# **CUSZ**: A High-Performance **GPU** Based **Lossy Compression** Framework for **Scientific Data**

**Jiannan Tian** Washington State University Sheng Di Argonne National Laboratory University of California, Riverside Kai Zhao Cody Rivera The University of Alabama Megan Hickman Fulp Clemson University Robert Underwood Clemson University Sian Jin Washington State University Xin Liang Oak Ridge National Laboratory Jon Calhoun Clemson University Dingwen Tao Washington State University Argonne National Laboratory Franck Cappello

October 5, 2020 PACT '20, Virtual Event









Background	Introduction	Design	Evaluation	Conclusion
●000	00	000000	00000	00

#### Trend of Supercomputing Systems Gap Between Compute and I/O



The compute capability is ever growing while storage capacity and bandwidth are developing **more slowly** and not matching the pace.

supercomputer	year	class	PF	MS	SB	MS/SB	PF/SB
Cray Jaguar	2008	1 pflops	1.75 pflops	360 тв	240 дв/з	1.5k	7.3k
Cray Blue Waters	2012	10 pflops	13.3 pflops	1.5 рв	1.1 тв/ѕ	1.3k	13k
Cray CORI	2017	10 pflops	30 pflops	1.4 рв	1.7 тв/s•	0.8k	17k
IBM Summit	2018	100 pflops	200 pflops	>10 рв●●	2.5 тв/s	>4k	80k

PF: peak FLOPS MS: memory size SB: storage bandwidth

• when using burst buffer •• counting only DDR4

Source: F. Cappello (ANL)

 Table 1: Three classes of supercomputers showing their performance, MS and SB.

data scale

**20 PB** 

simulation

**Current Status of Scientific** 

**Applications: Big Data** 

Design

Evaluation

Conclusion



Argonne 🖂

passive solution (?)

use up FS

26 PB for Mira@ANL

to reduce **10x** in need

CESM climate simulation

cosmology simulation

application

HACC

20% vs **50%** of h/w budget for storage 2013 vs 2017

per one-trillion-particle

5h30m to store NSF Blue Waters, I/O at 1 TBps

**10x** in need

**APS-U** High-Energy X-Ray **Beams Experiments**  hundreds of **PB** brain initiatives

100-PB buffer or, connection at 100 GBps

**100x** in need

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 3 / 20

**Error-Bounded Lossy Compression Matters** 

lossless-compress scientific datasets

**2:1** (FP-type) **10:1** or higher reduction ratio in need

distinct in design goals

industry lossy compressor

e.g., JPEG, MPEG

absolute error bound ( $L^{\infty}$  norm) pointwise relative error bound RMSE error bound ( $L^2$  norm) fixed bitrate

despite high reduction rate, not suitable for HPC

need diverse compression modes

SZ [Di and Cappello 2016; Tao et al. 2017; Xin et al.2018]

- prediction-based lossy compressor framework for scientific data
- strictly control the global upper bound of compression error

Design

Evaluation

Conclusion





Lossy compression for scientific data at varving reduction rate (10:1 to 250:1, left to right) Figure from Peter Lindstrom (LLNL)



aithub.com/szcompressor/SZ

Background	Introduction	Design	
0000	00	000000	
C7 Eromo	work		

#### SZ Framework (Error-Bound Workflow)

Evaluation

Conclusion





# Motivation, Challenge, and Contribution

#### **Research Objective and Contribution**

CUSZ is the first strictly error-bounded lossy compressor on GPU for scientific data.

#### Challenge

- Tight data dependency (loop-carried RAW) hinders parallelization.
   solution Eliminate dependency and parallelize it.
   DUAL-QUANTIZATION: {PRE, POST}QUANTIZATION
- Host-device communications only considering CPU/GPU suitableness.
   solution All tasks are done on GPU. Histograming [Gómez-Luna et al.]

Customized Huffman codec (corse-grained)

**EV**a 00

Design

Evaluation









# **Loop-Carried Read-After-Write** (P+Q) Procedure in SZ

- Lossless compression and decompression (codec) are mutually reversed procedures.
- Simlarly, SZ makes to-be-decompressed (reconstructed) data show during compression and make it under error control
- Error control is conducted during quantization and reconstruction:

 $\mathbf{e}^{\circ}/(2 \cdot \mathbf{e}\mathbf{b}) \times (2 \cdot \mathbf{e}\mathbf{b}) - \mathbf{e}^{\circ} < \mathbf{e}\mathbf{b}.$ 

This introduces loop-carried read-after-write dependency.



