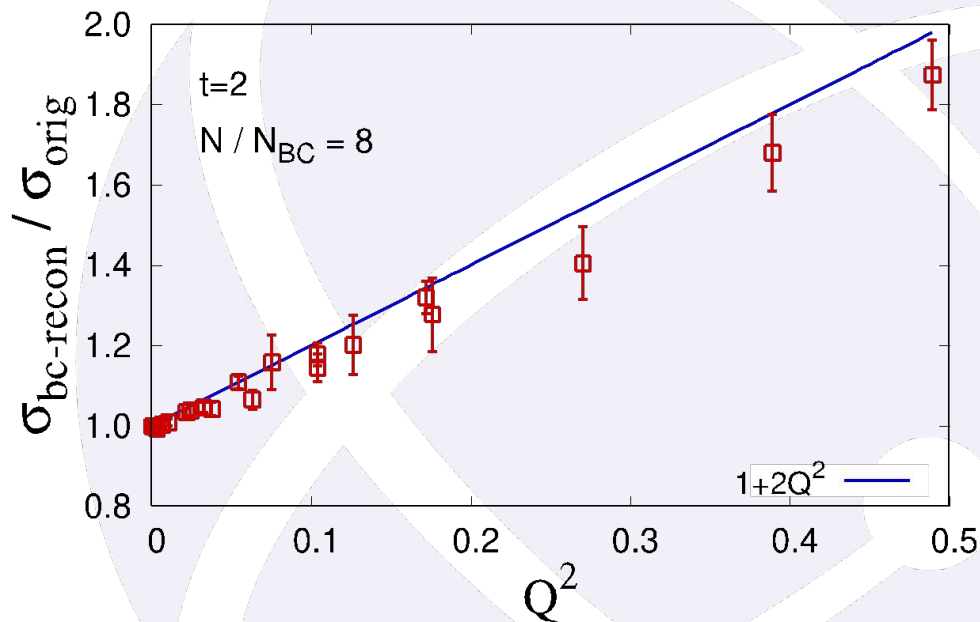  - Correlation pattern:

# Effect of Bias Correction

- 16 components of the data compressed using Autoencoder with $N_z = 1$

- Reconstruction without bias correction is biased

- Bias correction (yellow and blue) removes bias but increases statistical error

- The more bias correction data ($N_{bc}$) gives the smaller statistical error



AE, $N_{code}=1$

Legend:
- Original
- Recon
- Recon-BC 1/32
- Recon-BC 8/32

Y-axis: $C_{3pt}^V(t) / \mathrm{Norm}(t)$

X-axis: $t$

# Statistical Error Increase for Different Compression Qualities

- The statistical error increase is proportional to $Q^2 \equiv \frac{1}{D}\sum_{i=1}^{D} \frac{\sigma^2_{X_i - X'_i}}{\sigma^2_{X_i}}$

- For *independent data*, the error increase ratio due to bias correction is expected to be $1 + \frac{N}{2N_{bc}}Q^2$

