

# Compressed Floating-Point Arrays for High-Performance Computing

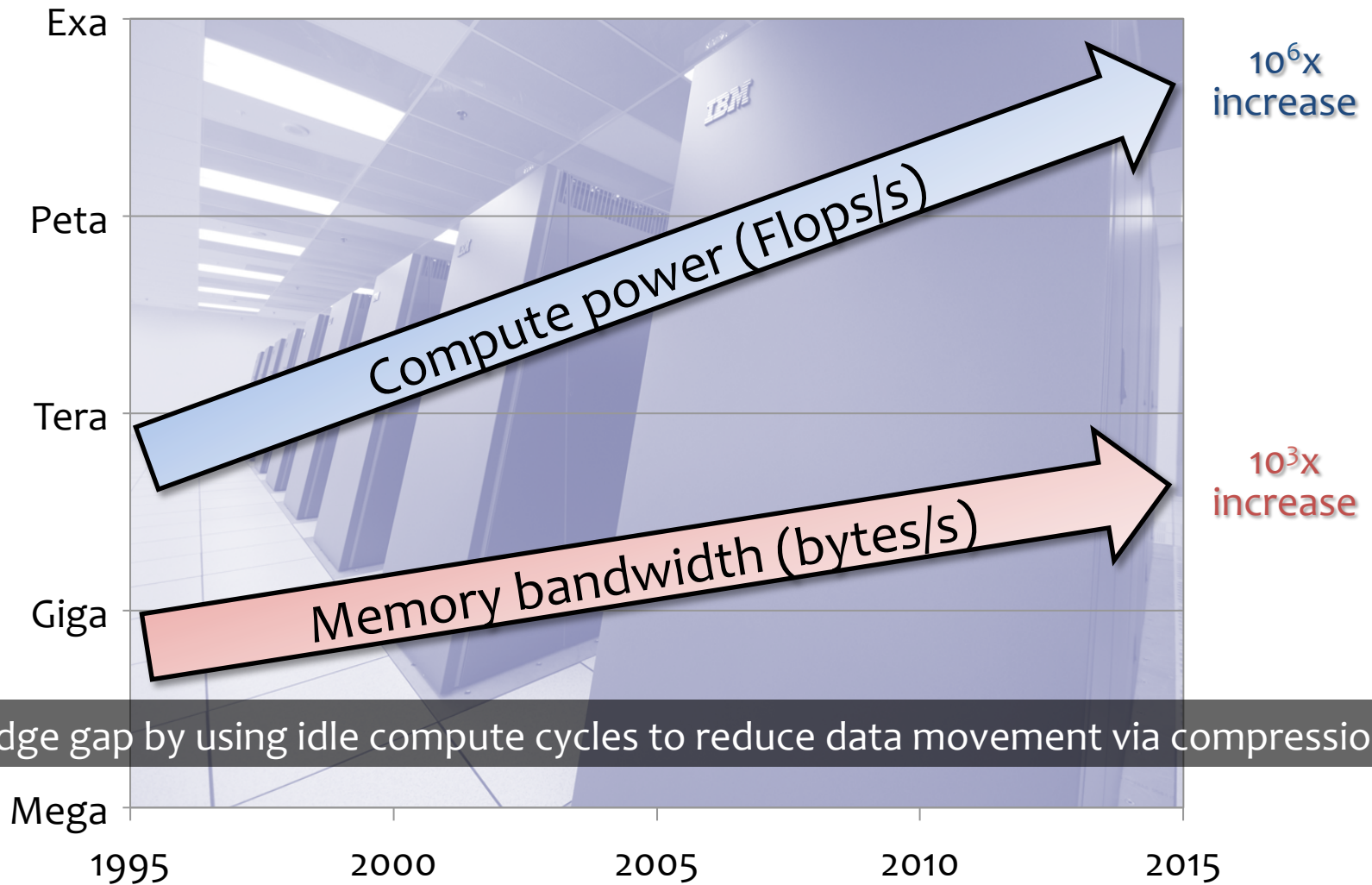
SIAM Conference on Imaging Science

Peter Lindstrom

May 24, 2016



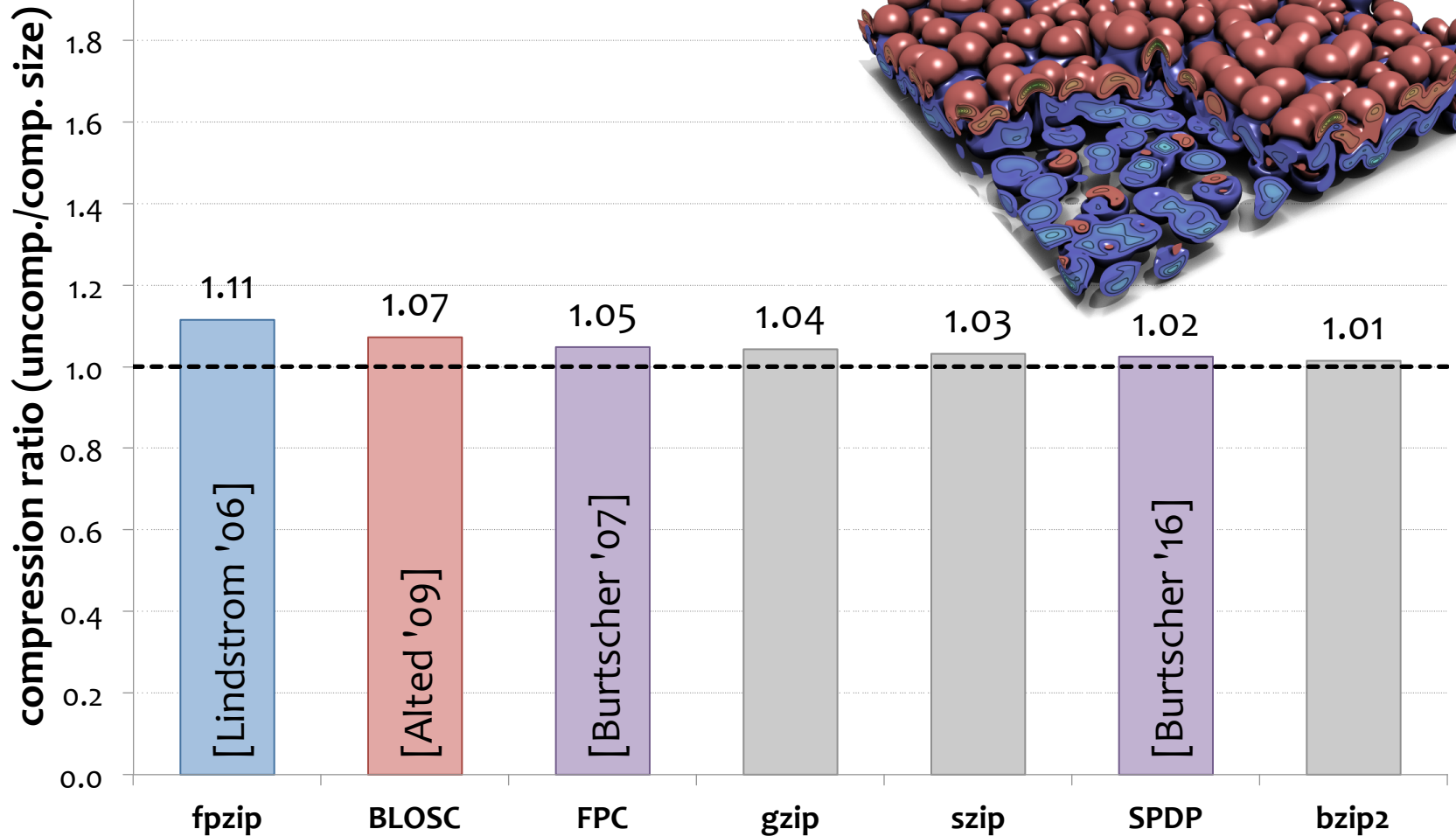
# Trend points to widening gap between compute power and memory bandwidth



Idea: Bridge gap by using idle compute cycles to reduce data movement via compression

# Floating-point data is difficult to compress— lossless compression is often not sufficient

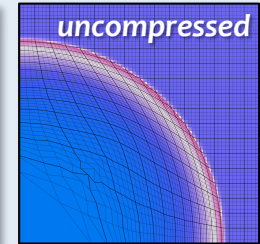
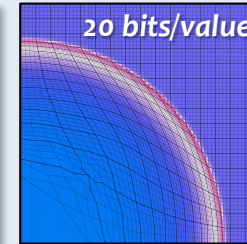
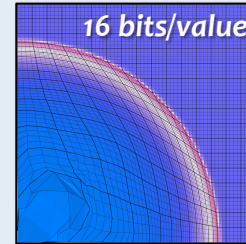
2.0 Lossy compression does *far* better, but is often met with skepticism



# We have shown that 4x lossy compression of simulation state variables can be tolerated

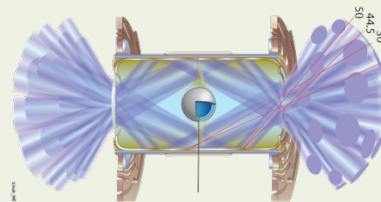
## LULESH: Lagrangian shock hydrodynamics

- QoI: *radial shock position*
- 25 state variables compressed over 2,100 time steps
- At **4x compression**, relative **error < 0.06%**



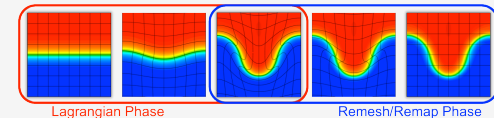
## pf3D: Laser-plasma multi-physics

- QoI: *backscattered laser energy*
- At **4x compression**, relative **error < 0.1%**



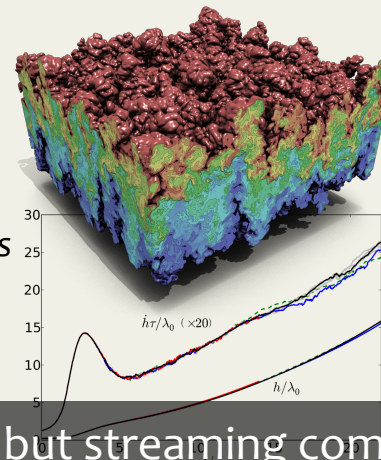
## MFEM: Cubic finite elements

- QoI: *function approximation*
- **6x compression** with ZFP **error < 0.7%** relative to FEM error



## Miranda: High-order Eulerian hydrodynamics

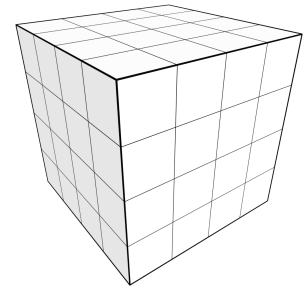
- QoI: *Rayleigh-Taylor mixing layer thickness*
- 10,000 time steps
- At **4x compression**, relative **error < 0.2%**



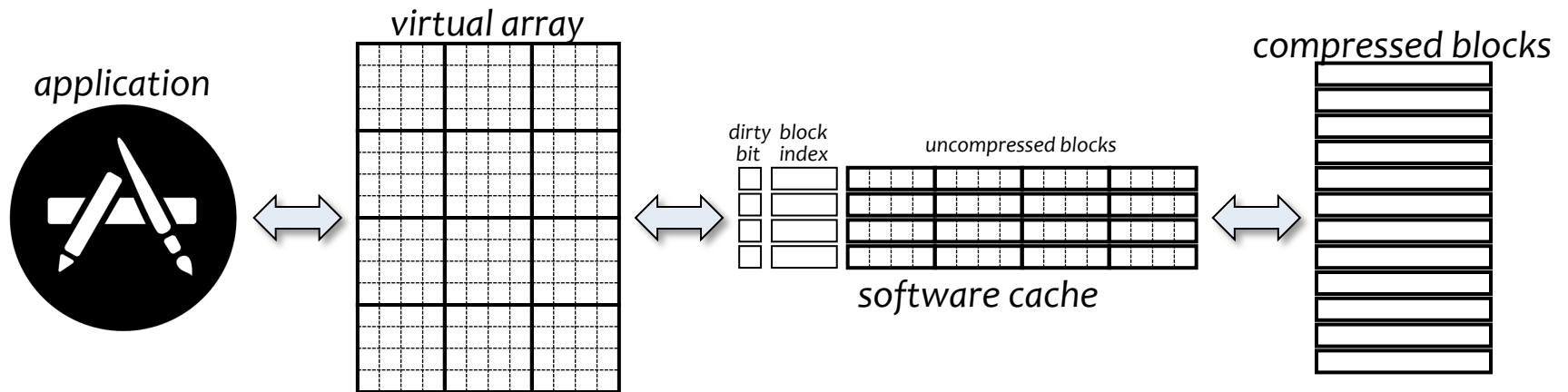
Lossy compression of state is viable, but streaming compression increases data movement

