Evaluating SZ3 Compressor Performance on High Energy Physics Data

Amy Byrnes, Serhan Mete, Jahred Adelman, Peter van Gemmeren, Michael E. Papka

Motivation

- The Large Hadron Collider (LHC) has created hundreds of petabytes of data. The High-Luminosity LHC (HL-LHC) will produce orders of magnitude more.
 - The HL-LHC is a major upgrade to the LHC, aiming to be operational by the end of 2030. It will produce more data in one year than the LHC produced during its first 10 years of operation.
- Significant R&D efforts are needed to handle the incoming tide of data.

Without using new technologies, the data volume of the HL-LHC will rapidly outgrow CERN's compute and storage capacities. This means everything is on the table, even historically disfavored methods like lossy compression.

Evaluation

- ALGO_NOPRED reduces data sizes by 66-83%, while keeping the data within 0.1-2% of its original values, on average. This is better than any of SZ3's other (lossy) compression methods.
 ALGO_NOPRED achieves the highest compression ratios outright, and the least error and fastest compression/decompression times on average.
- Events are independent, negating the assumptions that make prediction effective. Data are recorded per-collision-event, and the outcomes of these events are independent from one another. There is no meaningful relationship between adjacent data points. ALGO_NOPRED is the only lossy method in SZ3 that skips prediction entirely. Its superior performance supports our intuition that this data is not suited for prediction-based methods.
 Some data is more compressible than other data, but more investigation is needed to understand why. The data which achieves the highest compression ratio of 6 has, predictably, the lowest entropy. But other data which compress relatively well (ratios of 4.5, 5.4) have entropies equivalent to data that compress relatively poorly (2.9, 3).

ATLAS Data

- ATLAS (A Toroidal LHC ApparatuS) is one of two generalpurpose detectors at the LHC. ATLAS and CMS (the other general-purpose detector) are responsible for the observation of the Higgs boson in 2012.
- The data collected by ATLAS represents physics objects resulting from particle collisions. Hadronic jets represent partons (quarks and gluons), which cannot be measured directly.
- Quantities collected include: Energy (E), Transverse momentum (pt), Pseudo-rapidity (eta), and Azimuthal angle (phi)
- Many TB of real and simulated ATLAS data are available on the CERN Open Data portal – including the data used in this work. The data are provided under a Creative Commons C0 license.

SZ3 Compression Framework



 SZ3 is "a modular, error-bounded lossy compression framework for scientific

F	Compression ratio													-
L L														6
it	ALGO NOPRED	4.460	3.490	3.070	5.470	3.060	4.540	6.040	4.100	2.950	2.960	4.170		5.5
ression algor														5
	ALGO LORENZO REG	3.920	2.910	2.770	3.710	2.720	3.910	5.390	3.660	2.570	2.590	2.600		4.5
	/													4
	LGO INTERP LORENZO	3.630	2.700	2.610	3.930	2.530	3.620	4.950	3.450	2.390	2.640	2.720		35
														2.5
	ALGO_INTERP	3.605	2.675	2.600	3.500	2.520	3.595	4.910	3.430	2.365	2.310	2 265		3
												21205		2.5
d	ALCO LOSSIESS (zetd)	1 1 1 0	1 1 5 0	1 000	1 410	1 080	1 1 1 0	1 1 2 0	1 1 2 0	1 000	1 080	1 390		2
3	ALGO_LOJJLLJJ (23td)	1.110	1.150	1.050	1.410	1.000	1.110	1.150	1.150	1.050	1.000	1.550		1.5
^o		J _O ,	Jo.	J _O ,	Jo.	J _{Or}	J _O *	J _O ,	12	12	12	12		
\cup		- ~ (<u>, </u>	4, 5	$\sum_{x \in X} x \in X$		$b_{\lambda} \stackrel{c}{\searrow} \lambda$	$b_x \sim c_x$	Sx S		so 'G	90° 00	5	

datasets". SZ3's flagship features are its modularity, making it straightforward to build a custom compression pipeline, and its use of bestfit predictors. It's source code is available on GitHub under a BSD license.

Predictors are generally used in compression pipelines to reduce the entropy of the input

- data. This ultimately leads to smaller outputs. All
 predictors rely on the assumption that the value
 currently being compressed can be predicted with
 reasonable accuracy from some small number of
- previous values.
- SZ3's predictor attempts to fit the current data to a (linear, quadratic) curve with the immediately preceding data points. The initial SZ compressor was developed with scientific simulation data in mind. This prediction model is well suited to data where indices correspond to contiguous physical coordinates (e.g. temperature)
 - contiguous physical coordinates (e.g. temperature at a point (x, y, z)).

The "loss" incurred during lossy compression is the result of quantization. Scientists are



Compressed Data

understandably hesitant about losing any amount of precision from their data. Quantization is a welldefined process, so maximum error bounds can be imposed. SZ3 defaults to lossless compression with zstd if it detects that lossy compression is not possible with the given error bound.

Acknowledgements

 This work is supported by the Chicagoland Computational Traineeship in High Energy Particle Physics (C²P²) under US Department of Energy Office of Science grant DE-SC0023524.

