

THE UNIVERSITY OF

Optimizing Huffman Decoding for Error-Bounded Lossy Compression on GPUs

Cody Rivera (cjrivera1@crimson.ua.edu) Sheng Di (sdi1@anl.gov) Jiannan Tian (jiannan.tian@wsu.edu) Xiaodong Yu (xyu@anl.gov) Dingwen Tao (dingwen.tao@wsu.edu) Franck Cappello (cappello@mcs.anl.gov)



Background: Use Cases

- > HPC Application Needs
- HPC applications are generating **increasingly large amounts of data**
- Applications include large scale simulations, deep neural networks
- e.g., Hardware/Hybrid Accelerated Cosmology Code (HACC) (S. Habib et. al.) [1], which generates roughly 22 petabytes per simulation run
- In-Memory Caching
- Fast memory, especially on GPUs, is a scarce resource
- Can cache data more economically by compressing and decompressing intermediate data
- e.g., Wu *et. al.*'s work on quantum circuit simulation [17], where compression reduces total RAM usage from 32 exabytes to 768 terabytes







Background: Lossy Compression/SZ

- > Lossy compression on scientific data
- Offers much higher compression ratios than lossless compression by trading a little bit of accuracy
- An example: **SZ**, a prediction-based lossy compression that achieves high compression ratios [5]
 - Actively developed and researched
 - CPU, GPU (cuSZ), and domain-specific (DeepSZ, PastriSZ) versions
- We focus on SZ/cuSZ (over ZFP) for the following reasons:
 - Less distortion/higher PSNR at a given bitrate
 - Compression error can be explicitly **bounded** by the user







An Overview of cuSZ







> Prediction

 Predict data points using a data-fitting Lorenzo predictor (Ibarria et. al.) [7]

> Quantization

- Determine the prediction error for each point and classify it as an integer we call a quantization code
- > Huffman Coding
- Losslessly compress quantization codes



Background: Huffman Coding

- > Huffman Coding
- Classic lossless variable-length compression technique introduced by David Huffman in 1952
- Example: ABAACDAA (16 bits at 2 bits per character)
 - Encoded Text: 1010110110011 (13 bits)
 - Compression Ratio: 16/13 ≈ 1.23



Example Huffman Tree and Codebook



Motivation: Why Optimize Decoding?

cuSZ's current Huffman coding

- Encoding performance: average 25.7 GB/s in production [9], 138.3 GB/s experimentally (J. Tian et. al., IPDPS '21) [10]
- Decoding performance: average 32.3 GB/s

Research Focus

- Decompression is needed for data post-analysis as well as retrieval from in-memory caches
- Huffman **decoding**, however, is the primary bottleneck for cuSZ, taking up 83% of the time in a recent version





Scattering Outliers



Increasing Parallelism in Decoding





- Coarse-Grained Parallelism
- Many points per thread, few threads
- cuSZ's current Huffman decoder
- Does not map well to GPU architectures

- Fine-Grained Parallelism
- Few points per thread, many threads
- Maps more effectively to the GPU's massive parallelism



An Insight from Information Theory

> The Self-Synchronization Property

- Tendency for Huffman codes to **correct** themselves if a few bits were skipped, first written about by Ferguson and Rabinowitz [11]
- Example: 111000010111000
 - Correct decoding
 - (11)(10)(00)(010)(11)(10)(00)
 - CBADCBA
 - Incorrect decoding
 - 1(11)(00)(00)(10)(11)(10)(00)
 - CAABCBA
 - Is eventually correct

Symbol	Codeword
А	00
В	10
С	11
D	010
E	011

Self-Synchronizing Codebook

- Consider the message
- BACACCBDBAAEBBA

Symbol	Codeword
А	00
В	10
С	11
D	010
E	011



> Beginning of the procedure

- Legend:
- Thread Position
- 1 Synchronization Point
- \checkmark Verified Synchronization Point

Symbol	Codeword
А	00
В	10
С	11
D	010
E	011





Each thread decodes a subsequence Legend:

- Thread Position
- 1 Synchronization Point
- Verified Synchronization Point

Symbol	Codeword
A	00
В	10
С	11
D	010
E	011





 Synchronization points are initialized Legend:

- Thread Position
- Synchronization Point
- \checkmark Verified Synchronization Point

Symbol	Codeword
A	00
В	10
С	11
D	010
E	011





 Each synchronization point is verified by the previous thread Legend:

- Thread Position
- Synchronization Point
- Verified Synchronization Point

Index 0, Index 1, Index 2, Index 3

Symbol	Codeword
А	00
В	10
С	11
D	010
E	011



Subsequence 2, 3





- Once this is done, each thread will decode parts of the following correctly
- BACACCBDBAAEBBA

Legend:

- Thread Position
- Synchronization Point
- \checkmark Verified Synchronization Point

Symbol	Codeword
А	00
В	10
С	11
D	010
E	011





Another Approach to Fine-grained Parallelism

> Gap Arrays

- Determining synchronization points requires redundant decoding
- Yamamoto et. al. propose a solution: **precompute** the start points for each thread at encoding time, put them in a gap array [13], and use them for fast decoding