# We have developed zFP: the first inline compressor for floating-point arrays

- Inspired by ideas from h/w texture compression
  - 1D, 2D, or 3D array divided into fixed-size 4×4×4 blocks
  - Each block is independently (de)compressed
    - e.g. to a user-specified number of bits or quality
  - Fixed-size blocks ⇒ **random read/write access**
  - (De)compression is done inline, on demand
  - Write-back cache of uncompressed blocks limits data loss
- Compressed arrays via C++ operator overloading
  - Can be dropped into existing code by changing type declarations
  - `double a[n]` ⇔ `std::vector<double> a(n)` ⇔ `zfp::array<double> a(n, precision)`

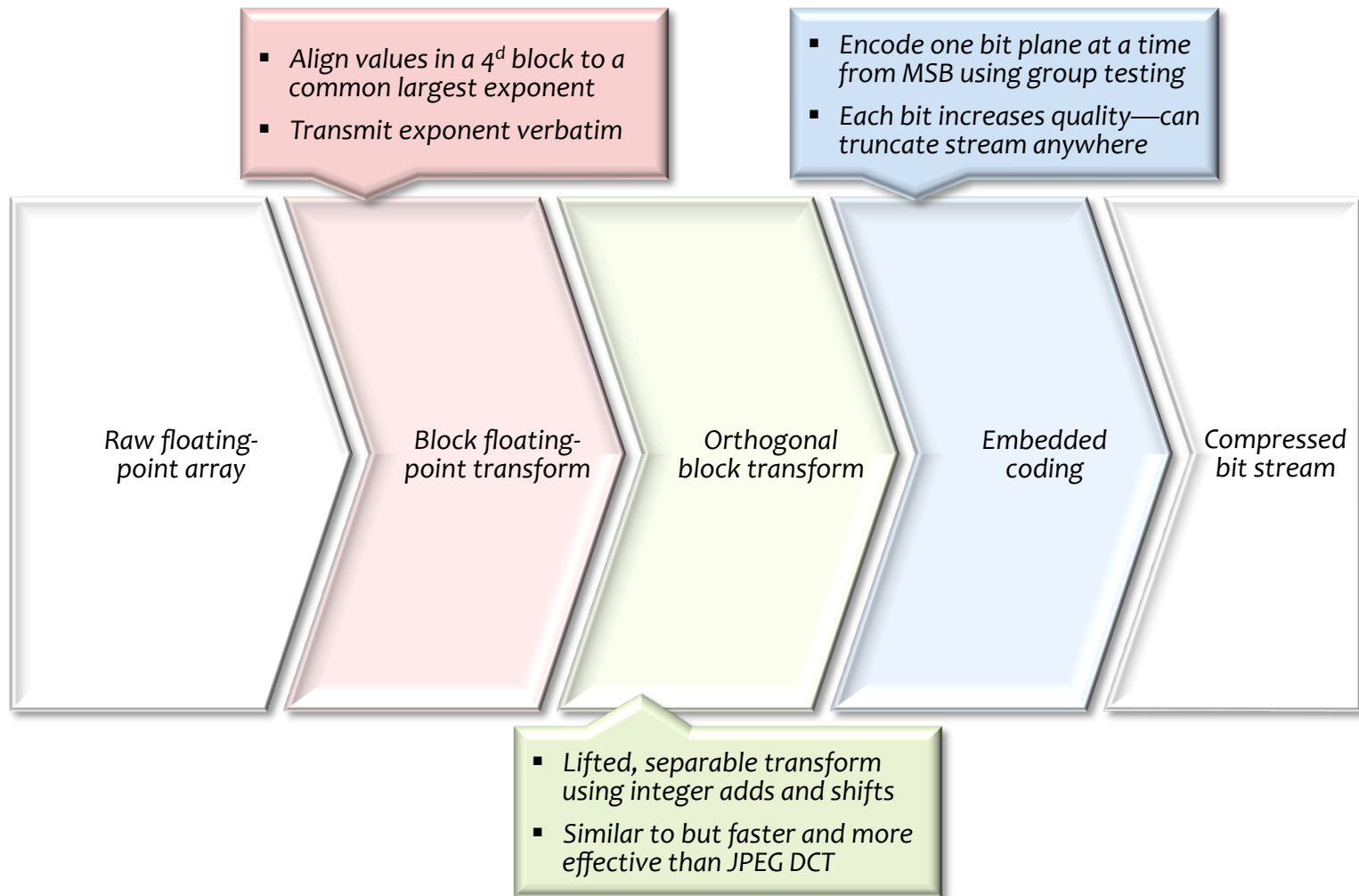


# Our compressed array limits data loss via a small write-back cache



Compress only “dirty” blocks when evicted from the cache

# Our ZFP compressor is comprised of three distinct components



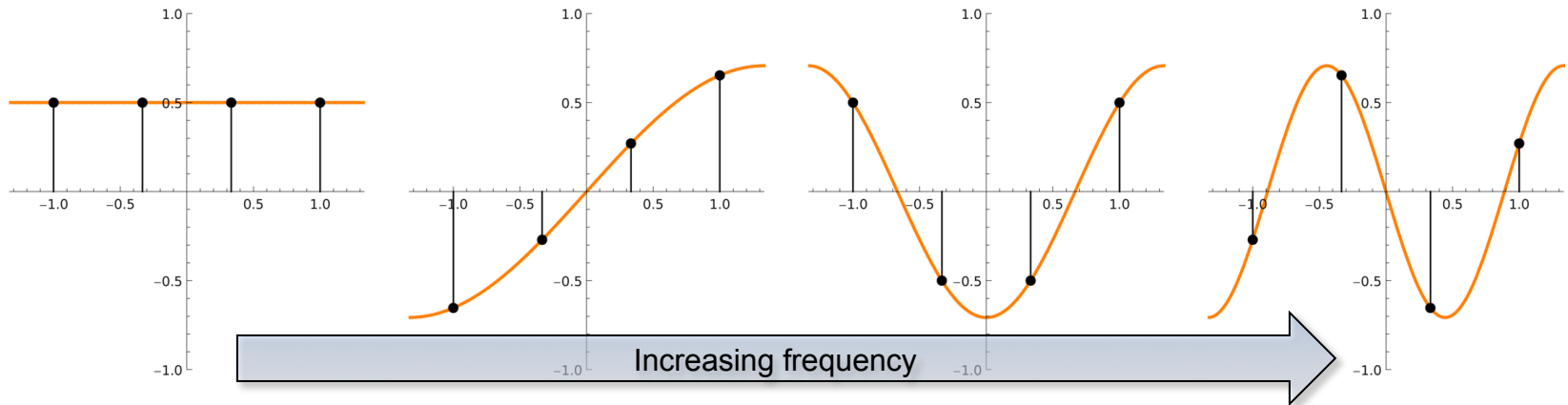
# We decorrelate values in a block using a separable orthogonal transform

$$\underbrace{\begin{pmatrix} \hat{f}_1 \\ \hat{f}_2 \\ \hat{f}_3 \\ \hat{f}_4 \end{pmatrix}}_{\text{coefficients}} = \frac{1}{2} \underbrace{\begin{pmatrix} 1 & 1 & 1 & 1 \\ c & s & -s & -c \\ 1 & -1 & -1 & 1 \\ s & -c & c & -s \end{pmatrix}}_{\text{orthogonal transform}} \underbrace{\begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{pmatrix}}_{\text{block}}$$

Free parameter  $t$

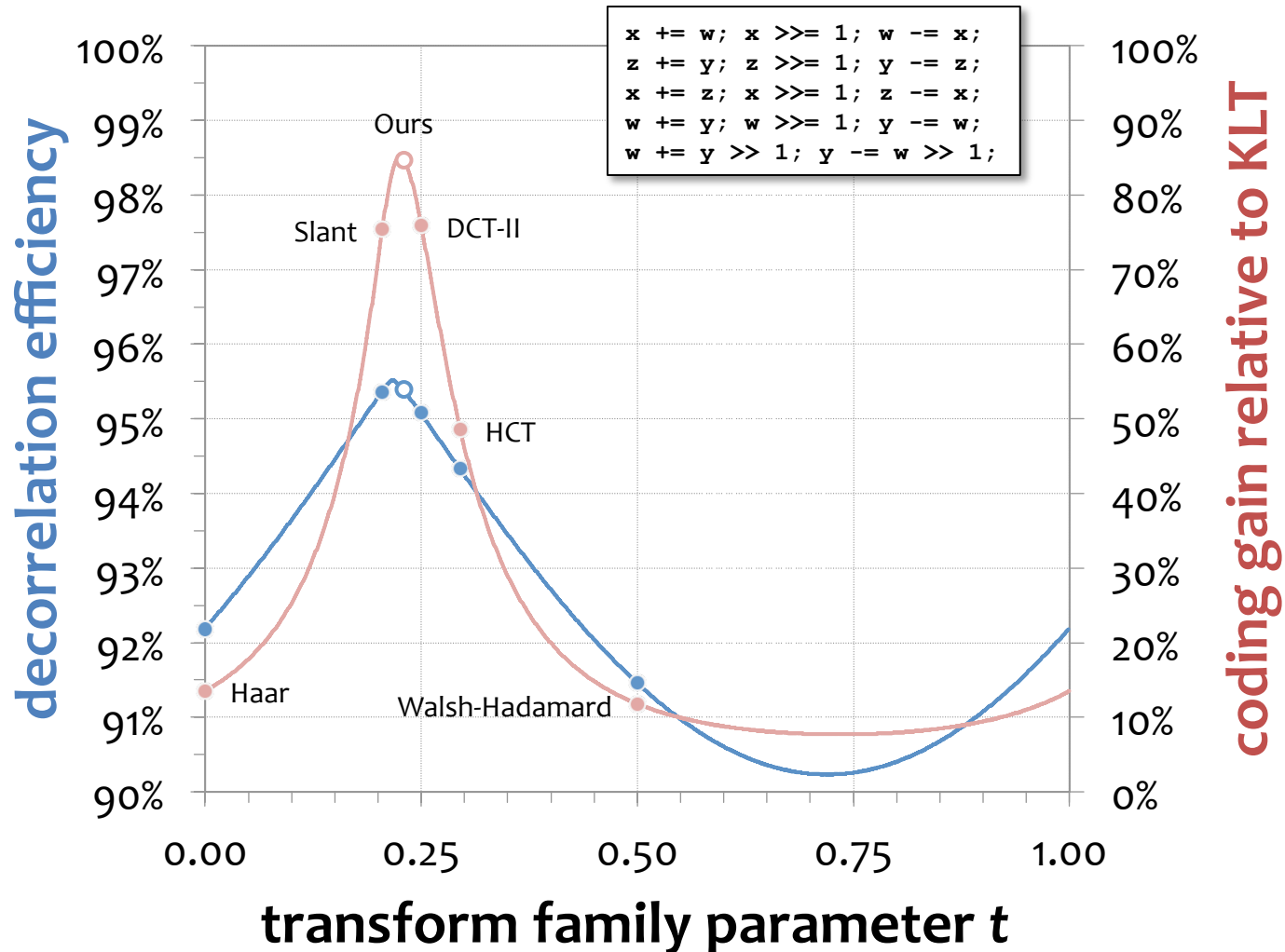
$$s = \sqrt{2} \sin \frac{\pi}{2} t \quad c = \sqrt{2} \cos \frac{\pi}{2} t$$

