

Efficient Parallel Software for Tucker Decompositions of Dense Tensors

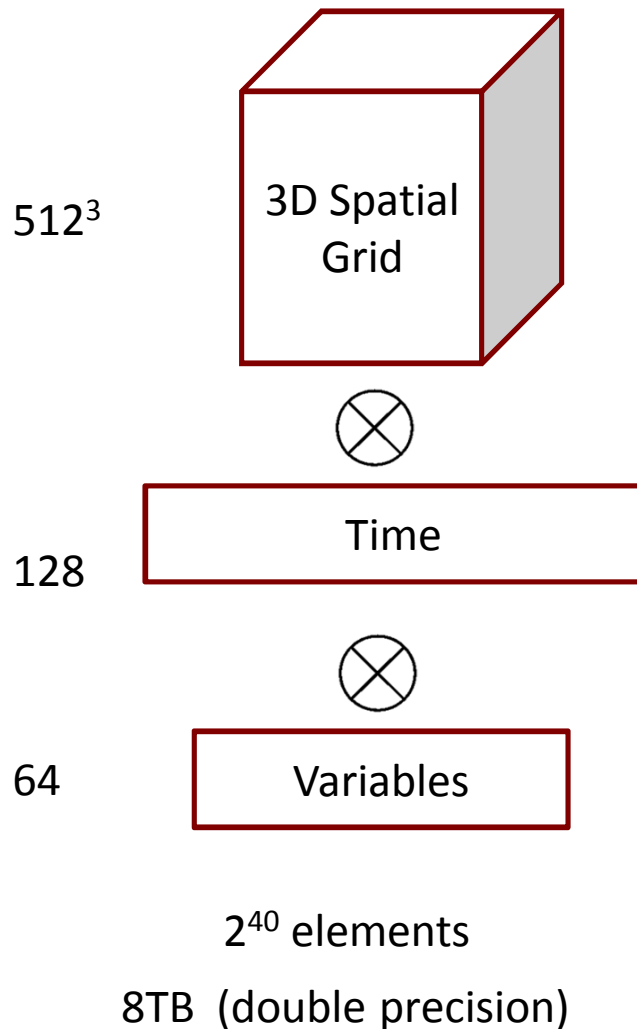
Alicia Klinvex¹, Grey Ballard², Tamara G. Kolda¹

¹Sandia National Laboratories

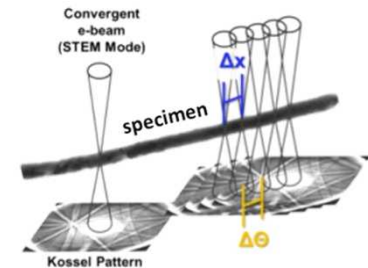
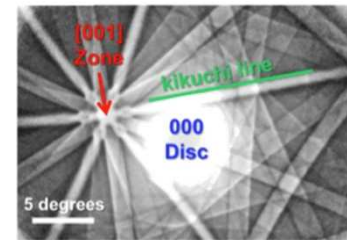
²Wake Forest University

