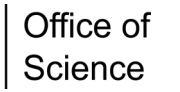


# LibPressio

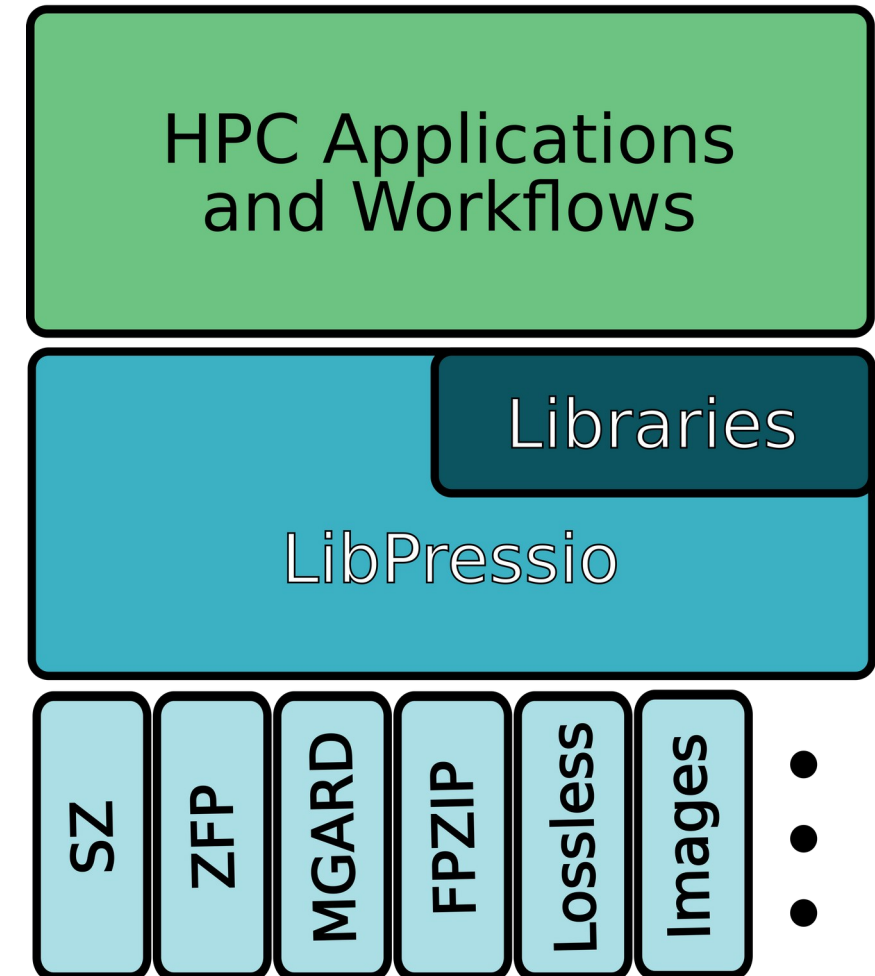


Robert Underwood – Clemson University



- Problem: Every library has its own API
  - Harder to learn and use
  - Requires rewrite to use new compressors
  - Limited collaboration and comparison
- Solution: One C/C++ API for all compressors
- Who: Application Users and Compression Developers
- Features:
  - Supports many compressors
    - SZ, ZFP, MGARD, FPZIP
    - Also: Images (i.e. JPEG, WEBP), Lossless (BLOSC)
  - Safe, Consistent, Simple, Introspectable, Fully Documented
  - Tooling Interface for Analysis