Basis functions  
for  $t = 1/4$



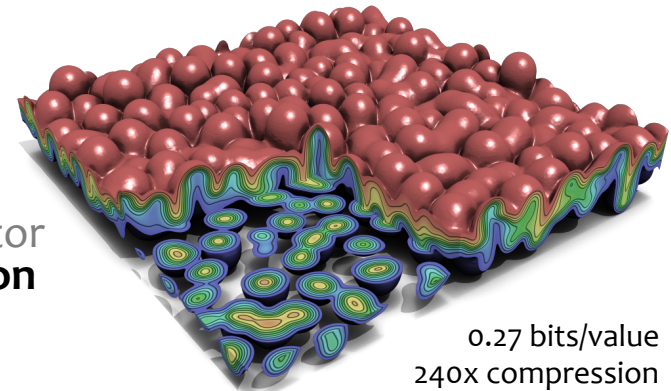


# Our integer transform is efficient, effective, and well-suited for h/w implementation

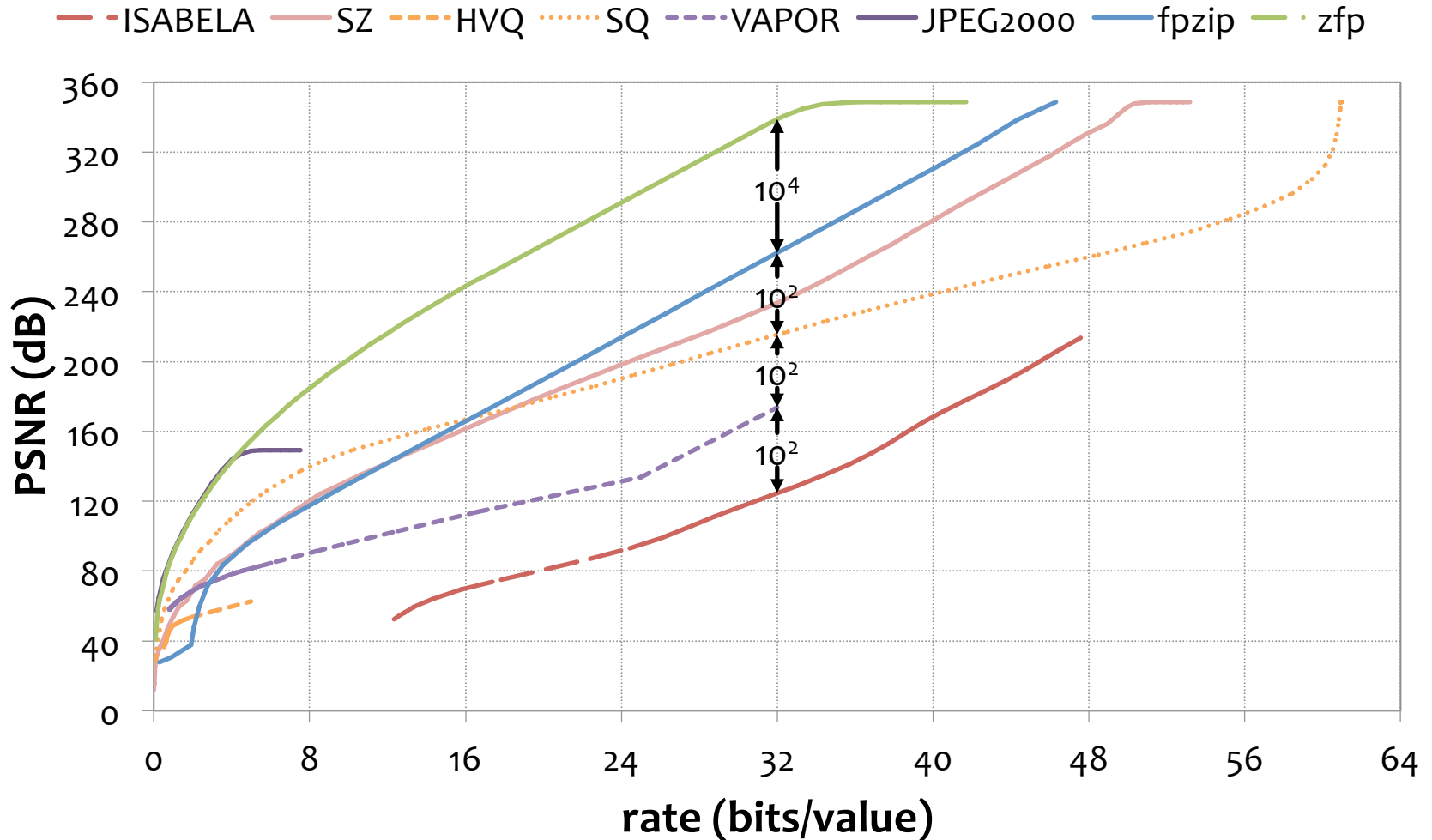


# ZFP is flexible and highly performant

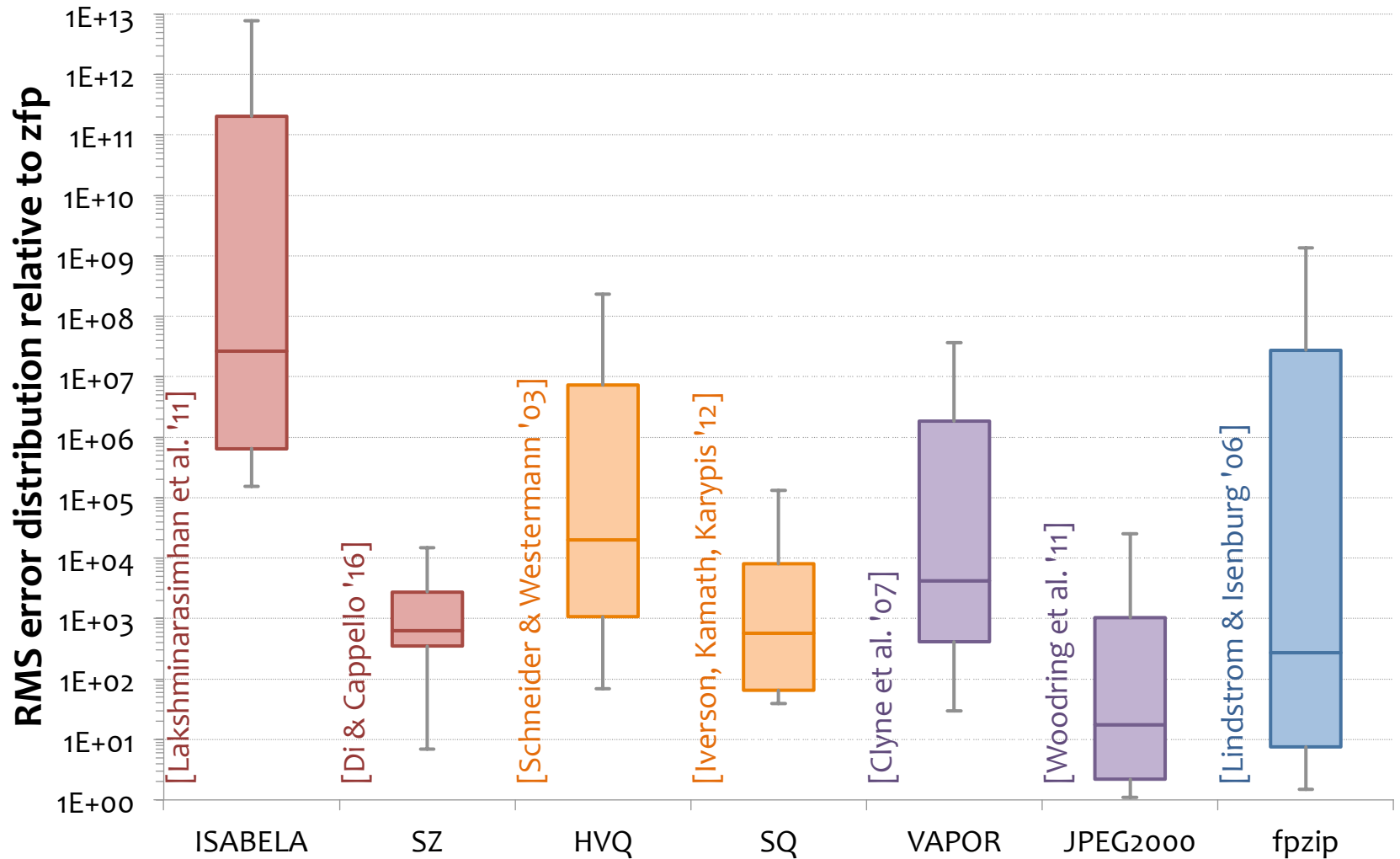
- Three compression modes in single CODEC
  - **Fixed rate:** Only inline (de)compressor with **read and write random access**
  - **Fixed precision:** Fixed no. mantissa bits ensures **relative error bound**
  - **Fixed accuracy:** User tolerance ensures **absolute error bound**
  - Also supports **progressive decompression**
- Very high quality across many science domains
  - **>100x more accurate** than closest competitor
- High, symmetric encoding & decoding speed
  - Up to **2 GB/s/core: 2-6x faster** than closest competitor
  - Simple algorithm is amenable to **h/w implementation**
    - Integer bit shifts, additions, bit-wise logical operations
- Small, independent blocks of compressed data
  - Fine granularity provides **adaptive quality, culling, queries, domain decomp, ...**
  - Supports streaming with a **tiny memory footprint** as small as a single block
  - **OpenMP, GPU parallelization** over (and within) blocks
  - **Resilience** to data corruption—a flipped bit affects only a single block
- C++ compressed arrays replace STL vectors with minimal code changes
  - LULESH miniapp integration took less than **20 minutes of coding effort**