SIAM CSE '17, Atlanta, March 2, 2017

Motivation: Advanced Simulations and Experiments Deluged by Data

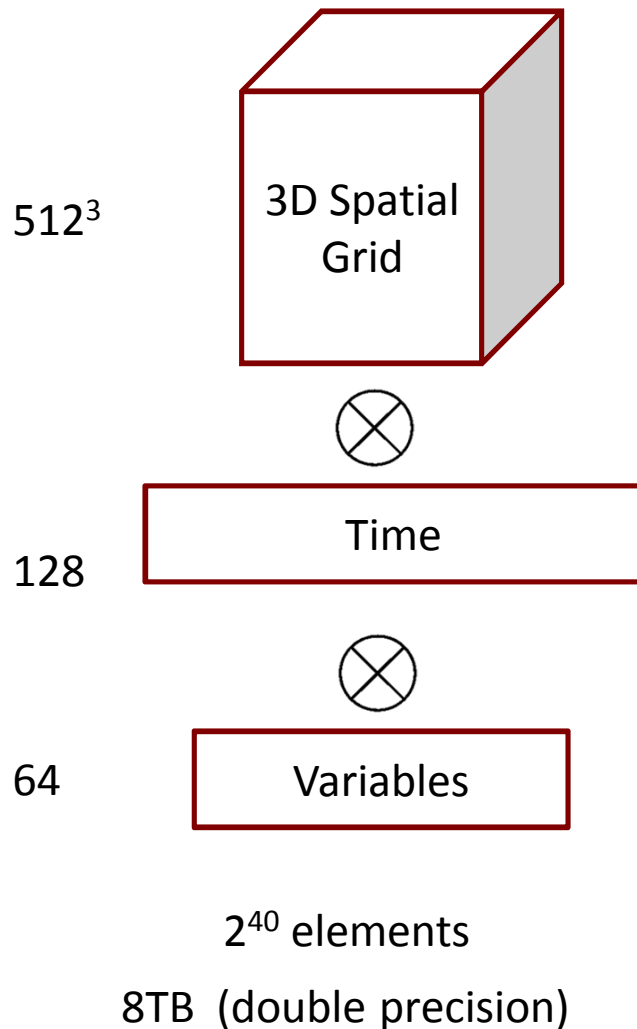


- Combustion simulations
 - S3D code uses direct numerical simulation
 - Gold standard for comparisons, but...
 - Single experiment produces terabytes of data
 - Storage limits spatial, temporal resolutions
 - Difficult to analyze or transfer data
- Electron microscopy
 - New technology produces 1-D spectra and 2-D diffraction patterns *at each pixel*
 - Single experiment produces terabytes of data
 - Usually 10-100 experiments per
 - Limited hardware utilization due to storage limits