## LibPressio Ecosystem



- Easy to Install/Use (Spack, Container, Cmake)
- Extensible with:
  - Compressors – [Writing a Compressor Plugin](#)
  - Tooling Modules – [Writing a Metrics Plugin](#)
- A simple, consistent workflow
  - Get a reference to a compressor
  - Configure and assign options
  - Describe input and output buffers
  - Compress
  - Decompress
  - Release Resources
- Resources
  - [Watch the "LibPressio Tutorial" on YouTube](#)
  - [Reference the extensive developer documentation](#)

```
//get the compressor
struct pressio* library = pressio_instance();
struct pressio_compressor* sz = pressio_get_compressor(library,
↳ "sz");

//configure, validate, and assign the options
struct pressio_options* config =
↳ pressio_compressor_get_options(sz);
pressio_options_set_integer(config, "sz:error_bound_mode",
↳ REL);
pressio_options_set_double(config, "sz:rel_err_bound", 0.01);
pressio_compressor_set_options(sz, config);

//read in an input buffer
size_t dims[] = {500,500,100};
struct pressio_data* description =
↳ pressio_data_new_empty(pressio_float_dtype, 3, dims);
struct pressio_data* input_data =
↳ pressio_io_data_path_read(description, "CLOUDf48.bin.f32");

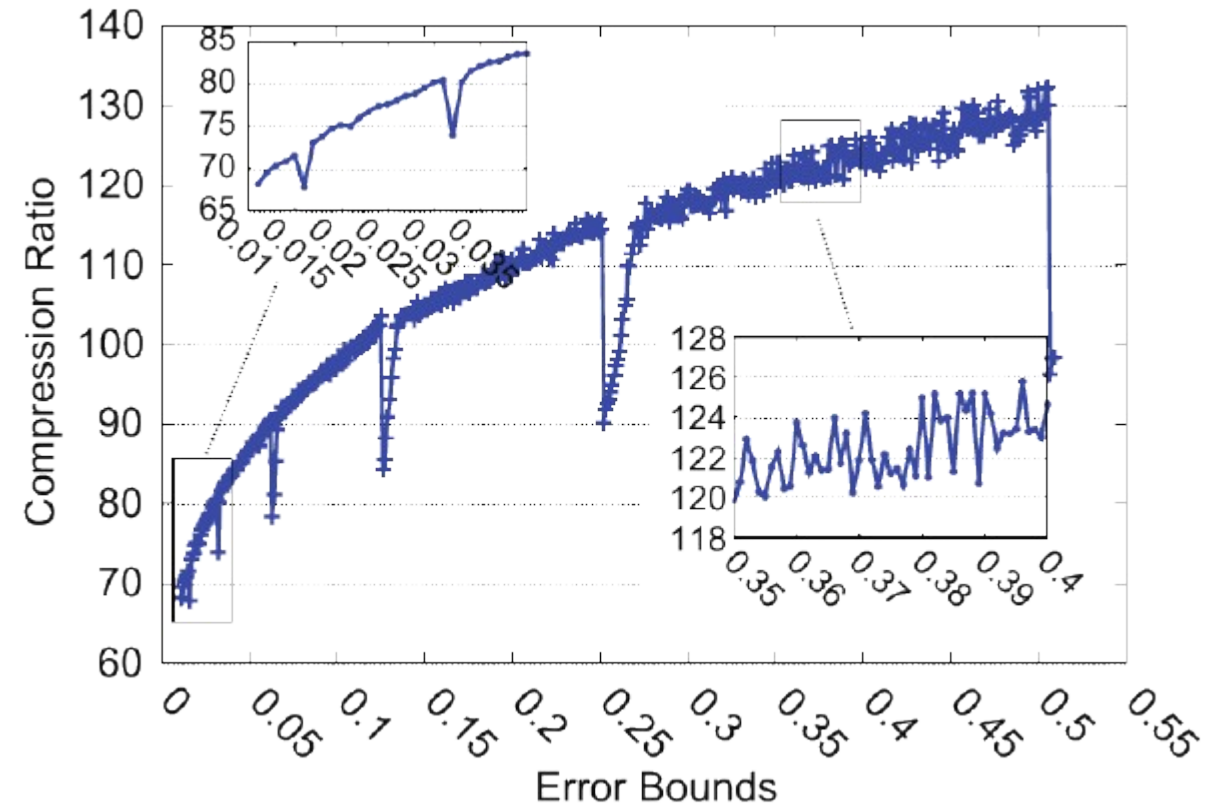
//create output buffers
struct pressio_data* compressed_data =
↳ pressio_data_new_empty(pressio_byte_dtype, 0, NULL);
struct pressio_data* decompressed_data =
↳ pressio_data_new_owing(pressio_float_dtype, 3, dims);

//compress and decompress the data
pressio_compressor_compress(sz, input_data, compressed_data);
pressio_compressor_decompress(sz, compressed_data,
↳ decompressed_data);
```