# ZFP fixed-accuracy mode dramatically improves quality over previous methods

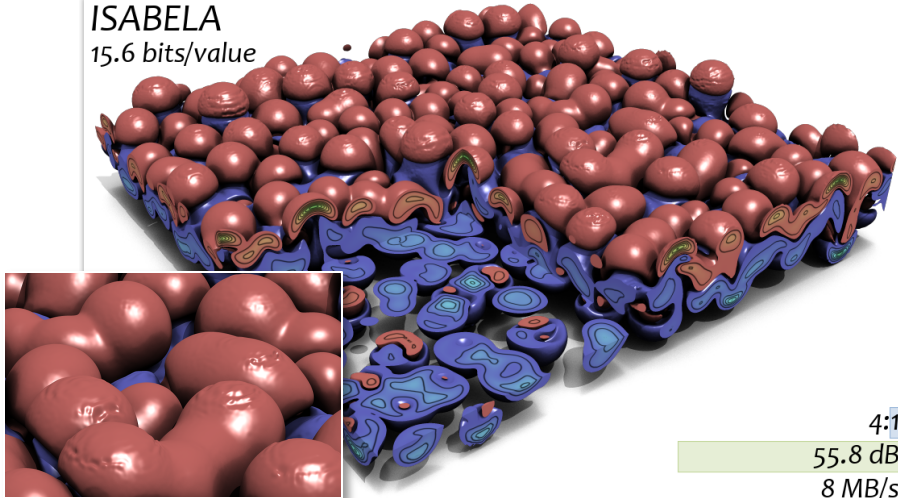


# ZFP consistently achieves the highest accuracy among compressors and data sets

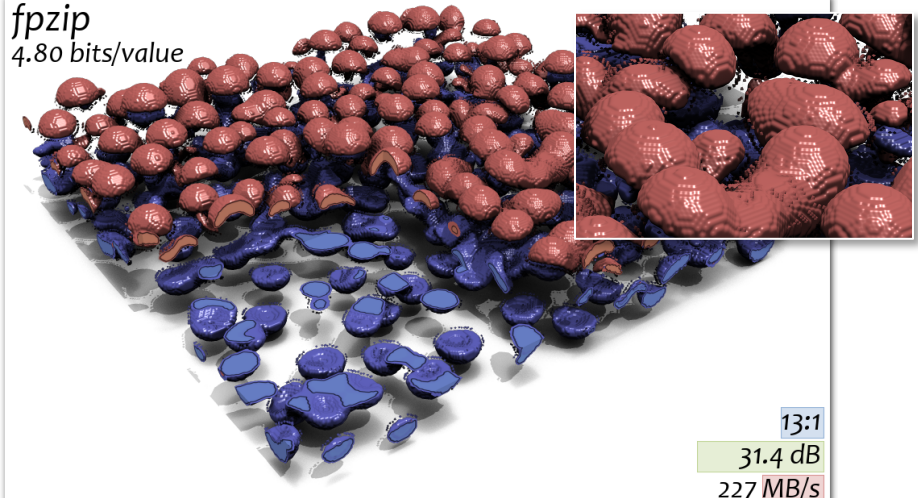


# ZFP provides >100x compression with imperceptible loss in quality

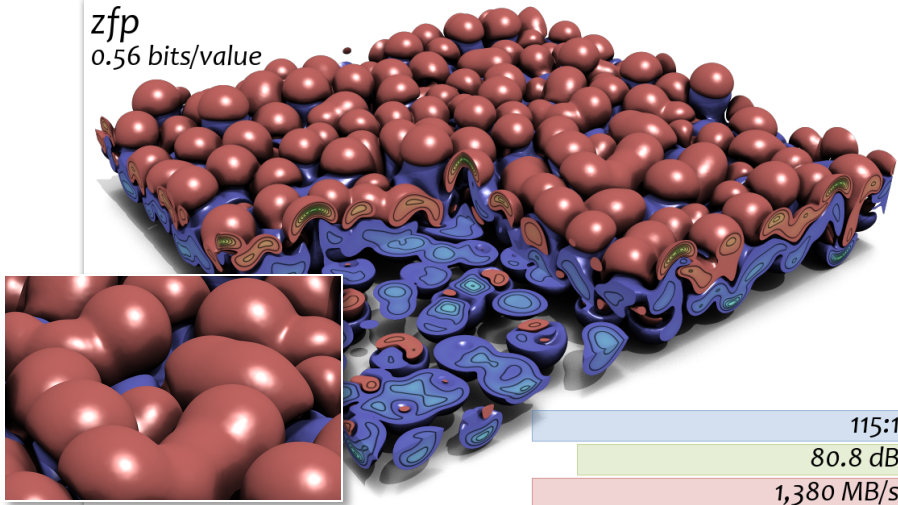
ISABELA  
15.6 bits/value



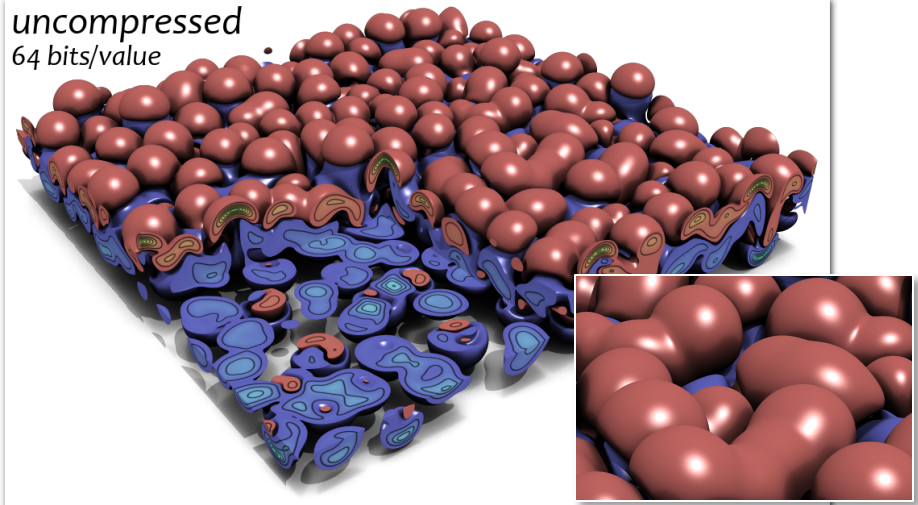
fpzip  
4.80 bits/value



zfp  
0.56 bits/value

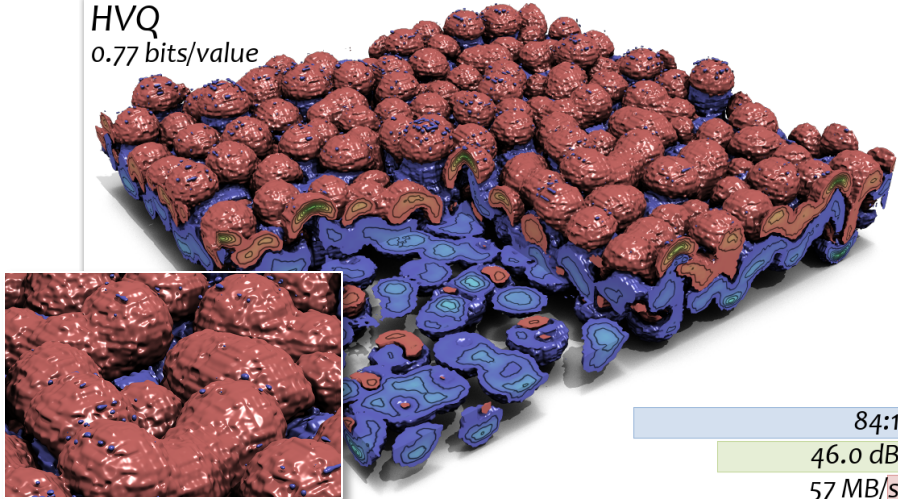


uncompressed  
64 bits/value

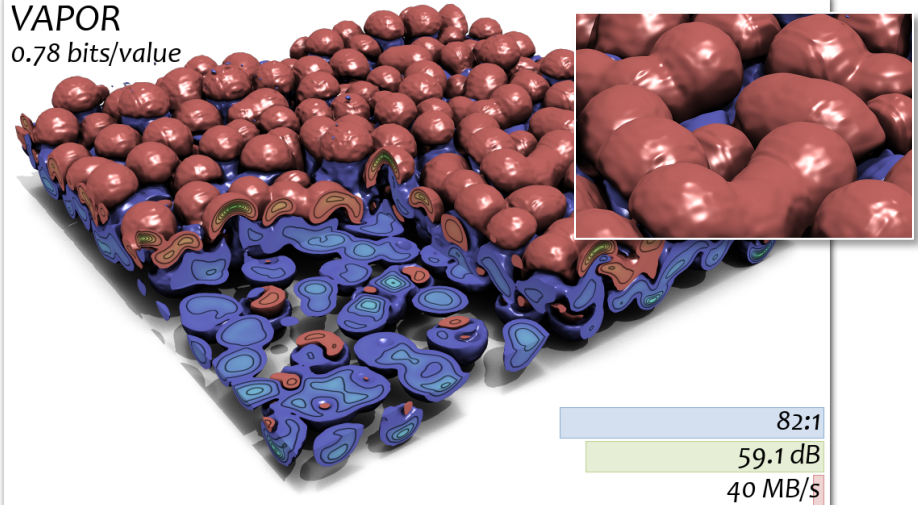


# ZFP provides >100x compression with imperceptible loss in quality

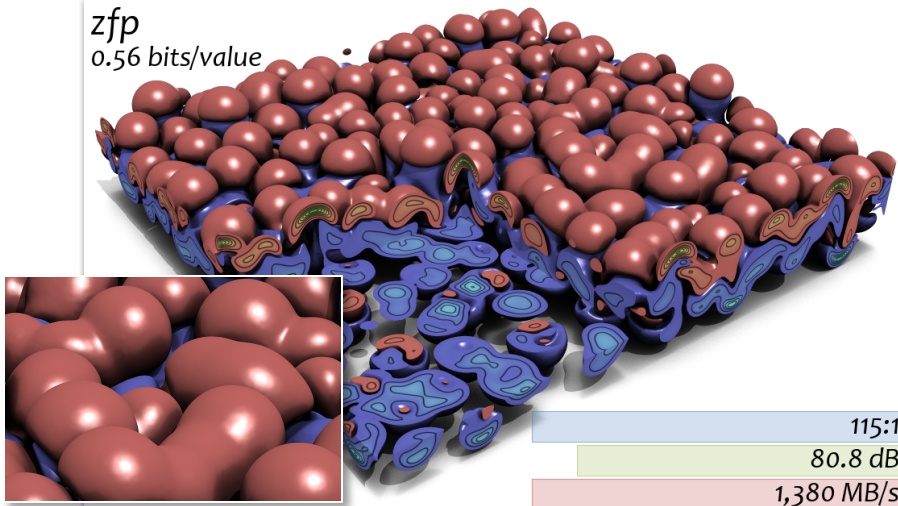
HVQ  
0.77 bits/value



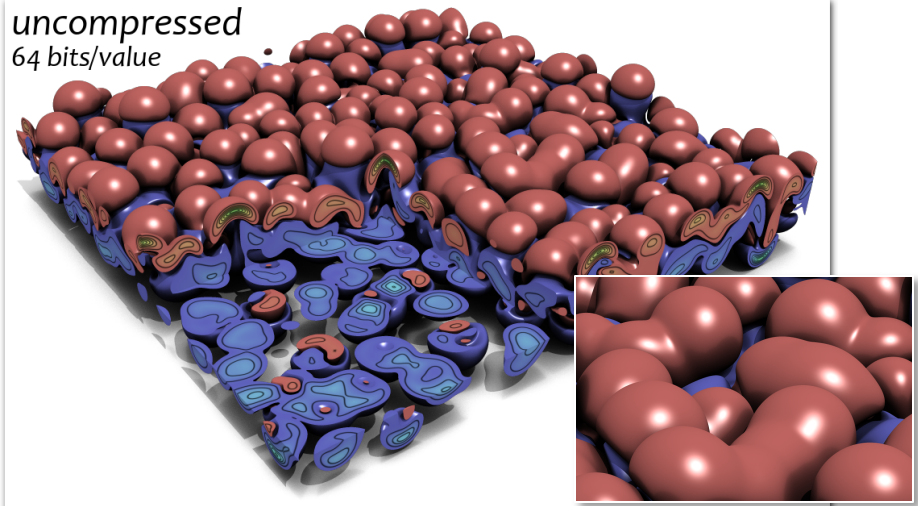
VAPOR  
0.78 bits/value



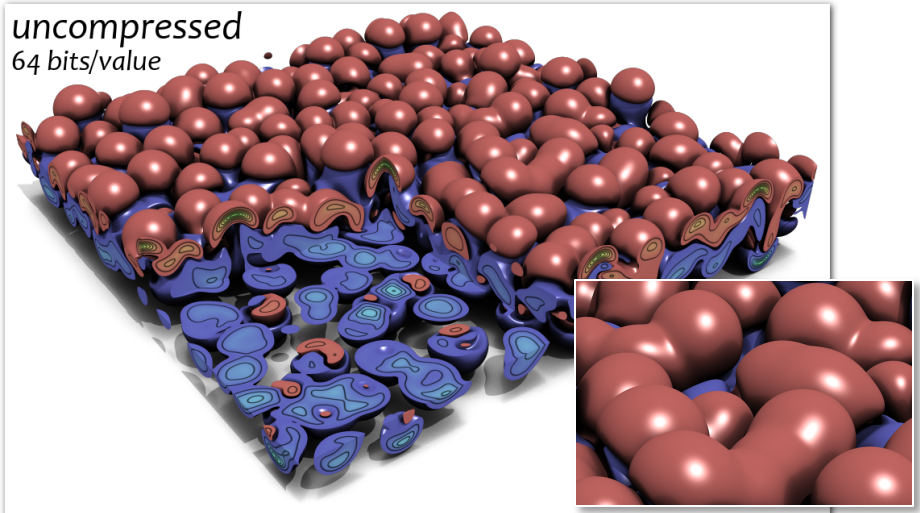
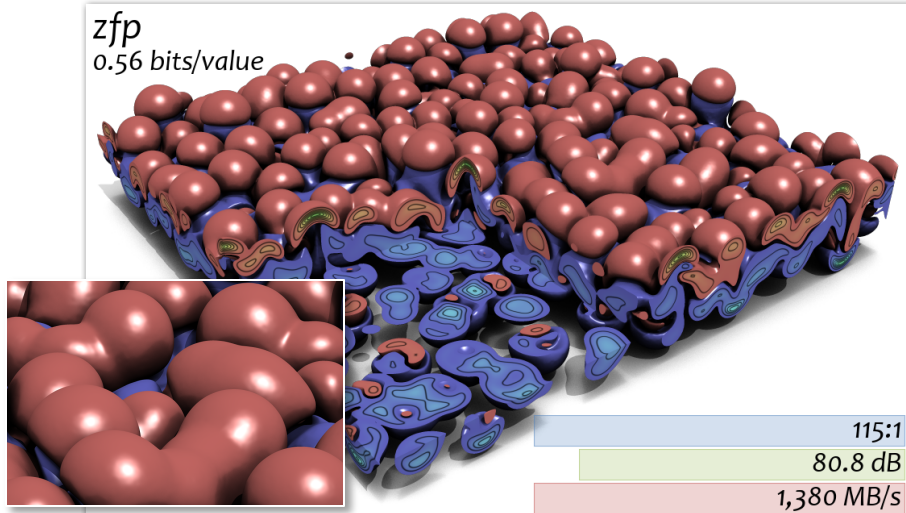
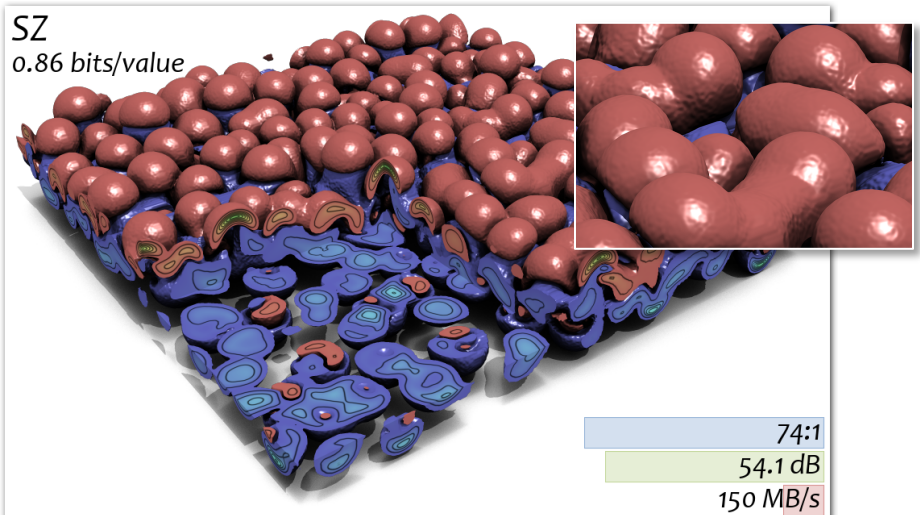
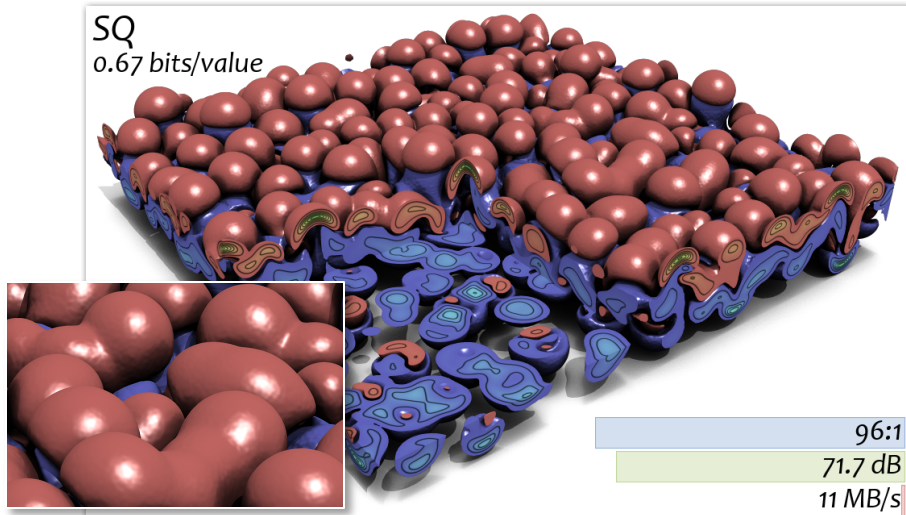
zfp  
0.56 bits/value