- Other applications
 - Telemetry
 - Cosmology Simulations
 - Climate Modeling

Motivation: Advanced Simulations and Experiments Deluged by Data

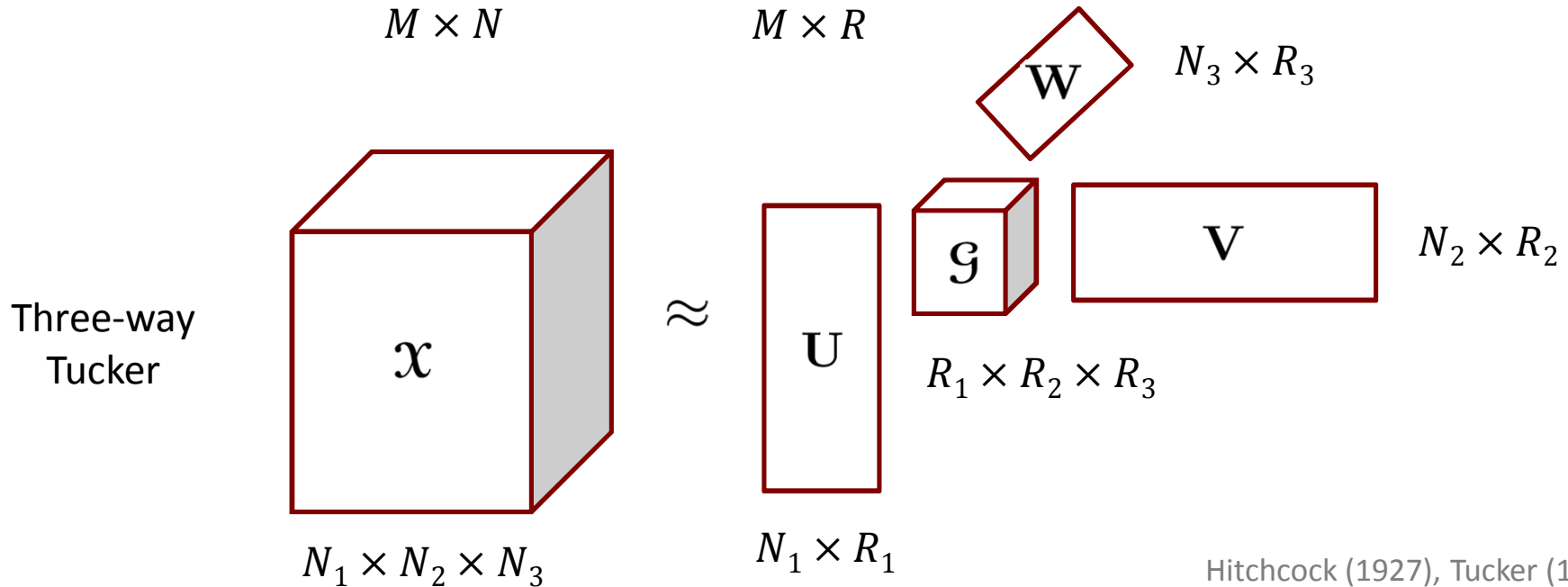
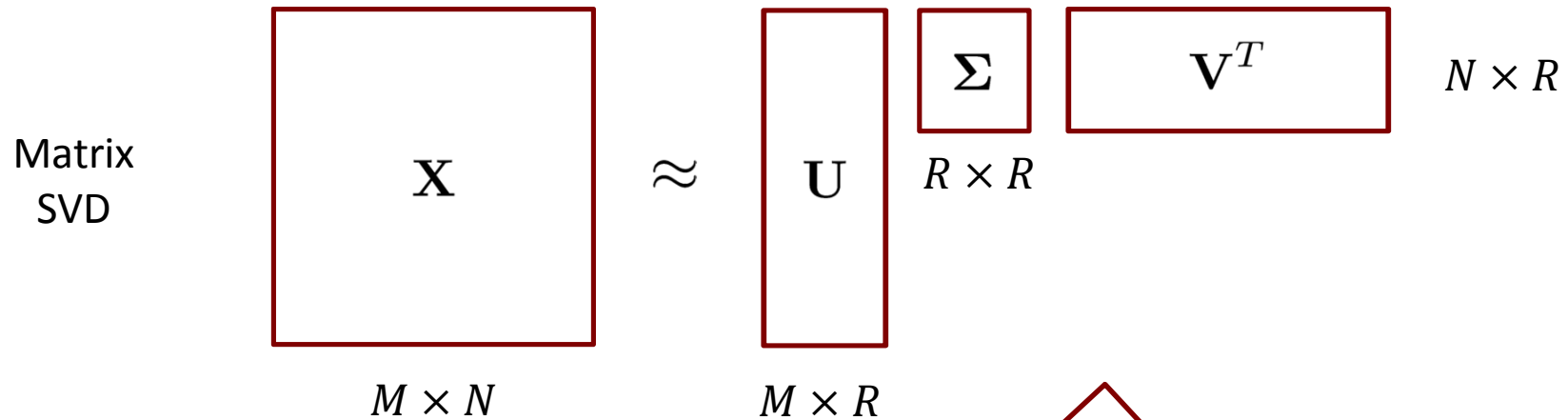


- We can compress this data using the power of math!

Outline

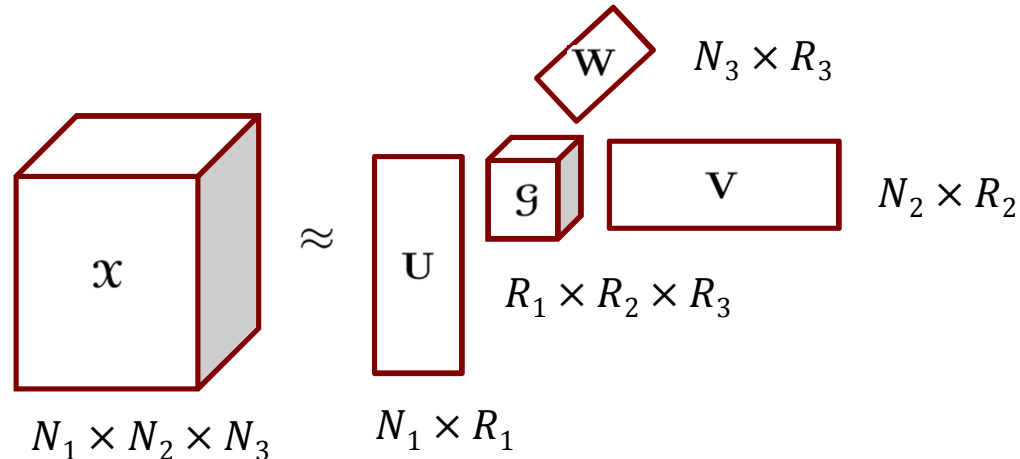
- Tucker decomposition definition
- ST-HOSVD algorithm
- TuckerMPI implementation
- Combustion simulation results

