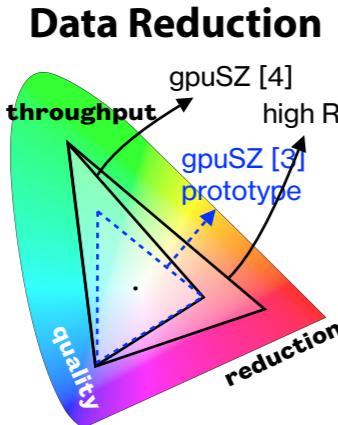


# Spline-Interpolation-Based Error-Bounded Lossy Compression for Scientific Data on GPUs

author : Jiannan Tian (IU)  
advisor : Dingwen Tao (IU)  
Sheng Di (ANL)  
Franck Cappello (ANL)

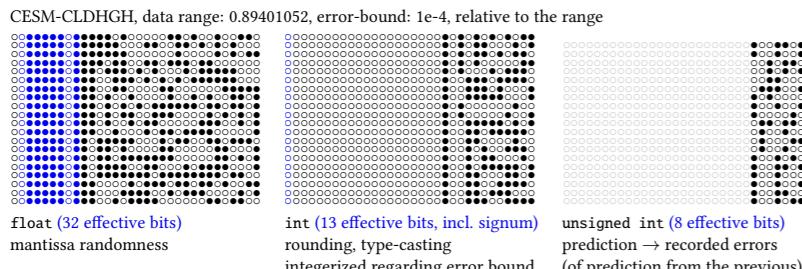
## Big-Data Scientific Application

application	data scale	to reduce
<b>HACC</b> cosmology simulation	<b>20 PB</b> per one-trillion-particle simulation	<b>10x</b> in need
<b>CESM</b> climate simulation	20% vs <b>50%</b> of h/w budget for storage 2013 vs 2017	<b>10x</b> in need
<b>APS-U</b> High-Energy X-Ray Beams Experiments	hundreds of PB brain initiatives	<b>100x</b> in need



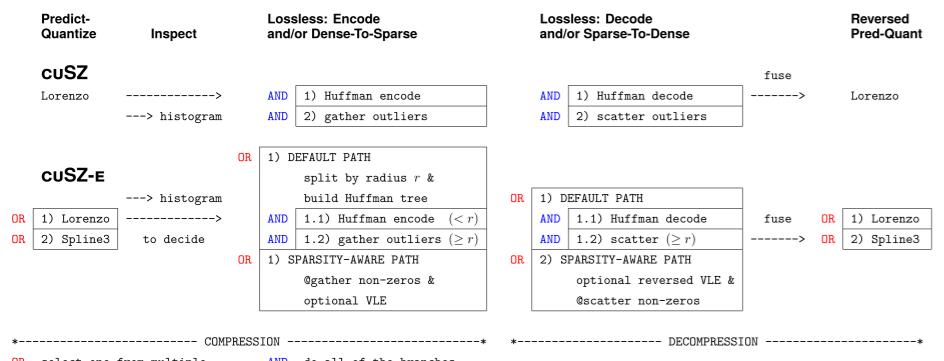
Data management is a real-world problem to address when we advance in scientific exploration.

## SZ [1, 2] Lossy Compression Essence



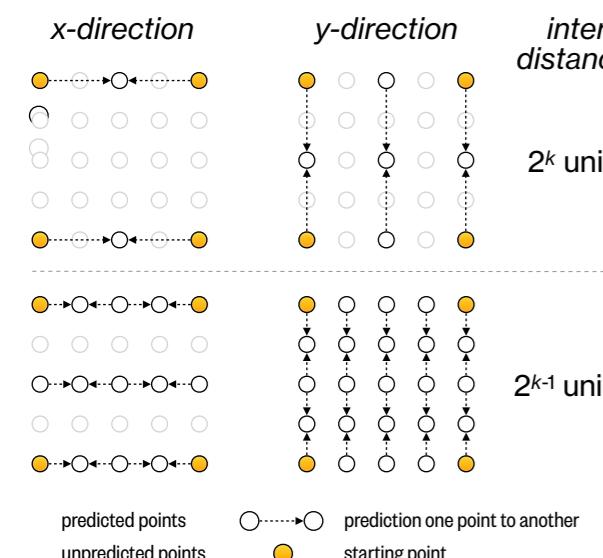
SZ, while guaranteeing the error-boundness, performs bit-level randomness elimination to increase the compressibility.

To lower bit randomness: prediction-based SZ.



When the prediction does well and results in highly sparse error-quantization code, gather (& scatter) are conducted to boost the compression ratio.

## Interpolative Prediction Data-Access Pattern



We use 2D interpolation for conceptual demonstration; the prediction direction alters each stage.