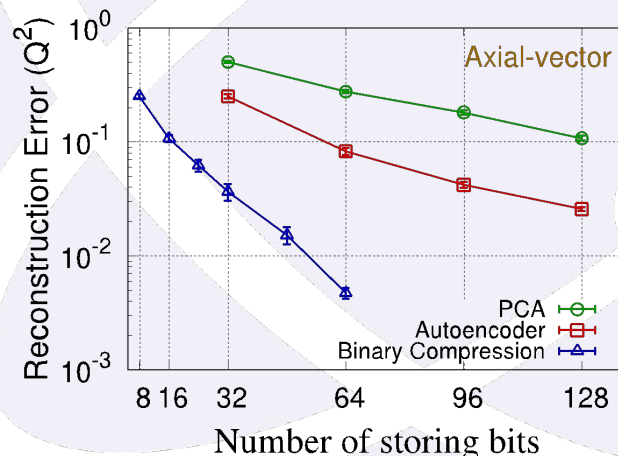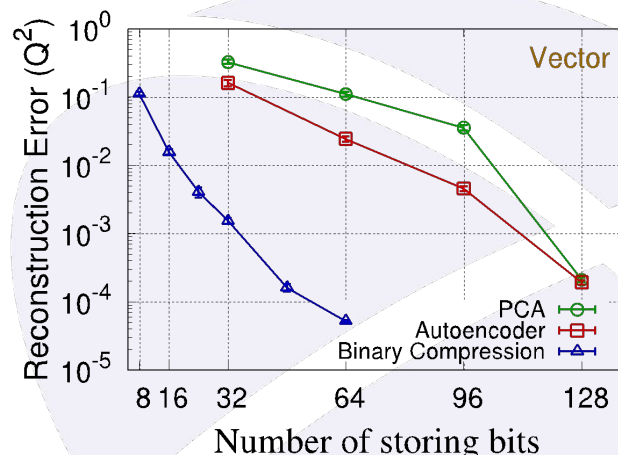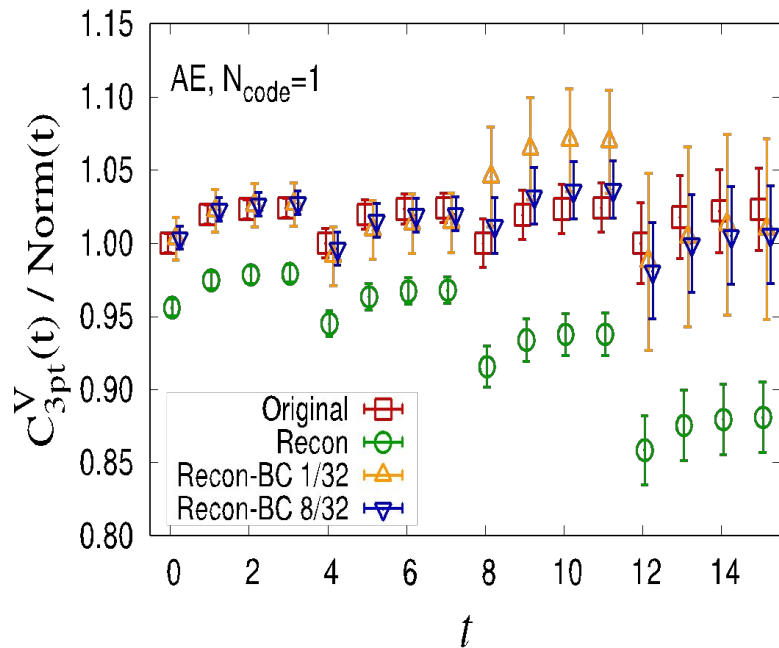- *Correlation between the data samples* makes it $1 + \alpha\frac{N}{2N_{bc}}Q^2$ with $0 < \alpha < 1$

# Expected Statistical Error Increase

- $\dfrac{\sigma_{bc-recon}}{\sigma_{orig}} = 1 + \alpha \dfrac{N}{2N_{bc}} Q^2$   with $0 < \alpha < 1$

- With 10% of bias correction data ($N/N_{bc}$ = 10) and $\alpha = 0.5$, expected error increase is $1 + 2.5Q^2$

- When $Q^2 = 10^{-2}$, expected error increase is 2.5%

- When $Q^2 = 10^{-3}$, expected error increase is 0.25%

- **Conclusion:**
  For good lossy compression algorithms, error increase due to bias correction is negligibly small

Figures: Compression of 16 floating-point numbers into $N_{bit}$ binary bits with three different compression algorithms

# Unbiased Reconstruction of Lossy-compressed Statistical Data



Reconstruction of lattice QCD data from a lossy compression with and without bias correction. Bias correction (yellow and blue) removes bias but increases statistical error accounting for the reconstruction quality of the compression algorithm. The more bias correction data gives the smaller statistical error.

*Project Description*

Demonstration of a novel bias correction algorithm for the reconstruction of lossy-compressed statistical data

*Project Outcomes*
- Showed that the bias correction algorithm removes the bias in the lossy reconstruction
- Demonstrated that statistical error increase due to the bias correction is small for good compression algorithms
- Paper in preparation

*PI: Boram Yoon (CCS-7)*
*Total Project Budget: $30K*
*ISTI Focus Area: Computational and Data Integrity*

# END