Evaluation

Conclusion





••••••

# Fully Parallelized (P+Q) Procedure in CUSZ

Design ○●○○○○ Evaluation

Conclusion





- ▶ Prioritize error control.
- Error control happens at the very beginning, prequantization:

 $\mathbf{d}^{\circ}/(2 \cdot \mathbf{e}\mathbf{b}) \times (2 \cdot \mathbf{e}\mathbf{b}) - \mathbf{d}^{\circ} \leq \mathbf{e}\mathbf{b},$ 

And postquantization is corresponding to quantization in SZ.





October 5, 2020  $\,\cdot\,$  PACT '20, Virtual Event  $\,\cdot\,$  cuSZ  $\,\cdot\,$  10 / 20

Canonical Codebook and Huffman Encoding

ca·non·i·cal adj.

A canonical encoding is then generated in which the numerical values of the codes are monotone increasing and each code has the smallest possible numerical value consistent with the requirement that the code is not the prefix of any other code.

[Schwartz and Kallick 1964]

- codebook transformed to a compact manner
- ► no tree in decoding
- tree build time: 4-7 ms (for now)
- canonize for 200 us (1024 symbols)

Design ○○○●○○ Evaluation

Conclusion



- Encoding/decoding is done in a coarse-grained manner.
- ► A GPU thread is assigned to a data chunk.
- Tune degree of parallelism to keep every thread busy.

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 11 / 20

Background

Introduction

Design 000000 Evaluation



Conclusion

### **Mixture of Different Parallelisms**



**Table 2:** Parallelism used for cuSZ's subprocedures
 (kernels) in compression and decompression.

atomic . •

.

Worth noting: in canonizing codebook

- problem size > max. block size (1024)
- utilize cooperative groups and arid.sync()
- syncthreads(): not able
- cudaDeviceSynchronize(): expensive

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 12/20

Backgr 0000	ound		<b>In</b> :	troductio	n		De oc	sign 0000●			Evaluatio	n		<b>Cor</b> 00	nclusion
Tur (De	ning eare	Coar e of F	'se- Para	<b>Gra</b>	inec m/C	l Hu onc	uffm urre	an ( ent T	Cod Thre	<b>lec</b> ead	Nun		UNIVERSITY		VERSIDE
	HACC			CESM			HURR	RICANE		NYX			QMCP	PACK	IONAL LABORATORY
CHUNK SIZE	1071.8 M	B 280,953,86	57 f32	24.7 MB #THREAD	6,480,000	f32	95.4 MB #тнреар	25,000,000	132	512 MB #тнреар	134,217,72	1NELATE	601.5 MB	157,684,3 DEELATE	20 f32
2 <sup>6</sup>				1.0e5	11.3	25.0									
27				5.1e4	15.5	37.8									
2 <sup>8</sup>				2.5e4	67.1	41.6	9.8e4	5.1	11.0						
2 <sup>9</sup>				1.3e4	55.6	30.7	4.9e4	10.2	9.4						
2 <sup>10</sup>				6.3e3	48.2	19.6	2.4e4	64.6	34.2	1.3e5	4.7	5.9	1.5e5	4.7	5.1
2 <sup>11</sup>	1.4e5	4.6	2.8				1.2e4	57.3	27.7	6.6e4	5.7	6.3	7.7e4	5.2	6.2
2 <sup>12</sup>	6.9e4	5.1	5.1				6.1e3	50.7	17.8	3.3e4	25.1	16.1	3.8e4	12.9	11.1

**Table 3:** Throughputs (in GB/s) versus different numbers of threads launched on V100. The optimal thread number in terms of inflating and deflating throughput is shown in bold.