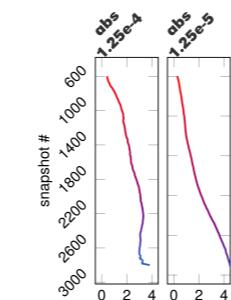
An interpolation iteration contains  $ndim$  stages. Stages feature altering interpolation direction. When an iteration finishes, the interpolation distances shrinks by a factor of 2.

Starting points ● are distant from each other and beyond GPU thread blocks can handle; they become **anchor points** and are saved directly.

## Preliminary Evaluation (Seismic Data)

### Kernel Throughput

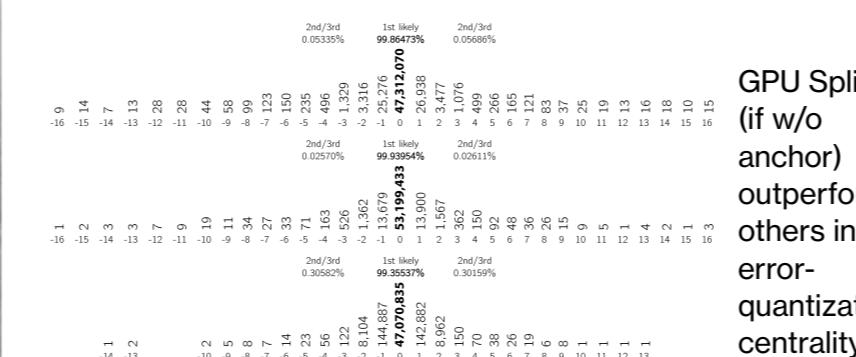
### PSNR Advantage over Default



Spline interpolation outperforms the default predictor (Lorenzo).

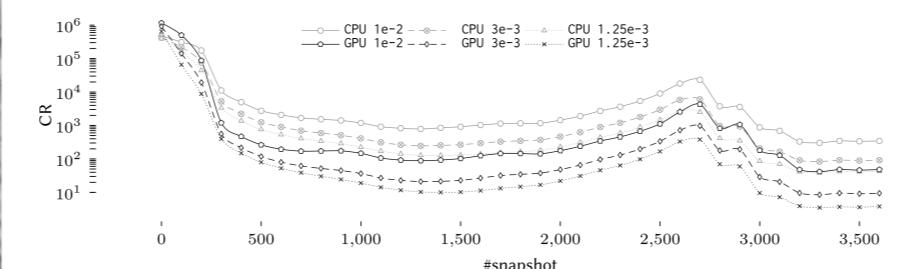
Need further in-depth optimization.

## Error-Quantization Distribution



GPU Spline (if w/o anchor) outperforms others in error-quantization centrality.

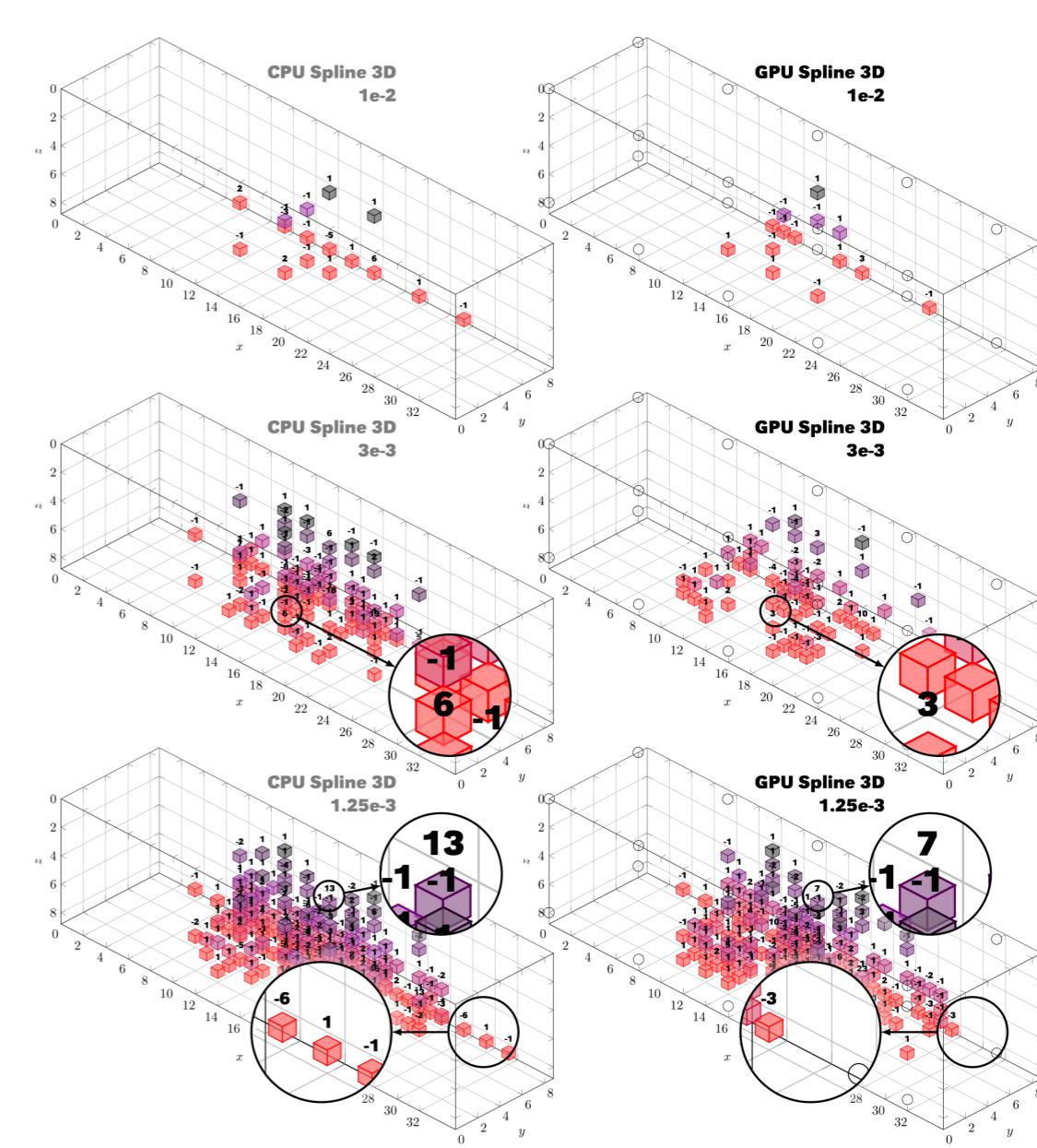
## Compression Ratio Throughout 3,000 RTM Snapshots