uncompressed  
64 bits/value



# ZFP provides >100x compression with imperceptible loss in quality

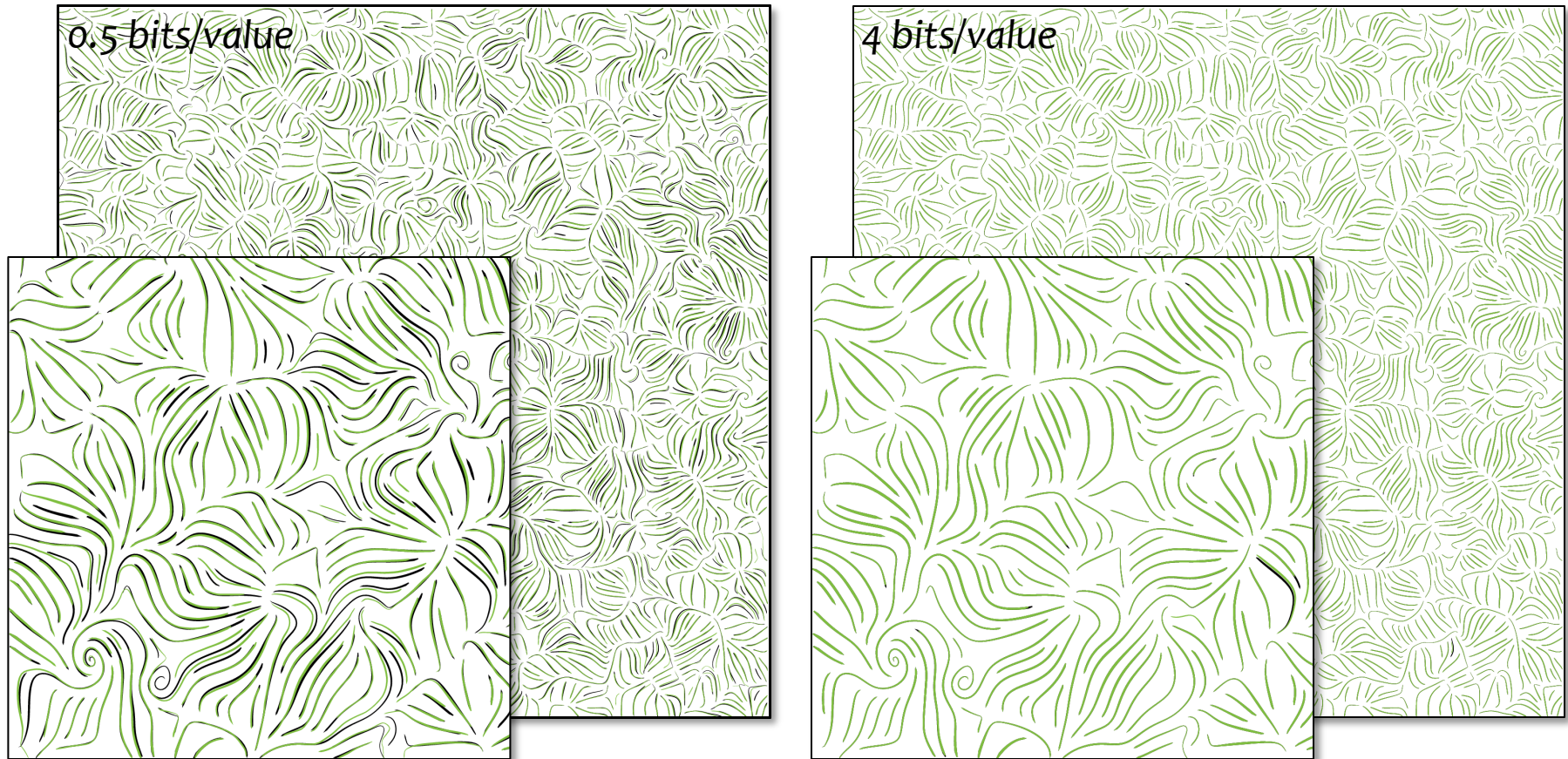


# Got artifacts?

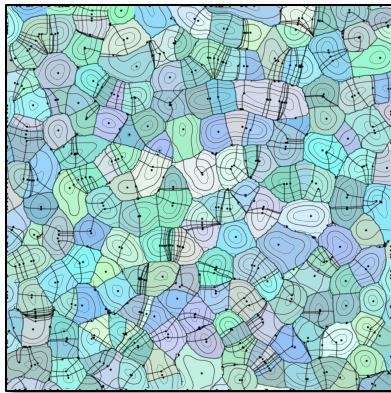




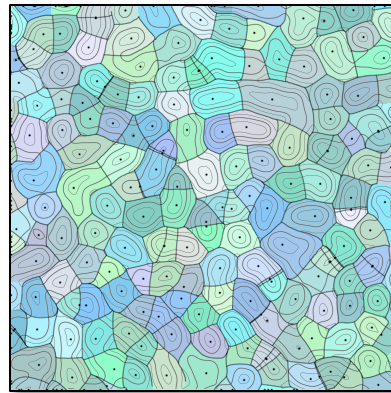
# Velocity field Runge-Kutta integration shows good agreement with uncompressed field



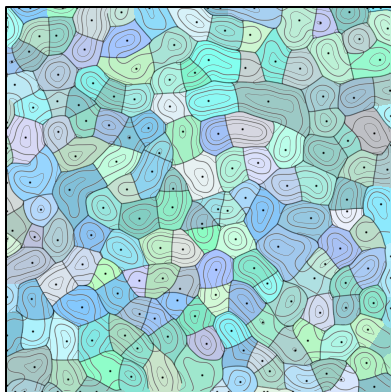
# Morse segmentation at 16x compression shows lack of blocking artifacts