- Since last year:
- Several new and improved compressor plugins
- Better integration with IO libraries
- Improved metrics and metrics execution plugins
- Significant improvements to language bindings
- Current Uses:
  - Language Bindings (Bash, Python, Julia, R, Rust)
  - FRaZ/Opt – Autotuning Frameworks for EBLC
  - Z-Checker – Error Analysis Framework
  - Fault Injection Workflow
  - Distributed Compression Benchmarking
  - And many other research workflows!

# LibPressioOpt – Why is this difficult?

- What: Tune compressors using user's metrics
- Why is this needed:
  - Users need lossy compression:
    - to reduce storage footprint
    - to achieve “best fit” compression
    - to manage streaming volume
  - Users care how their analysis are affected
  - Many user metrics are hard to bound analytically
  - Sometimes we can improve over analytical methods
- Why is this hard?
  - The relationships between bounds and metrics are complex



# Formulating the Optimization Problem

- **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.  $0 \leq e \leq U$

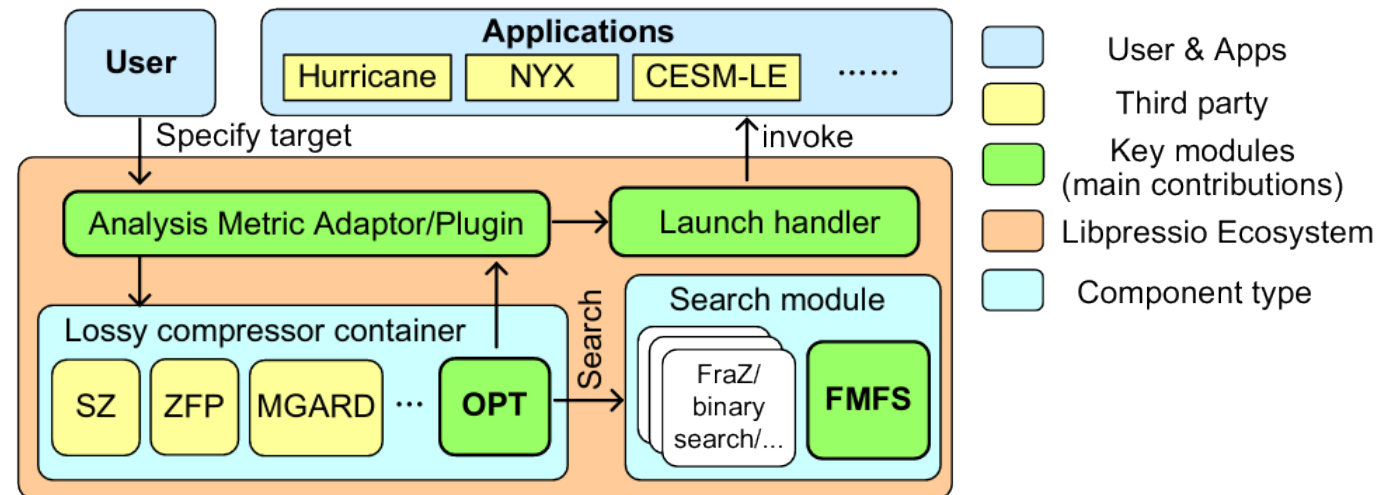
if  $(\rho_r(D_{f,t}, e, \theta) - \rho_t(D_{f,t}))^2 \leq \epsilon^2 \rho_t(D_{f,t})$ , terminate

- **Many Algorithms preform poorly:**

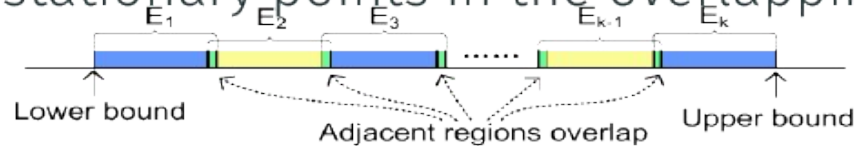
We don't have a analytic forms for  $\rho_r$ ,  $\rho_r'$ , or  $\rho_r''$

