# FRaZ: A Generic High-Fidelity Fixed-Ratio Lossy Compression Framework for Scientific Floating-point Data

**Robert Underwood** †‡, Sheng Di ‡, Jon C. Calhoun †, and Franck Cappello ‡
Questions: robertu@g.clemson.edu
Video Presentation: https://youtu.be/oXpZAEEywHg
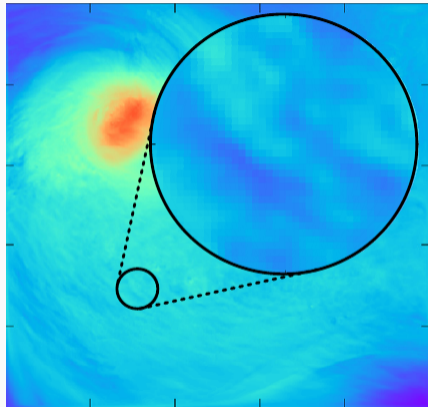Code: https://www.github.com/CODARCode/FRaZ

May 18, 2020

Clemson University †, Argonne National Laboratory ‡

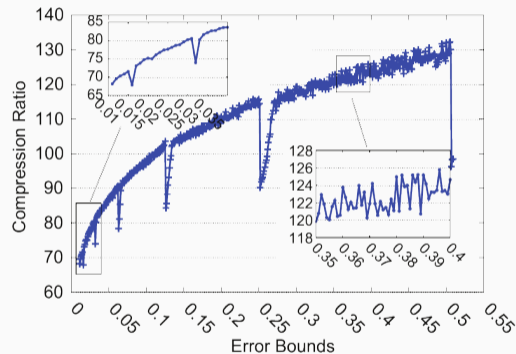# Why do we need Fixed Ratio Lossy Compression?

1. To reduce the storage footprint
   - The ORNL Summit limit: 50 TB/project
   - Many Scientific codes such as HACC or CESM produce 100s of TB if not PB of data
2. To achieve "best fit" compression
   - Users want to store as they can in their available storage
   - Without fixed-ratio, they either suffer a loss in quality or result to trial and error
3. Streaming applications
   - Scientific instruments such as the APS and LCLS-II may generate image data rates exceeding 250GB/s.
   - However, the backing storage is limited to 25GB/s



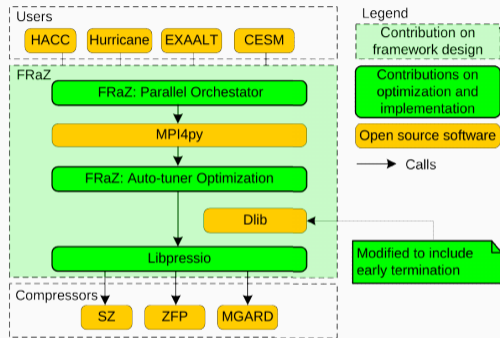Hurricane, dataset used in paper with zoom-in view

- Current compressors don't implement fixed-ratio compression or implement an similar "fixed-rate" mode which isn't error bounded (see paper)

- The relationship between error bound and compression ratio is not monotonic and non-convex for all compressors and datasets

- This is especially true of compressors like SZ which have a dictionary encoding stage

- White-box approaches (where the compressor is deeply known) quickly fall out of date



Non-monotonicity in the Hurricane dataset

3

- Formulated fixed-ratio compression as an optimization problem in a way that converges quickly
- Evaluated several different optimization algorithms to find one that works on all of our test cases, and then modified it to improve performance for our FRaZ
- Implemented and ran parallel search to improve the throughput of the technique



Overview of FRaZ Architecture and Contributions

- **Given:**
  Original Dataset $D_{f,t}$
  Decompressed Dataset $D'_{f,t}$
  Fixed Compression Parameters $\theta$
  Acceptable Compressor Error Bound $U$
  Real compression ratio $\rho_r(D_{f,t}, e, \theta)$
  Target compression ratio $\rho_t(D_{f,t})$
  Target compression ratio relative tolerance $\epsilon$
  Let: Compressor Error Bound $e$
- **Minimize over $e$:**
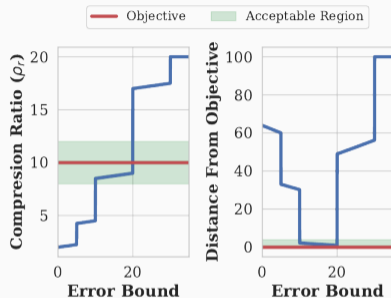  $(\rho_r(D_{f,t}, e, \theta) - \rho_t(D_{f,t}))^2$ s.t. $0 \leq e \leq U$
  if $(\rho_r(D_{f,t}, e, \theta) - \rho_t(D_{f,t}))^2 \leq \epsilon^2 \rho_t(D_{f,t})$, terminate
- **Many Algorithms preform poorly:**
  We don't have a analytic forms for $\rho_r$, $\rho_r\prime$, or $\rho_r\prime\prime$
  Numerical derivatives are costly, $O(sec) - O(min)$
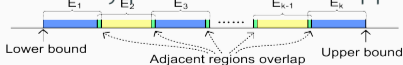  Empirically, $\rho_r$ often is non-convex many local optima