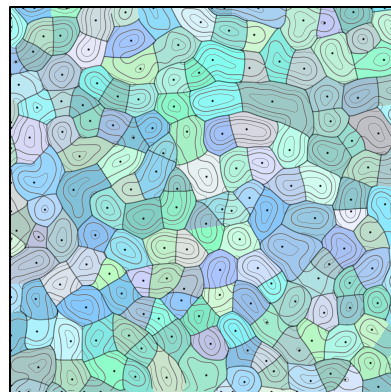
1 bit/value



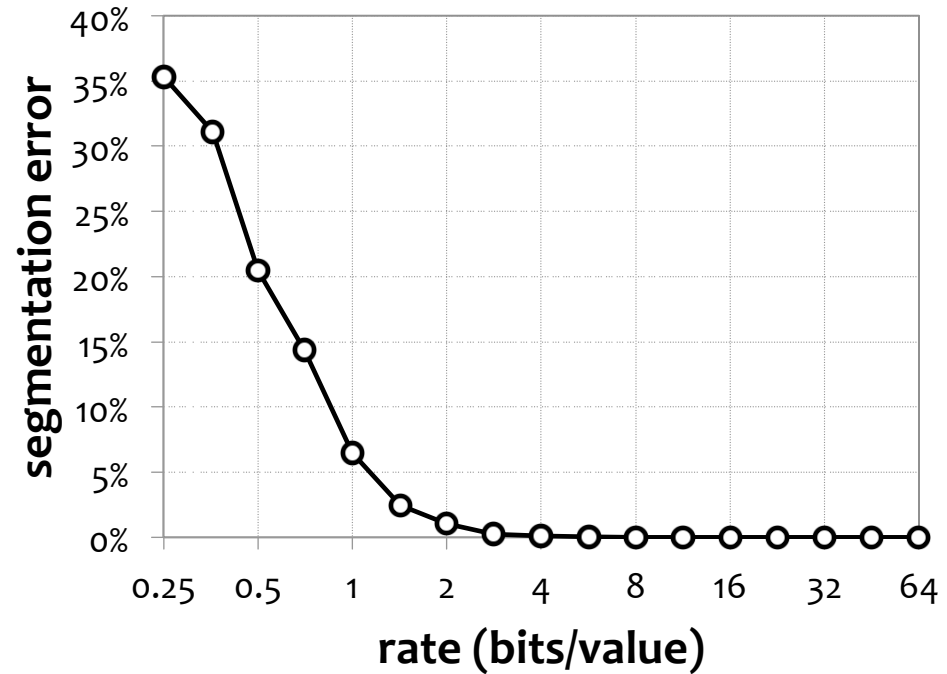
2 bits/value



4 bits/value

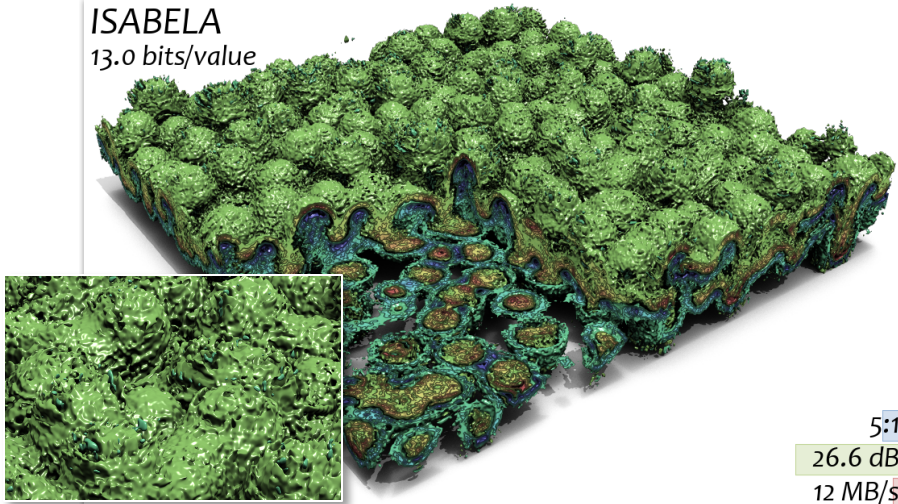


64 bits/value

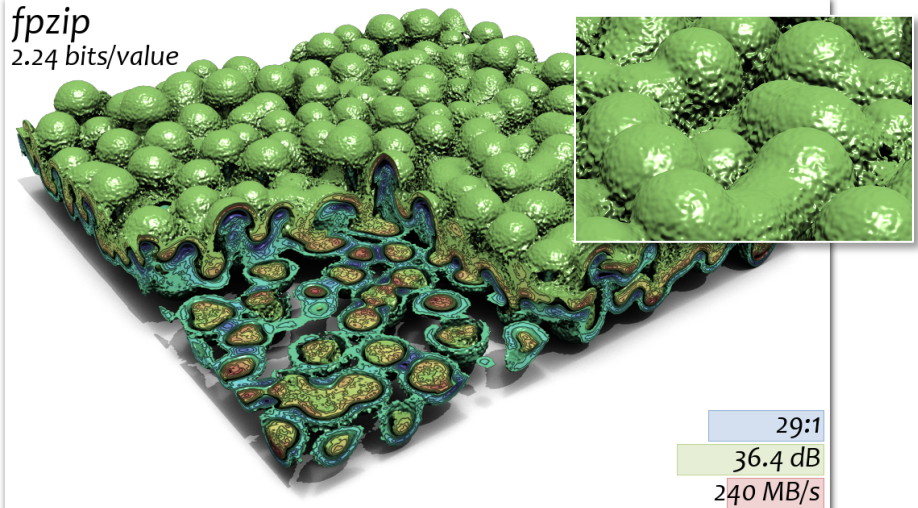


# ZFP shows no artifacts in derivative computations (velocity divergence)

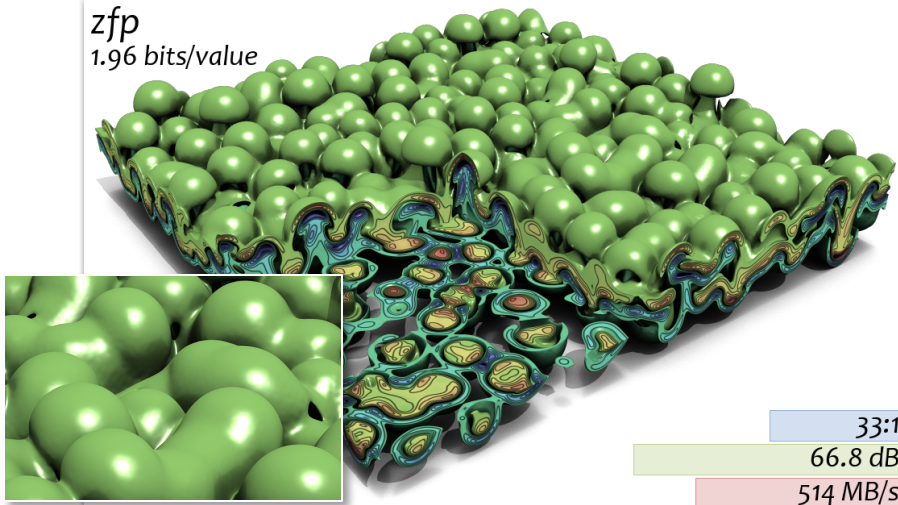
ISABELA  
13.0 bits/value



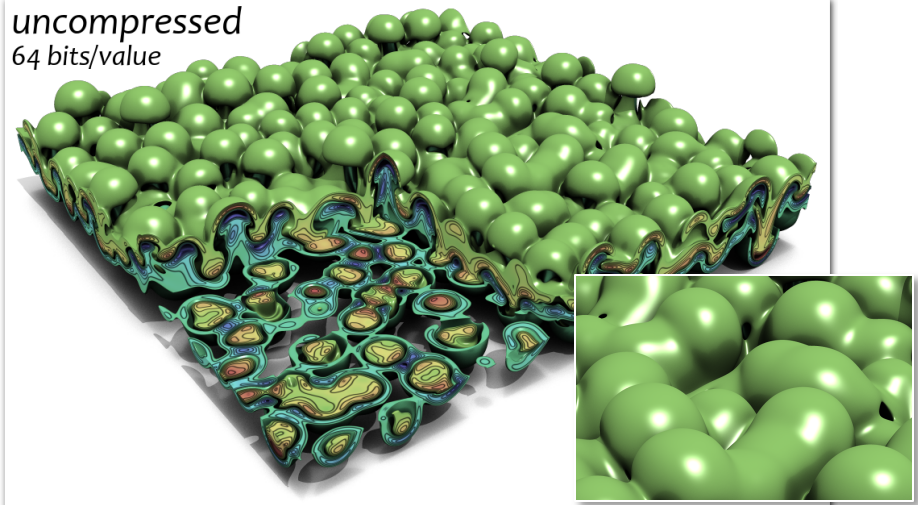
fpzip  
2.24 bits/value



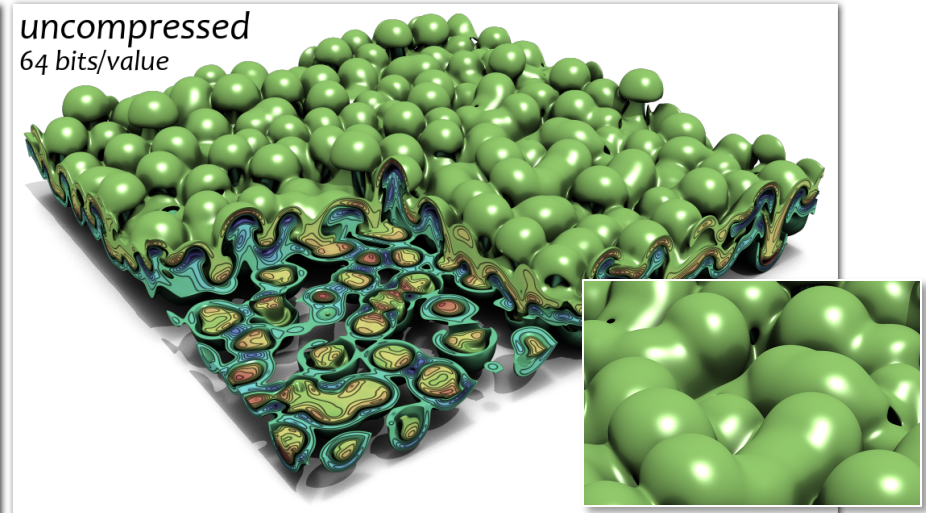
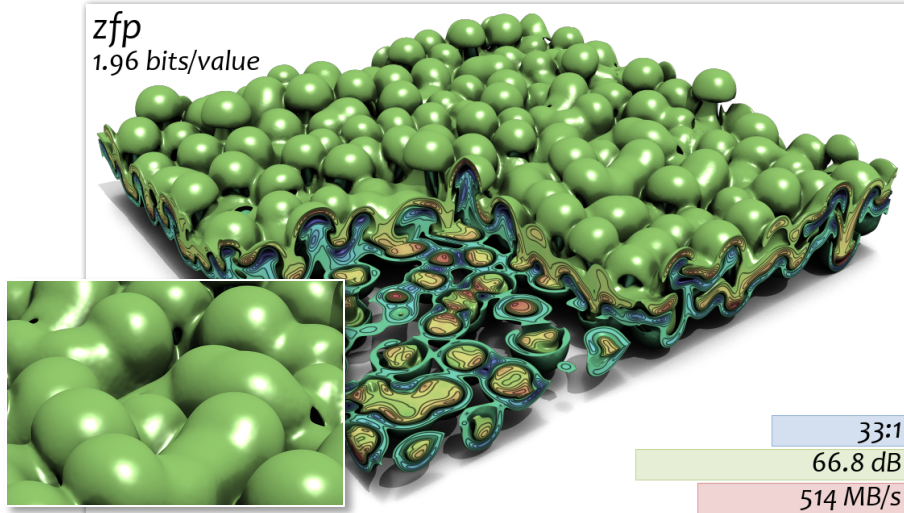
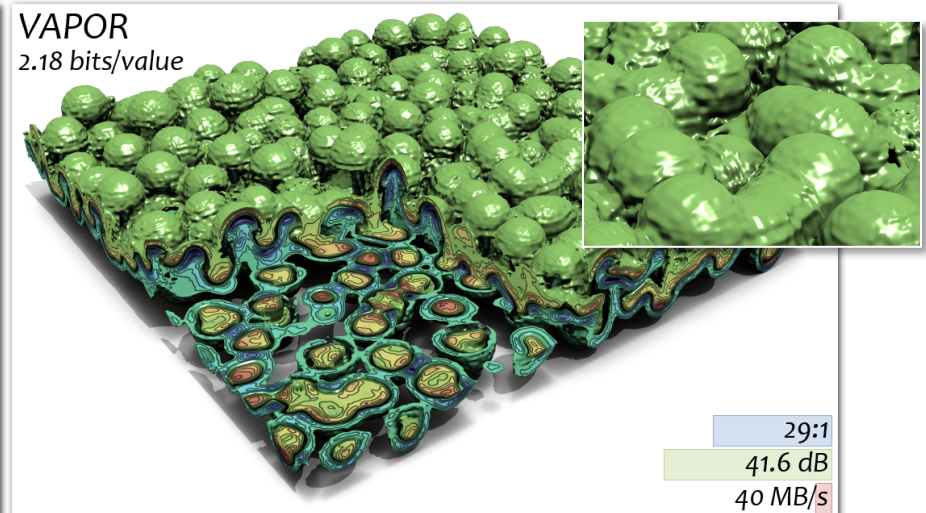
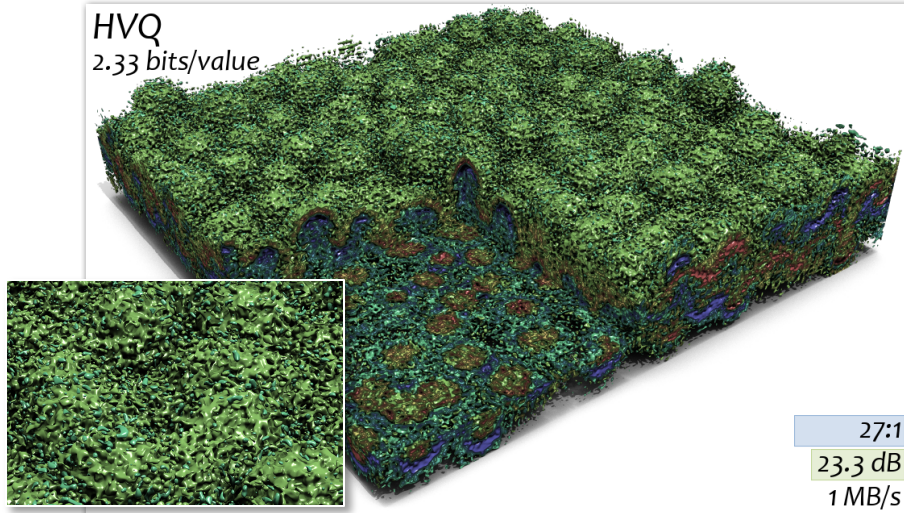
zfp  
1.96 bits/value



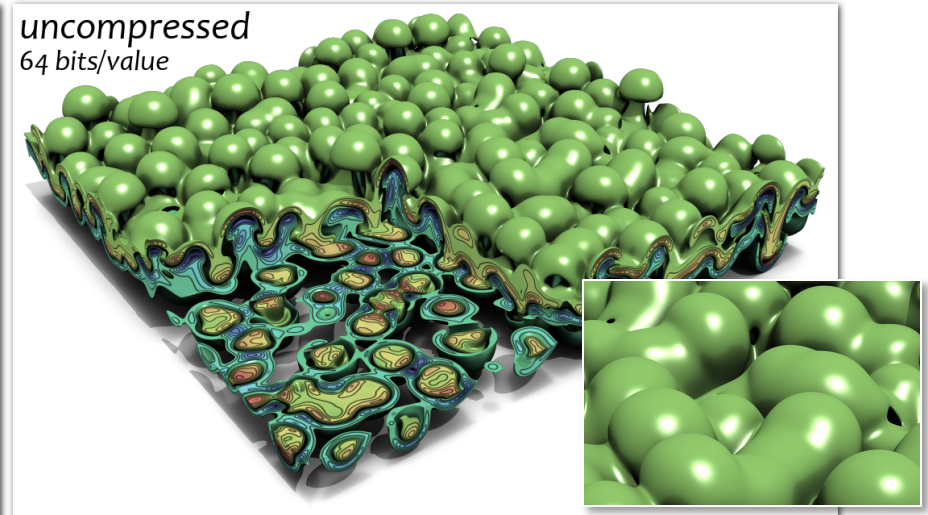
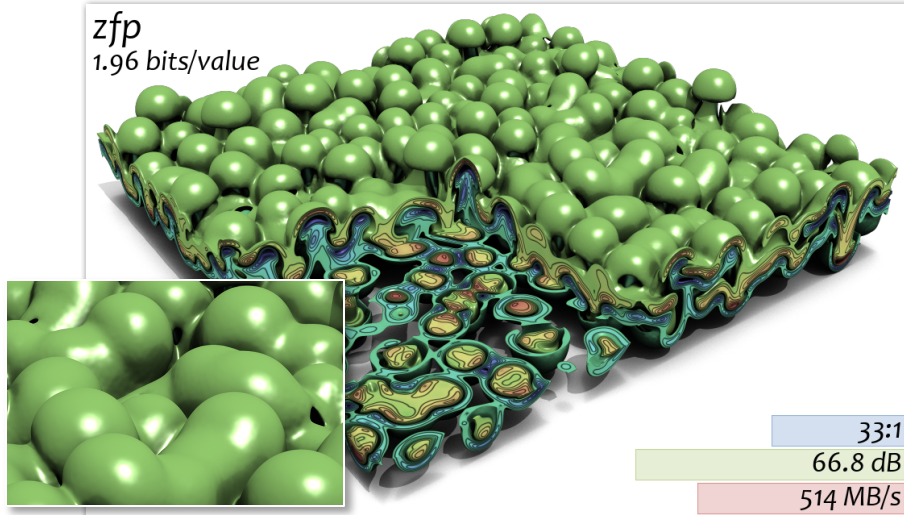
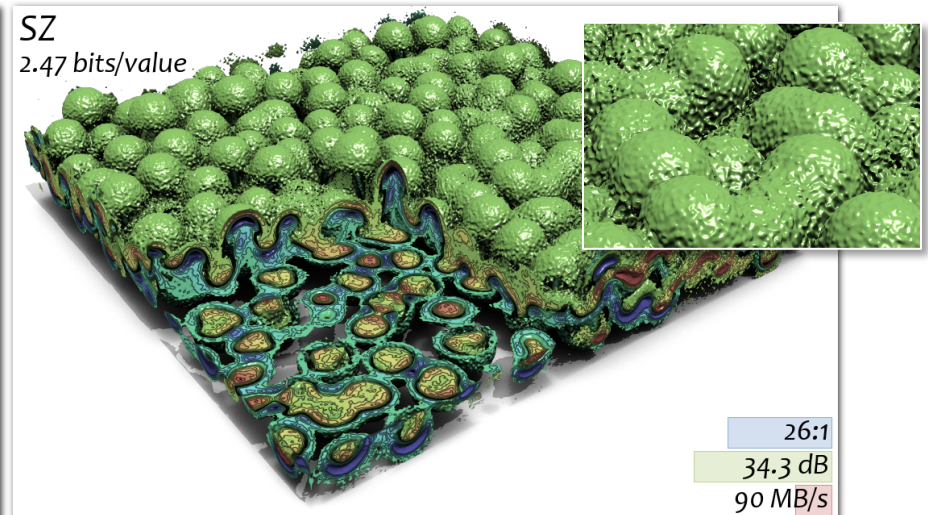
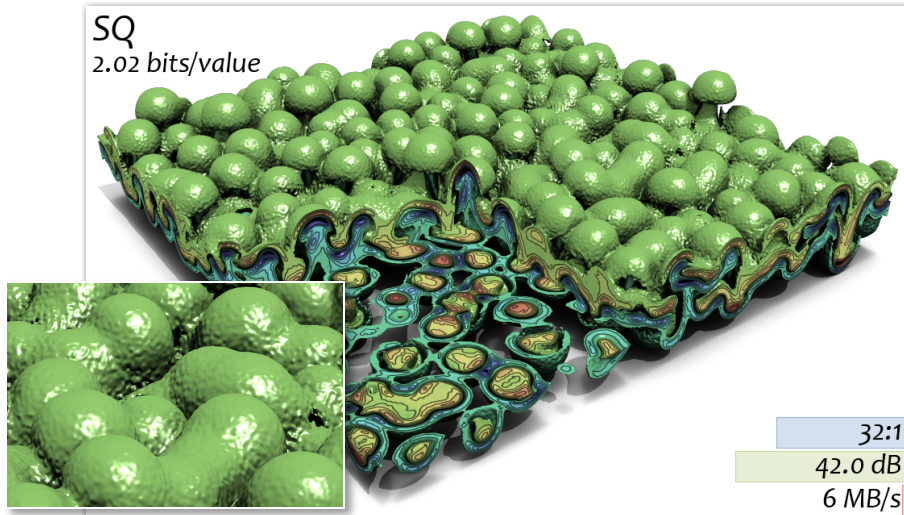
uncompressed  
64 bits/value