Numerical derivatives are costly,  $O(sec) - O(min)$

Empirically,  $\rho_r$  often is non-convex many local optima



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



## Algorithm 2 TRAINING

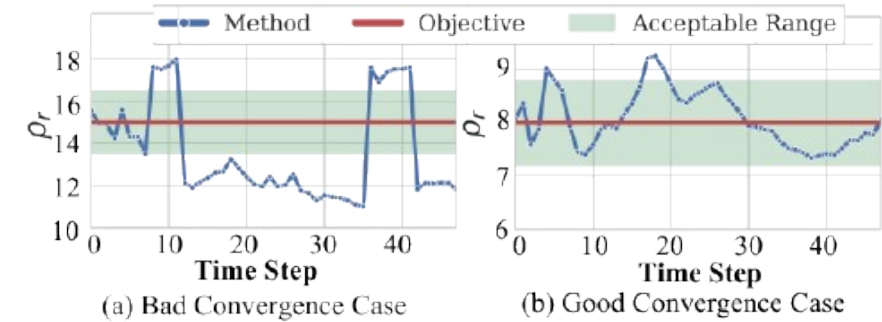
**Input:** target compression ratio  $\rho_t(D_{f,t})$ , acceptable error  $\epsilon$ , dataset  $D_t$ , max allowed compression error  $U$

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

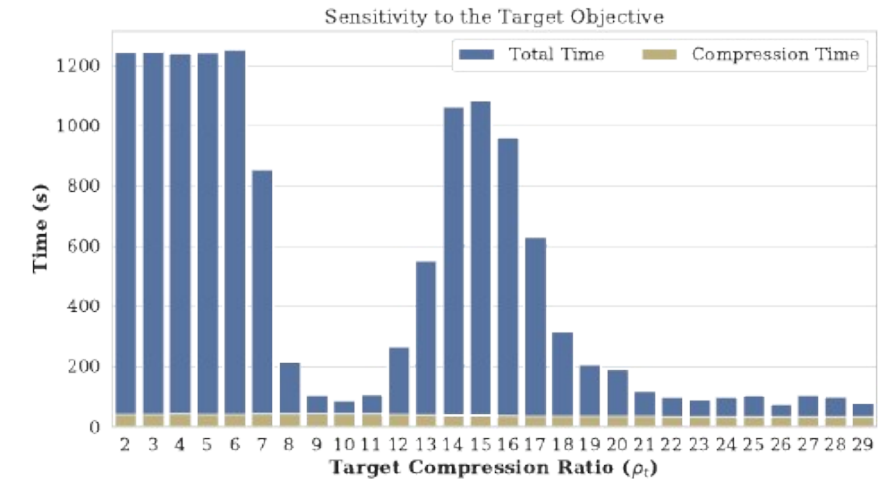
```
1: tasks[N]
2: done ← false
3: for (i, (l, u)) ∈ make_error_bounds(U) do
4:   tasks[i] ← launch_task(D_t, l, u, ρ_t(D_{f,t}), ε, h)
5: end for
6: while notdone do
7:   last_task ← next_completed(tasks)
8:   candidate ← compression_ratio(last_task)
9:   if ρ_t(D_{f,t})(1 - ε) ≤ candidate ≤ ρ_t(D_{f,t})(1 + ε) then
10:    done ← true
11:    for task ∈ tasks do
12:      cancel_if_not_finished(task)
13:    end for
14:  end if
15:  done ← has_next(completed)
16: end while
17: ρ_r(D_{f,t}, e) = ∞
18: for task ∈ tasks do
19:   if finished(task) then
20:     ρ ← compression_ratio(task)
21:     if (ρ_r - ρ)² < (ρ_t - ρ)² then
22:       ρ_r = ρ
23:     end if
24:   end if
25: end for
26: return ρ_r(D_{f,t}, e), error_bound(task)
```