Tucker Compression: Extends the Matrix SVD to Multiway Arrays



Hitchcock (1927), Tucker (1966)

Tucker Compression (3-way)



$$\mathcal{X} \approx \mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W}$$

\mathcal{G} = “Core Tensor” = Reduced representation, determined by factor matrices

$\mathbf{U}, \mathbf{V}, \mathbf{W}$ = “Factor Matrices” = Orthogonal matrices spanning high-variance subspaces

\times_k = Tensor-times-matrix in mode k

Tucker Compression (d-way)

$$\underbrace{\mathcal{X}}_{N_1 \times N_2 \cdots \times N_d} \approx \underbrace{\mathcal{G}}_{R_1 \times R_2 \cdots \times R_d} \times_1 \underbrace{\mathbf{U}^{(1)}}_{N_1 \times R_1} \times_2 \underbrace{\mathbf{U}^{(2)}}_{N_2 \times R_2} \cdots \times_d \underbrace{\mathbf{U}^{(d)}}_{N_d \times R_d}$$

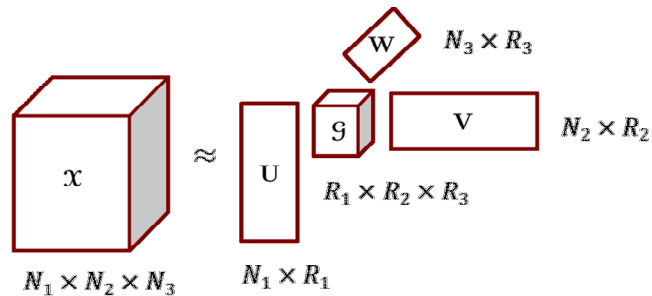
\mathcal{G} = “Core Tensor” = Reduced representation, determined by factor matrices

$\mathbf{U}^{(k)}$ = k th “Factor Matrix” = Orthogonal matrix spanning high-variance subspaces

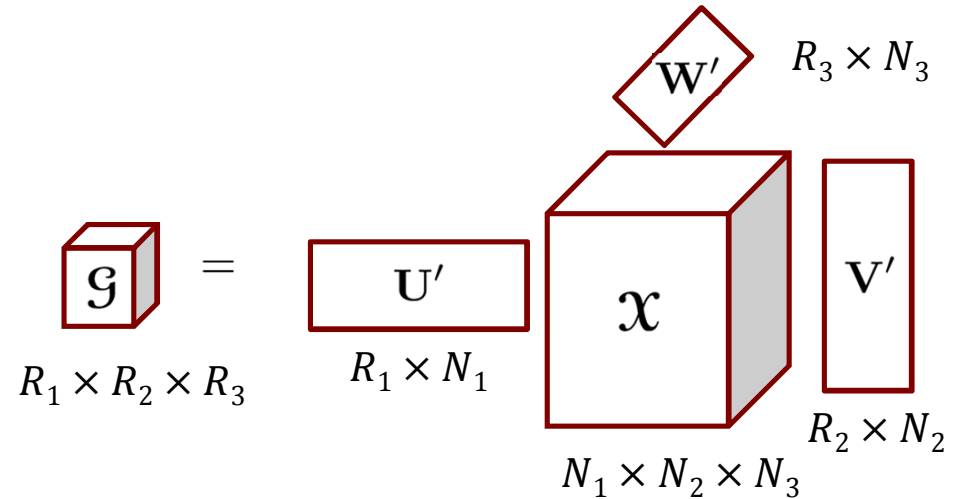
\times_k = Tensor-times-matrix in mode k