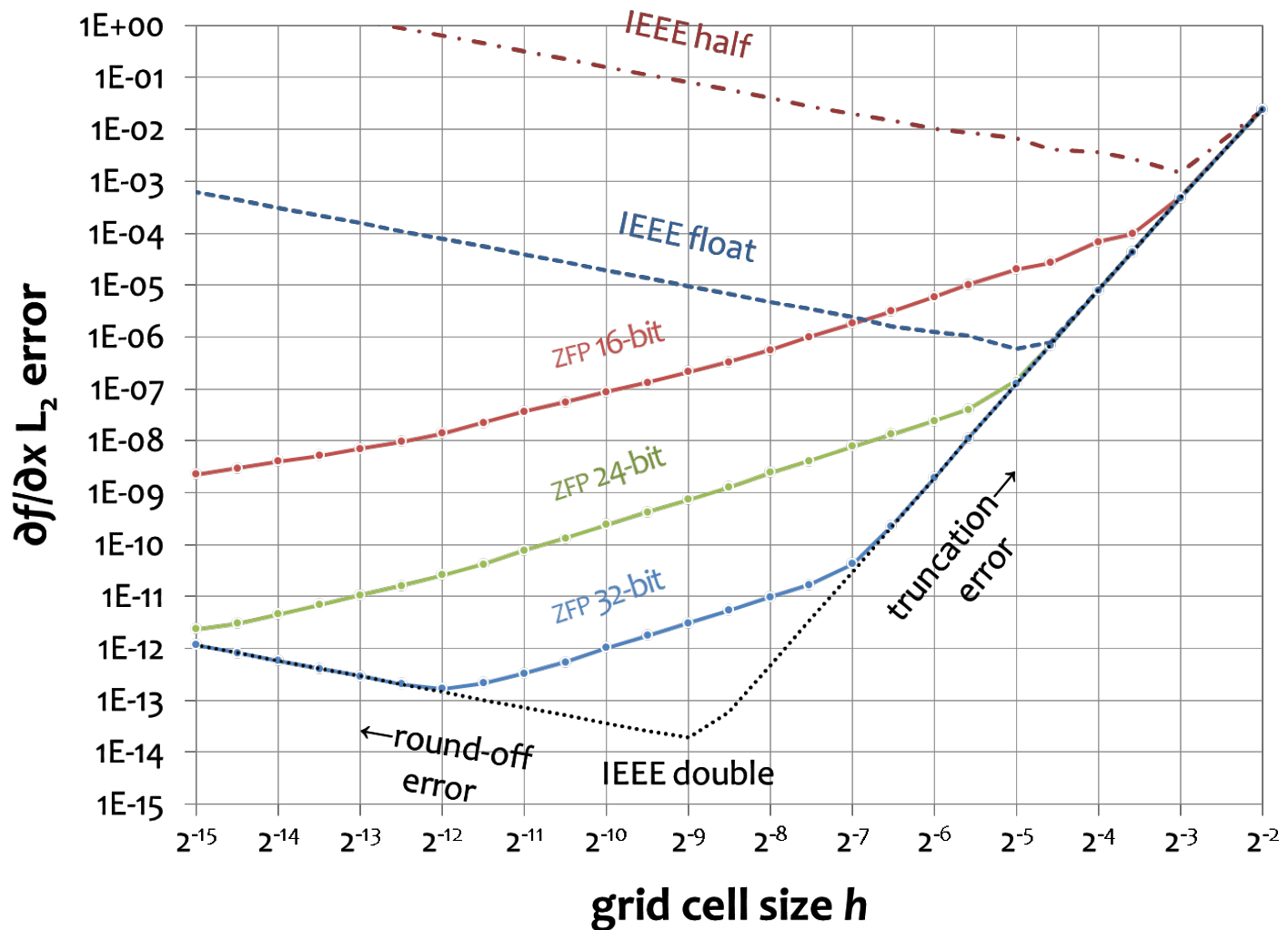
# ZFP shows no artifacts in derivative computations (velocity divergence)



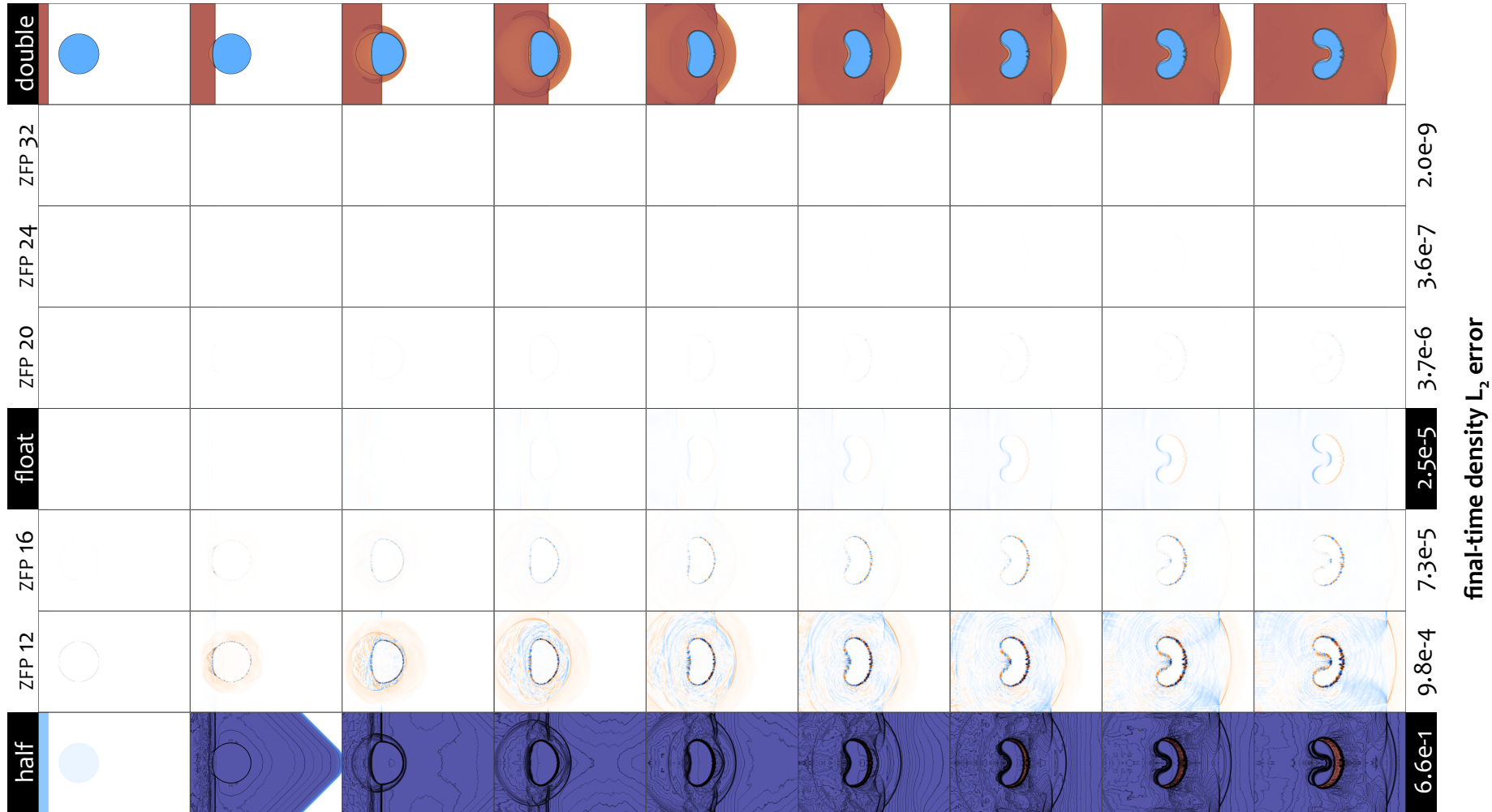
# ZFP shows no artifacts in derivative computations (velocity divergence)