Worker Algorithm

- 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

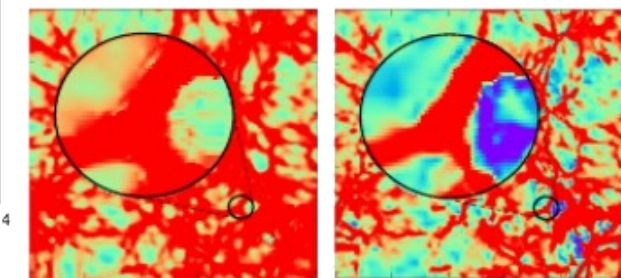
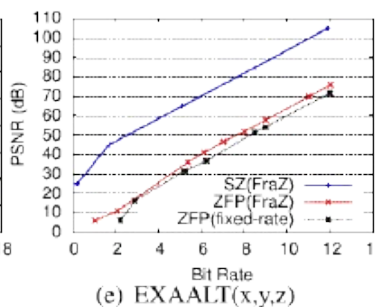
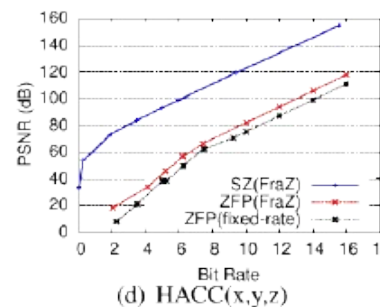
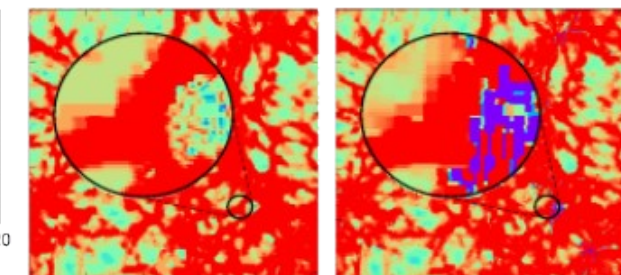
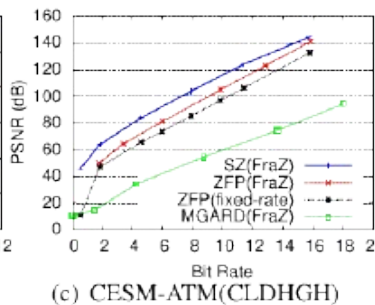
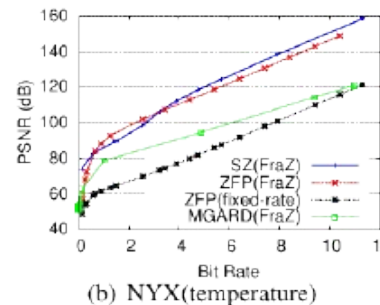
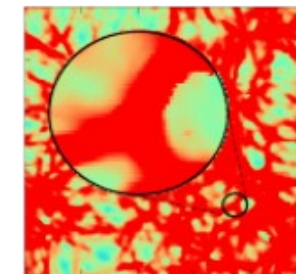
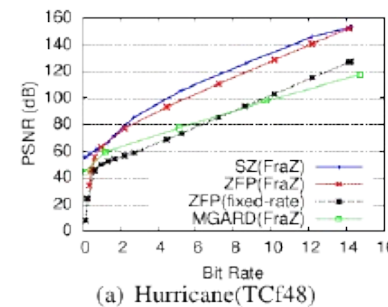


# Results: Quality of Solution

- 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



1

2

# Thanks

*This research was supported by the Exascale Computing Project (17-SC-20-SC), a joint project of the U.S. Department of Energy's Office of Science and National Nuclear Security Administration, responsible for delivering a capable exascale ecosystem, including software, applications, and hardware technology, to support the nation's exascale computing*



**RESEARCH SPONSORED BY**  
The Exascale Computing Project

A Collaborative effort of the U.S. Department of Energy, Office of Science and the National Nuclear Security Administration

(17-SC-20-SC)