Compression Ratio: $C \approx \prod_{k=1}^d \frac{N_k}{R_k}$

Choosing Tucker Ranks to Retain Accuracy



Find orthogonal matrices \mathbf{U} , \mathbf{V} , \mathbf{W} that reduce the size of tensor but retain its “mass”



For a given relative error ϵ , choose projection ranks R_1 , R_2 , and R_3 such that:

$$\|\mathcal{X} - (\mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W})\| \leq \epsilon \|\mathcal{X}\|$$

Core tensor satisfies: $\mathcal{G} = \mathcal{X} \times_1 \mathbf{U}' \times_2 \mathbf{V}' \times_3 \mathbf{W}'$

$$\|\mathcal{X}\|^2 - \|\mathcal{G}\|^2 \leq \epsilon^2 \|\mathcal{X}\|^2$$

Algorithm: ST-HOSVD (3-way)

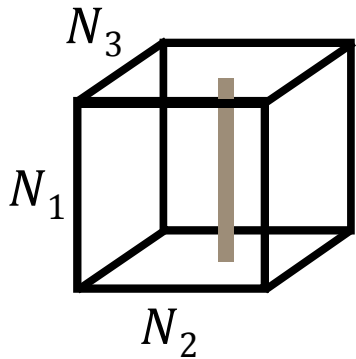
Vannieuwenhoven, Vandebril, Meerbergen (SISC 2012)

1. Choose \mathbf{U} with projection rank R_1 such that: $\|\mathbf{X}_{(1)}\|^2 - \|\mathbf{U}'\mathbf{X}_{(1)}\|^2 \leq \epsilon^2\|\mathbf{X}\|^2/3$
 - a) Compute gram matrix: $\mathbf{X}_{(1)}\mathbf{X}_{(1)}'$
 - b) Use eigendecomposition of $N_1 \times N_1$ matrix to choose R_1
 - c) Set $\mathbf{U} = R_1$ leading eigenvectors of gram matrix
2. Shrink to size $R_1 \times N_2 \times N_3$: $\mathcal{Y} = \mathbf{X} \times_1 \mathbf{U}'$
3. Choose \mathbf{V} with projection rank R_2 such that: $\|\mathbf{Y}_{(2)}\|^2 - \|\mathbf{V}'\mathbf{Y}_{(2)}\|^2 \leq \epsilon^2\|\mathbf{X}\|^2/3$
 - a) Compute gram matrix: $\mathbf{Y}_{(2)}\mathbf{Y}_{(2)}'$
 - b) Use eigendecomposition of $N_2 \times N_2$ matrix to choose R_2
 - c) Set $\mathbf{V} = R_2$ leading eigenvectors of gram matrix
4. Shrink to size $R_1 \times R_2 \times N_3$: $\mathcal{Z} = \mathcal{Y} \times_2 \mathbf{V}'$
5. Choose \mathbf{W} with projection rank R_3 such that: $\|\mathbf{Z}_{(3)}\|^2 - \|\mathbf{W}'\mathbf{Z}_{(3)}\|^2 \leq \epsilon^2\|\mathbf{X}\|^2/3$
 - a) Compute gram matrix: $\mathbf{Z}_{(3)}\mathbf{Z}_{(3)}'$
 - b) Use eigendecomposition of $N_3 \times N_3$ matrix to choose R_3
 - c) Set $\mathbf{W} = R_3$ leading eigenvectors of gram matrix
6. Shrink to size $R_1 \times R_2 \times R_3$: $\mathcal{G} = \mathcal{Z} \times_3 \mathbf{W}'$