2<sup>13</sup>

2<sup>14</sup>

2<sup>15</sup>

2<sup>16</sup>

3.4e4

1.7e4

8.6e3

4.3e3

13.6

63.1

65.8

45.9

12.1

35.0

28.1

14.3

1.6e4

8.2e3

4.1e3

69.7

72.4

50.0

52.4

42.6

23.1

1.9e4

9.6e3

4.8e3

72.7

75.9

56.0

40.3

29.0

16.1

# Evaluation Setup: Platform and Dataset

#### Evaluation Platform (UA PantaRhei cluster)

- NVIDIA V100 (SXM2, 16 GB)
- Dual 20-core Intel Xeon Gold 6148 CPUs
- 16-lane PCIe 3.0 interconnect
- Comparison Baselines
  - (algorithmic) SZ-1.4.13.5: 16-bit quantization
  - ► (performance) cuZFP 0.5.5
- Test Datasets (from SDRB)
  - ► 1D HACC, cosmology particle simulation
  - 2D CESM-ATM, climate simulation
  - 3D HURRICANE, ISABEL simulation
  - 3D NYX, cosmology simulation
  - 4D QMCPACK, quantum Monte Carlo simulation

-	-			

Evaluation

•0000



DATASETS	ТҮРЕ	DATUM SIZE DIMENSIONS	#FIELDS EXAMPLE(S)
COSMOLOGY	fp32	1,071.75 MB	6 in total
HACC		<b>280,953,867</b>	x, vx
CLIMATE	fp32	24.72 MB	79 in total
CESM-ATM		1,800×3,600	CLDHGH, CLDLOW
CLIMATE	fp32	95.37 MB	20 in total
Hurricane		100×500×500	CLOUDf48, Uf48
cosmology	fp32	512.00 MB	6 in total
Nyx		512×512×512	baryon_density
QUANTUM	fp32	601.52 MB	2 formats in total
QMCPACK		288×115×69×69	einspline

Table 4: Real-world datasets used in evaluation.

Design

Background	
0000	

Design

Evaluation 00000

Conclusion 00

UC RIVERSIDE

WASHINGTON STATE

UNIVERSITY

# **Breakdown Evaluation** of

of P	Perfor	manc	e							EMSON AI	gonne
	1	P)REDICT. (Q)UANT.	н	UFFMAN CODING		KERNEL COMP.	GPU2CPU VALREL@10 <sup>-4</sup>	OVERALL COMPRESS	HUFFMAN DECODING	REVERSED (P+Q)	KERNEL DECOMP.
		MB/S		MB/S				MB/S	MB/S	MB/S	MB/S
CPU-SZ	ИАСС	137.7		328.6		-	-	94.1	196.0	659.3	151.1
	CESM-ATM	105.0		459.1		-	-	85.5	502.2	451.9	237.9
	HURRICANE	93.8		504.0		-	-	78.5	524.5	306.8	185.0
	NYX	98.5		648.7		-	-	84.7	670.4	300.5	201.8
	QMCPACK	97.5		396.2		-	-	80.8	660.3	313.4	211.1
			HISTO.	DICT.	ENC.				CANONICAL		
		GB/S	GB/S	MS	GB/S	GB/S	GB/S	GB/S	DEC. GB/S	GB/S	GB/S
cuSZ	HACC	207.7	602.8	5.16	54.1	40.0	53.2	22.8	35.0	16.8	11.4
	CESM-ATM	252.1	345.3	4.33	57.2	41.1	81.9	27.4	41.6	58.5	24.3
	HURRICANE	175.8	418.0	4.81	55.2	38.2	40.8	19.7	34.2	43.9	19.2
	NYX	200.2	427.6	3.84	58.8	41.1	134.1	31.6	52.4	29.7	19.0
	QMCPACK	189.6	346.1	4.09	61.0	40.7	99.2	28.9	40.3	22.4	14.4
cuZFP	HACC	-	-	-	-	-	-	_	-	-	-
	CESM-ATM	-	-	-	-	47.6	27.7	17.5	-	-	113.1
	HURRICANE	-	-	-	-	83.7	27.7	20.8	-	-	102.2
	NYX	-	-	-	-	71.3	56.3	31.7	-	-	103.1
	QMCPACK	-	-	-	-	72.6	42.5	26.8	-	-	115.5

Table 5: Breakdown comparison of kernel performance among CPU-SZ, cuSZ, and cuZFP. "-" for N/A.

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 15 / 20

**Design** 

Evaluation



Conclusion

### Performance Evaluation: Serial, OpenMP and CUDA



October 5, 2020 · PACT '20, Virtual Event · cuSZ · 16 / 20



Design

**Bitrate** 

Introduction

Background

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 17 / 20

Evaluation

Conclusion



**Evaluation of Compression Quality** 

# in Rate-Distortion (2/2)



Conclusion





Design

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 18 / 20

Background Introduction

#### Acknowledgement (Exascale Computing Project)

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 DE-AC02-06CH11357, and also supported by the National Science Foundation under Grants SHF-1617488, SHF-1619253, OAC-2003709, OAC-1948447/2034169, and OAC-2003624.