A gap array: {0, 0, -2, -1}



Implementing and Optimizing Decoders

> Adaptation

- Change from single-byte input to multi-byte input
- Thread-Level Optimization (for self-sync)
- With self-synchronization, adjacent threads may decode very different amounts of data-**divergence**
- Program with the thread hierarchy in mind

> Memory Optimization (for <u>both</u>)

- Use wider/vector loads and stores
- Use the GPU's **shared memory** to cache decoded results



CUDA C programming guide, [6]



Motivation for Memory Optimizations

> High-compression ratio data

- Often found in scientific computing/cuSZ workflows, especially where the data has been well-predicted
- Significant performance penalties for increased compression ratio ≈ decreased error-bound

Reason

- With high compression ratios, each thread writes more data
- Also, there is a larger stride between threads, an even worse access pattern





Shared Memory Optimization Details

The technique

- Each thread writes into the block-local shared memory
- The shared memory is cooperatively written out to global memory

> Allocating shared memory

- Proportional to the compression ratio of the data
- Different portions of the data need different amounts of shared memory
- Use **multiple kernel launches** to efficiently decompress different portions of the data

Algorithm 1: Decoding and writing using a shared memory buffer.

DecodeWrite — decode and write using shared memory

```
\triangleright The shared memory buffer of size n
 1 sharedBuffer[n]
 2 si <- outIndex[blockIdx.x · blockDim.x]</pre>
 3 ei <- outIndex[(blockIdx.x + 1) · blockDim.x]</pre>
 4 gid <- threadIdx.x + blockDim.x · threadIdx.x
 5 tempEnd <- ei</pre>
 6 while si < ei do
      start <- outIndex[gid] - si, end <- outIndex[gid + 1]</pre>
 8
      if si < start and end < si + n then
 9
          outArray[start ... end) <- Decode(inArray, startPoint[gid])</pre>
                                       ▷ If symbols can fit into the buffer, decode them
10
      else if start < si + n and end > si + n then
11
          tempEnd <- outIndex[gid]</pre>
        ▷ Executed by one thread if buffer is not large enough; results in another iteration
12
      end if
      outArray[si ... tempEnd) = sharedBuffer[0 ... tempEnd - si)
13
                          > This write is performed cooperatively by threads in the block
      si <- tempEnd</pre>
14
15 end while
```



Evaluation

Experimental Setup

- **Datasets**: Multidimentional data from a variety of scientific domains; data sources include the Scientific Data Reduction Benchmark [15], in addition to some other sources
- Platform: 2 Xeon Gold 6428 "Cascade Lake" CPUs, 20 cores; 8 Nvidia Tesla V100-32GB SXM2 GPUs (only 1 GPU was used for evaluation)

	datum size	#field
datasets	dimensions	examples(s
cosmology HACC	1,071.75 MB $280,953,867$	6 in tota xx, v
molecular dynamics EXAALT	$\begin{array}{c} 951.73 \ \mathrm{MB} \\ 2338{\times}106711 \end{array}$	6 in tota dataset2.
climate CESM-ATM	$^{642.70}_{26 imes 1800 imes 3,600}$	CLDICE, RELHUM
cosmology Nyx	$512 \text{ MB} \\ 512 \times 512 \times 512$	$^{6 \text{ in tota}}_{\text{baryon density}}$
climate Hurricane	$^{381.47~\mathrm{MB}}_{4 imes100 imes500 imes500}$	$\overset{13 \text{ in tota}}{\text{CLDICE, QRAIN}}$
quantum circuits QMCPack	$^{601.52}_{115 \times 69 \times 69 \times 288}$	2 in tota einspline, einspline.pr
petroleum exploration RTM	$\substack{180.73\text{ MB}\\449\times449\times235}$	1 in total (3600 snapshots snapshot-100
quantum chemistry GAMESS	$306.19 \mathrm{MB} \\ 80.265.168$	3 in tota dddd, ffdd, ff



WASHINGTON STATE Arg

Bridges2 cluster at Pittsburg Supercomputing Center [14]

Our tested datasets



Evaluation

- Evaluation on Decoding Alone
- Outperforms coarse-grained cuSZ decoder
 - Average **2.74x** for self-synchronization
 - Average **3.64x** for gap arrays
- Predictably, gap arrays are faster than self-synchronization









Evaluation

> Evaluation on cuSZ's decompression, overall

- Considering decoding took up 83% of cuSZ's time, not surprising to see a speedup
- Average 2.08x for self-synchronization
- Average **2.43x** for gap arrays



Performance of Optimized Decoders in cuSZ Decompression



Discussion & Conclusion

> Discussion

- Self-Synchronization
 - Benefits: Huffman encoder and decoder are decoupled
 - **Drawbacks:** Redundant computations
- Gap Arrays
 - Benefits: Higher performance in decoding
 - Drawbacks: Encoder and decoder must be coupled; higher overhead on encoder

Conclusion

- Apply two algorithms for finer-grained parallel Huffman decoding to significantly speed up cuSZ's decompression
- **Optimize** these algorithms to more effectively take advantage of GPU architectures
- Future work to be done
 - Incorporate optimizations into a Huffman coding library
 - **Extend** work to other applications of Huffman decoding



Thank you!

All the questions and ideas are welcomed

Contact:

Dingwen Tao: <u>dingwen.tao@wsu.edu</u> Cody Rivera: <u>cjrivera1@crimson.ua.edu</u>