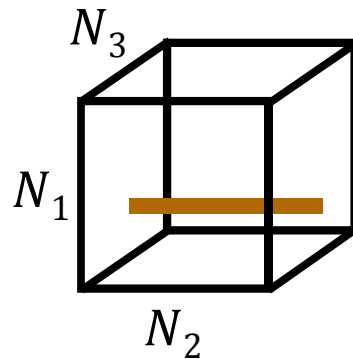
Key Kernels

- Three main steps for mode k
 - **GRAM** – Compute matrix product of unfolded tensor with itself.
 - Unfolded matrix size is $N_k \times (R_1 \cdots R_{k-1} N_{k+1} \cdots N_d)$.
 - Result is $N_k \times N_k$
 - **EVECS** – Compute eigenvalues and eigenvectors of $N_k \times N_k$ gram matrix. Use this to determine R_k
 - Call LAPACK, since matrix is small
 - **TTM** – Tensor-times-matrix to shrink mode k from size N_k to size R_k .
 - Input is size $R_1 \times \cdots \times N_k \times N_{k+1} \times \cdots \times N_d$
 - Result is size $R_1 \times \cdots \times R_k \times N_{k+1} \times \cdots \times N_d$
- These can be viewed as matrix operations

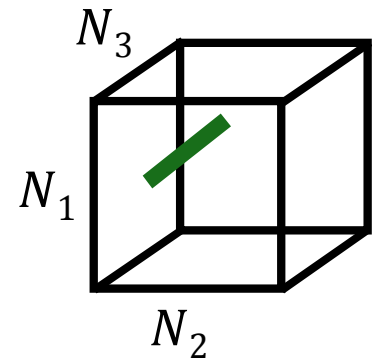
Matricization/Unfolding: Rearranging a Tensor as a Matrix



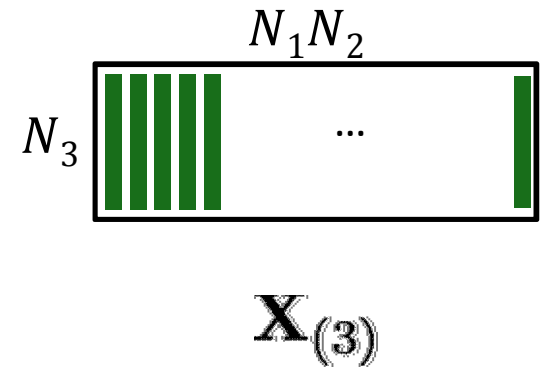
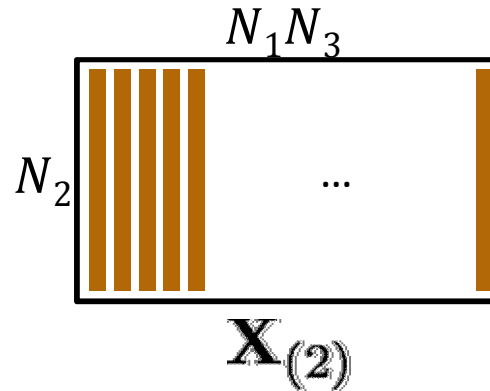
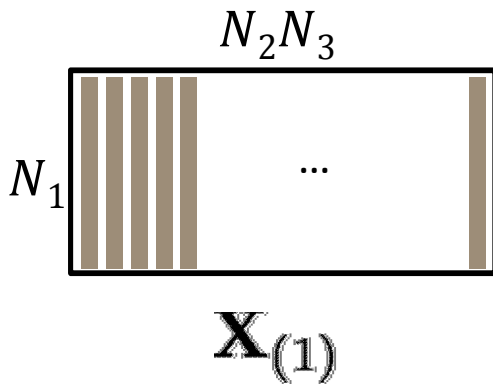
mode-1 fibers



mode-2 fibers



mode-3 fibers

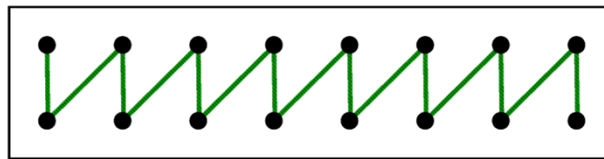


Mathematically convenient expression, but do not want to explicitly form the unfolded matrix due to the expense of the *data movement*.