Design





Conclusion





EXASCALE COMPUTING PROJECT

October 5, 2020 · PACT '20, Virtual Event · cuSZ · 19 / 20





#### github.com/szcompressor/cuSZ



# BackUp (*l*-Predictor)



► Gaussian-like, with signum altering to Manhattan distance to the (polarized) current point (■).

$$\mathbf{G}_{5\times5} = \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix} \quad \ell_{5\times5} = \begin{bmatrix} -1 & 4 & -6 & 4 & -1 \\ 4 & -16 & 24 & -16 & 4 \\ -6 & 24 & -36 & 24 & -6 \\ 4 & -16 & 24 & -16 & 4 \\ -1 & 4 & -6 & 4 \end{bmatrix}$$

Works for arbitrary dimension: from line to cube, to hypercube...



#### BackUp (GPU Building Huffman Tree)



Sequential GPU building Huffman tree is too slow (1024 symbols for 4 ms).

#### Improvement introduced by regularize memory access

- Our preliminary improvement by switching to thrust is 4×, from 4 ms to 1 ms.
- It is worthy maintaining workload on GPU for simplifying workflow.
- Our customized Huffman coding serves for HPC scenarios.
  - Snapshots show high similarity across consecutive timestamps.
  - So, we may only need a quasi-optimal tree for a large group of snapshots.
  - Hence, under some circumstances, tree building can be hidden.

#### BackUp (State-of-the-Art SIMD)



As we can see, SZ is a log(n) (linear-time) algorithm. Due to the low computational intensity, the performance is generally bounded by memory bandwidth. Our focus is to develop GPU version of SZ.

		RTX 2060S	RTX 5000	Tesla V100	Tesia A100
specification	compute (FP32 TFLOPS)	7.18	11.15	14.13	19.49
	#multiprocessor (SM)	34	48	80	108
	memory bandwidth (св/s)	448	448	897	1555
dual-quant	absolute perf. (gв/s)	47.4 (100.0%)	63.0 (132.9%)	252.1 (531.9%)	?
	normalized (#SM)	1.39 (100.0%)	1.31 ( 94.2%)	3.15 (226.6%)	?
	normalized (#SM+mem.bw)	3.11 (100.0%)	2.93 ( 94.2%)	3.51 (112.9%)	?

#### **Case Study of Compression Quality: Statistical Information**

FIELD	SZ-1.4	cuSZ	FIELD	SZ-1.4	cuSZ
CLOUDf48	84.99	94.18	QSNOWf48	84.31	93.36
CLOUDf48.log10	84.51	87.17	QSNOWf48.log10	84.87	84.93
Pf48	84.79	84.79	QVAPORf48	84.79	84.80
PRECIPf48	85.35	92.86	TCf48	84.79	84.79
PRECIPf48.log10	84.82	84.77	Uf48	84.79	84.79
QCLOUDf48	85.03	98.91	Vf48	84.79	84.79
QCLOUDf48.log10	85.22	95.21	Wf48	84.79	84.79
QGRAUPf48	88.21	97.02	baryon_density	89.71	98.25
QGRAUPf48.log10	84.90	84.82	dark_matter_density	86.57	87.77
QICEF48	84.61	95.51	temperature	84.77	84.77
QICEf48.log10	85.56	85.77	velocity_x	84.77	84.77
QRAINf48	85.36	97.37	velocity_y	84.77	84.77
QRAINf48.log10	84.93	84.56	velocity_z	84.77	84.77
Hurricane avg.	85.01	86.96	Nyx avg.	85.58	85.98

**Table 6:** Comparison of PSNR between cuSZ and
 SZ-1.4 on Hurricane (FIRST 20) and Nyx (LAST 6) under valrel =  $10^{-4}$ .





#### CLOUD<sub>f48</sub>



Table 7: Statistical information (percentile) of example fields having high PSNR under valrel =  $10^{-4}$ . The range of eb or even  $\frac{1}{10}$  eb at 0 or min value cover a majority of data in the fields.