# ZFP improves accuracy in finite difference computations using less precision than IEEE

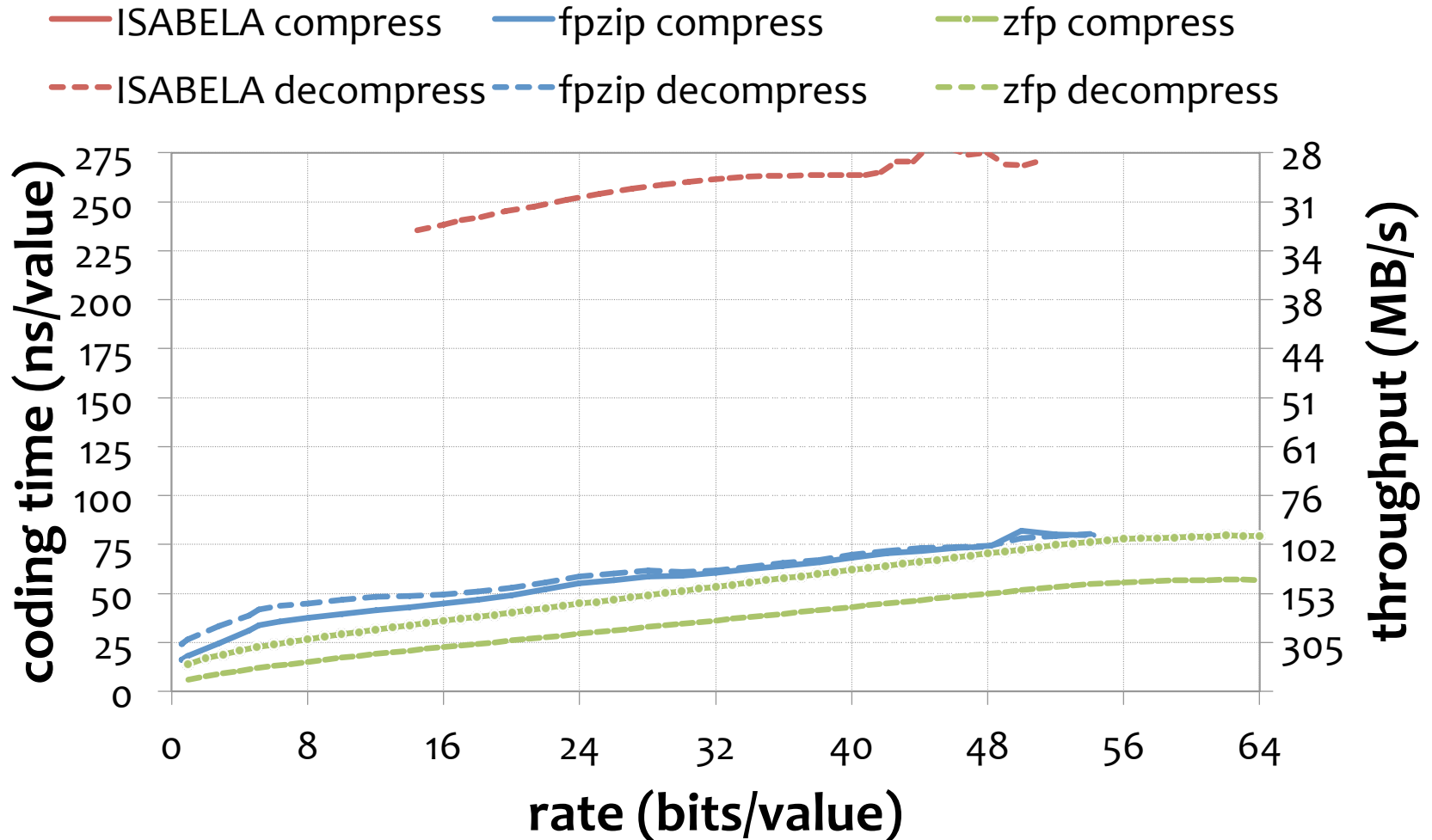


# ZFP achieves 5.3x inline compression with $<0.1\%$ error in shock-bubble interaction



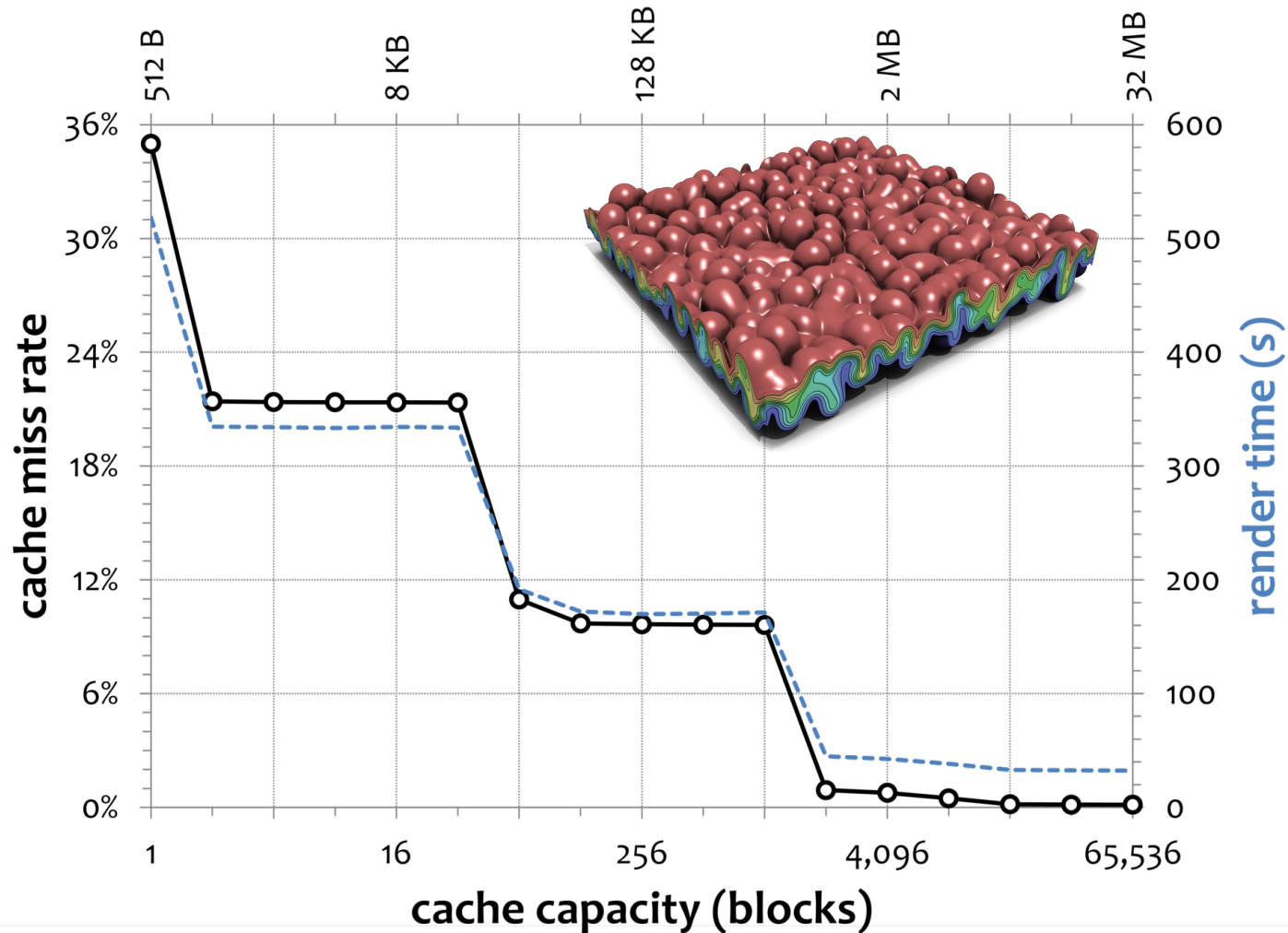
final-time density  $L_2$  error

# ZFP achieves up to 2 GB/s/core throughput and delivers predictable performance



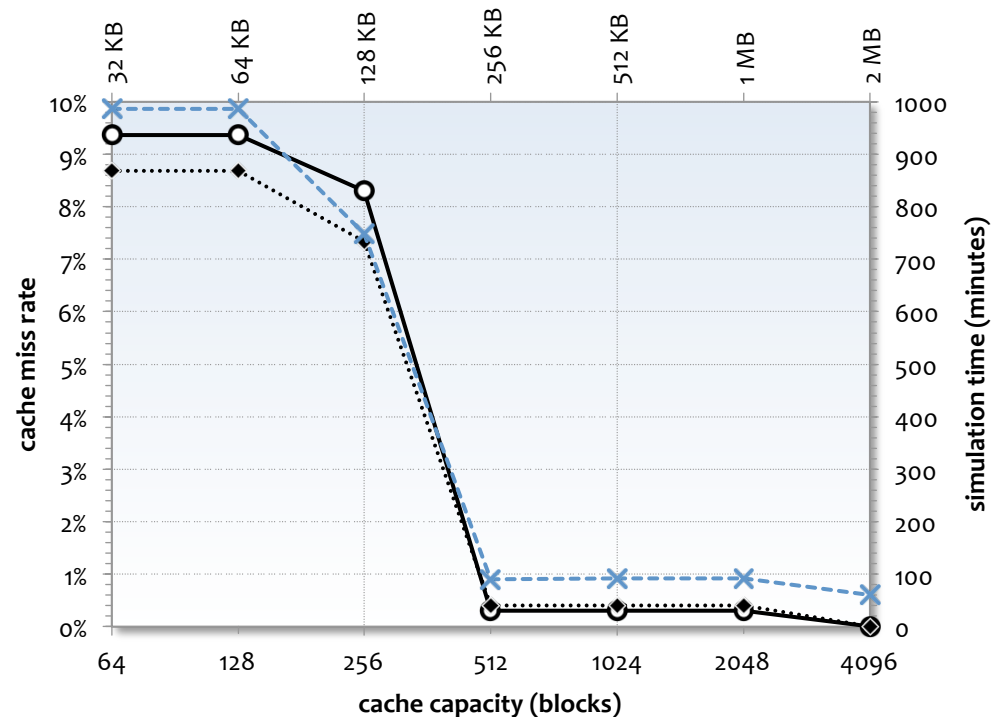
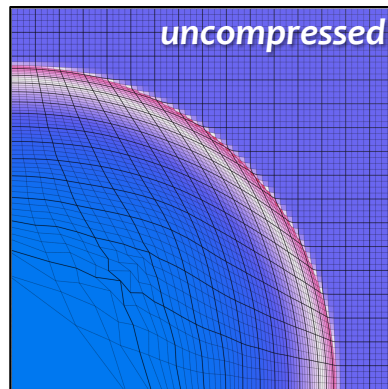
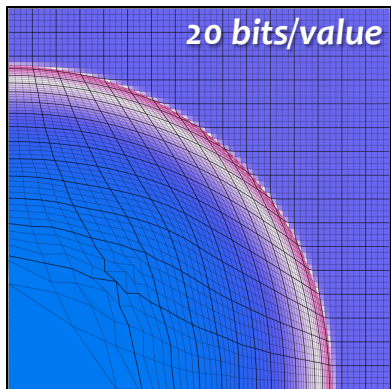


# Our s/w implementation is 17% away from performance break-even point on single core



# Substantial data movement reduction was achieved in LULESH mini app

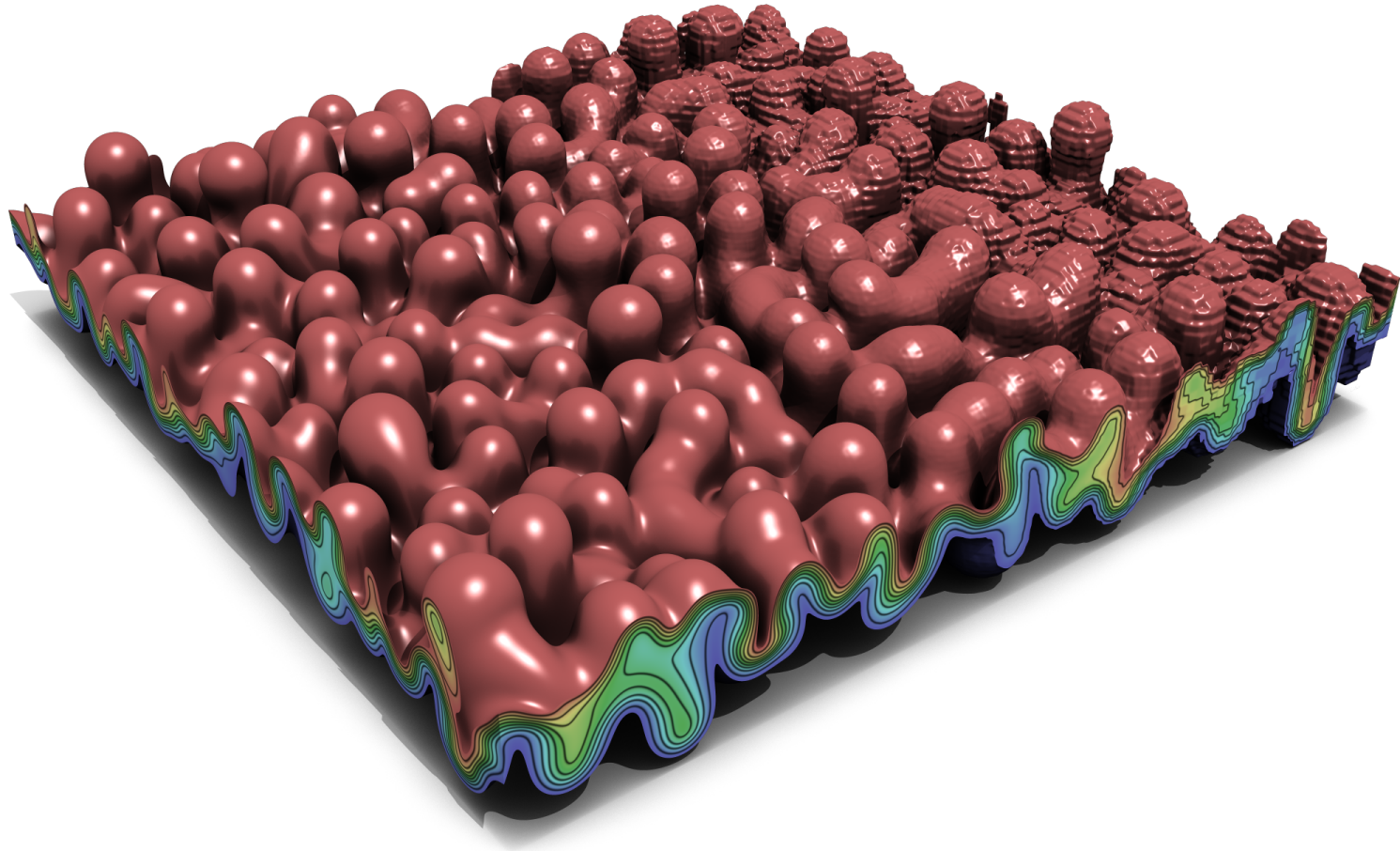
quantity	with compression	without compression	ratio
cache size; total array size	256 KB	2 MB	8x
decompressed; read accessed	51.4 GB	733 GB	14x
compressed; write accessed	15.9 GB	242 GB	15x
running time	90 min	15 min	6x
compression calls/ block/time step	1.17	0	
read miss rate	0.3%	N/A	



# Summary: ZFP has the potential to dramatically reduce data movement in HPC applications

- ZFP is a better floating-point representation than IEEE
  - Half the storage for the same accuracy
  - Four orders of magnitude higher accuracy for the same storage
- ZFP is so far primarily used for read-only access
  - I/O compression: Turbulence, combustion, fusion, radiation transport, seismology, climate, weather, VFX, high-dynamic-range images, ...
  - In-memory table compression: Equation of state, opacity tables
- Inline compression reduces memory footprint, bandwidth
  - Alleviates many-core cache and memory bandwidth contention
  - ZFP's symmetric performance ensures fast compress and decompress
- Hardware compression is ubiquitous on GPUs, mobile devices
  - How come HPC community has not caught on?

# ZFP is released as open source



<http://computation.llnl.gov/projects/floating-point-compression>