## Compression Ratio & Quality Across {GPU, CPU}-Spline and Default (Lorenzo)

eb	quant.	GPU Spline Prediction			CPU Spline Prediction			GPU Lorenzo Prediction	
		freq.	CR	quality	freq.	CR	quality	CR	quality
1.25e-3	-1	0.2935%	67.0	G	0.5149%	35.4	G		
	0	99.2538%	59.3	PSNR	98.5888%	70.9	PSNR	69.7	G+VLE
	-1	0.2870%	118.5	0.9979	0.4973%	397.5	0.9983		
3.00e-3	-1	0.1150%	177.0	G	0.2382%	82.4	G		
	0	53.081227%	131.5	PSNR	47.078719%	164.8	PSNR	63.6	G+VLE
	-1	0.1156%	263.0	0.9937	0.2385%	938.8	0.9932		
1.00e-2	-1	0.0257%	827.0	G	0.0534%	369.2	G		
	0	99.9395%	316.2	PSNR	47.312079%	738.4	PSNR	56.2	G+VLE
	-1	0.0261%	632.5	0.9809	0.0569%	3606.7	0.962		

## Impact on Error-Quantization



Interpolating 32x8x8 data chunk. With anchor points, the error quantizations are smaller in amplitude throughout different error bounds. Anchor points are used for performance concern: fitting subproblem size to GPU hardware architecture.

## Acknowledgment

This R&D was supported by the Exascale Computing Project (ECP), Project Number: 17-SC-20-SC, a collaborative effort of two DOE organizations – the Office of Science and the National Nuclear Security Administration, responsible for the planning and preparation of a capable exascale ecosystem. This repository was based upon work supported by the U.S. Department of Energy, Office of Science, under contract DEAC02-06CH11357, and also supported by the National Science Foundation under Grants SHF-1617488, SHF-1619253, OAC-2003709, OAC-1948447/2034169, and OAC-2003624.

- [1] S. Di and F. Cappello, "Fast error-bounded lossy HPC data compression with SZ," in 2016 IEEE International Parallel and Distributed Processing Symposium, Chicago, IL, USA: IEEE, 2016, pp. 730–739.
- [2] D. Tao, S. Di, Z. Chen, and F. Cappello, "Significantly improving lossy compression for scientific data sets based on multidimensional prediction and error-controlled quantization," in 2017 IEEE International Parallel and Distributed Processing Symposium, Orlando, FL, USA: IEEE, 2017, pp. 1129–1139.
- [4] J. Tian et al., "cuSZ: An efficient gpu-based error-bounded lossy compression framework for scientific data," in Proceedings of the ACM International Conference on Parallel Architectures and Compilation Techniques, 2020, pp. 3–15.
- [5] J. Tian et al., "Optimizing Error-Bounded Lossy Compression for Scientific Data on GPUs." Proceedings of the 2021 IEEE International Conference on Cluster Computing, (Virtual Event) Portland, OR, September 7–10, 2021.

The color palette is used to increased the sense of space description.  
The larger z value, the more purple.