The Acceptable Region is where we can early terminate the search

- We choose Dlib's `find_global_min`
  – Lipschitz Optimization + NEWOUA,
  http://blog.dlib.net/2017/12/a-global-optimization-algorithm-worth.html

5

1. By Field – embarrassingly parallel
2. By Timestep
   - Do first timestep as normal
   - Guess next solution is same as last
   - If wrong, do full tuning again
3. By Error Bound Range
   - Split search range $[0, U]$ into $n$ similarly sized subranges run an independent search on each as hardware allows
   - a slight overlap (i.e. 10%) improves performance allowing for sufficient stationary points in the overlapping region



Lower bound    Adjacent regions overlap    Upper bound

**Algorithm 2** TRAINING

**Input**: target compression ratio $\rho_t(D_{f,t})$, acceptable error $\epsilon$, dataset $D_t$, max allowed compression error $U$

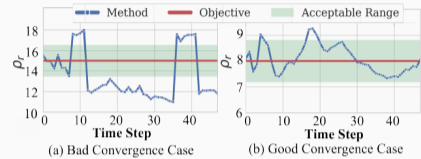**Output**: real compression ratio $\rho_r(D_{f,t}, e)$, recommended error bound setting $e$

1: $tasks[N]$
2: $done \leftarrow false$
3: **for** $(i, (l, u)) \in make\_error\_bounds(U)$ **do**
4:     $tasks[i] \leftarrow launch\_task(D_t, l, u, \rho_t(D_{f,t}), \epsilon, h)$
5: **end for**
6: **while** $not done$ **do**
7:     $last\_task \leftarrow next\_completed(tasks)$
8:     $candidate \leftarrow compression\_ratio(last\_task)$
9:     **if** $\rho_t(D_{f,t})(1 - \epsilon) \leq candidate \leq \rho_t(D_{f,t})(1 + \epsilon)$ **then**
10:         $done \leftarrow true$
11:         **for** $task \in tasks$ **do**
12:             $cancel\_if\_not\_finished(task)$
13:         **end for**
14:     **end if**
15:     $done \leftarrow has\_next(completed)$
16: **end while**
17: $\rho_r(D_{f,t}, e) = \infty$
18: **for** $task \in tasks$ **do**
19:     **if** finished(task) **then**
20:         $\rho \leftarrow compression\_ratio(task)$
21:         **if** $(\rho_r - \rho)^2 < (\rho_t - \rho)^2$ **then**
22:             $\rho_r = \rho$
23:         **end if**
24:     **end if**
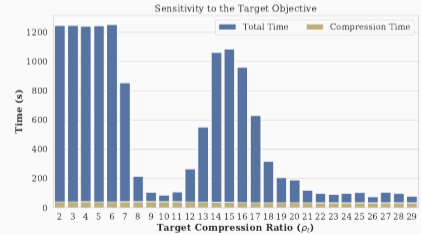25: **end for**
26: **return** $\rho_r(D_{f,t}, e), error\_bound(task)$

Worker Algorithm

6

- Runtime depends substantially if the requested target is feasible:
  - Good (feasible) Case: We terminate early most of the time
  - Bad (infeasible) Case: We alternate between a compression ratio which is too small or too large
- Very small compression ratios are often infeasible because there is a minimum compressed size
- There are also gaps between feasible and infeasible. For this figure $\rho_t(D_{f,t}) \in [14, 16]$ are infeasible for the specified $\epsilon$
- In the feasible case, overhead is often $\approx$ 2x just compressing with the correct error bound.



Solutions in good/bad case



Time to solution for many targets

7

- Fixed Ratio SZ/ZFP is generally better than ZFP Fixed Rate at each compression ratio:
  - Better Rate Distortion (higher PSNR per bit rate)
  - Higher SSIM
  - Higher PSNR
  - Better visual quality
- Figure 1: Rate Distortion for Several Datasets
- Figure 2: Visual Quality for Several Compressors
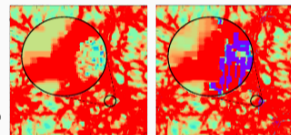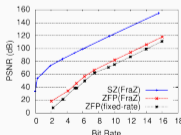


(a) Hurricane(TCf48)

(b) NYX(temperature)

(c) CESM-ATM(CLDHGH)

(d) HACC(x,y,z)

(e) EXAALT(x,y,z)

1



(a) original raw data

(b) ZFP (FraZ) (PSNR=76(c) ZFP (fixed-rate) (PSNR=56, SSIM=0.997, ACF(error)=0.516)SSIM=0.986, ACF(error)=0.383)

(d) SZ (FraZ) (PSNR=80.4(e) MGARD (FraZ) (PSNR=70, SSIM=0.999, ACF(error)=0.344)SSIM=0.977, ACF(error)=0.92)

2

# Conclusions

- Major Conclusions:
  - Fixed Ratio is better than existing Fixed Rate methods at preserving the data quality for equivalent compression ratios
  - Fixed Ratio Compression is higher performance when there are a large number of feasible compression ratios
  - We have relatively low overhead in the feasible case
- Future Work:
  - Arbitrary User Error Bounds – bounds that correspond with the quality of a scientist's analysis result relative to that on noncompressed data
  - Online Version – Develop an online version of this algorithm to provide in situ fixed-ratio compression for simulation and instrument data.
  - Algorithm Improvements – Further improve the convergence rate of our algorithm to make it applicable for more use cases

## Acknowledgements