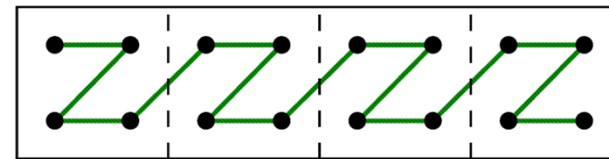
No Data Movement in Local Unfolding – Block Structure

- Assume local data is stored so that mode-1 unfolding is in column major order
- All unfoldings are blocked, with different block sizes
- Rather than rearrange data (standard practice), exploit structure with block operations

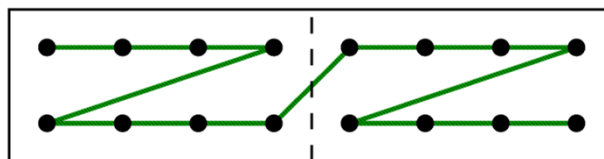
Local Layout: $2 \times 2 \times 2 \times 2$ (4-way Tensor)



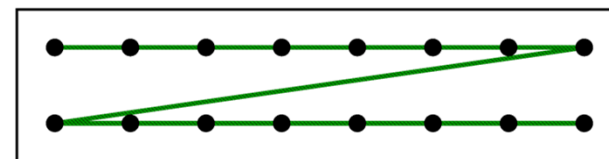
$k = 1$



$k = 2$



$k = 3$

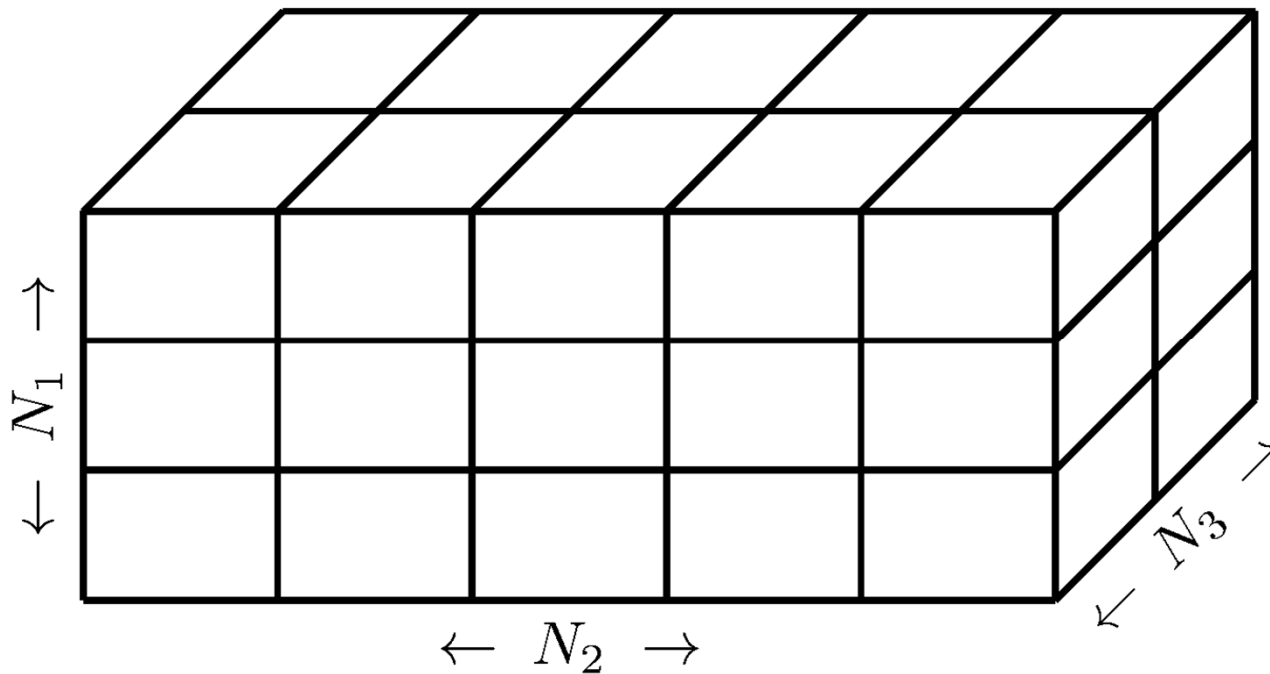


$k = 4$

Parallel Tucker Decomposition

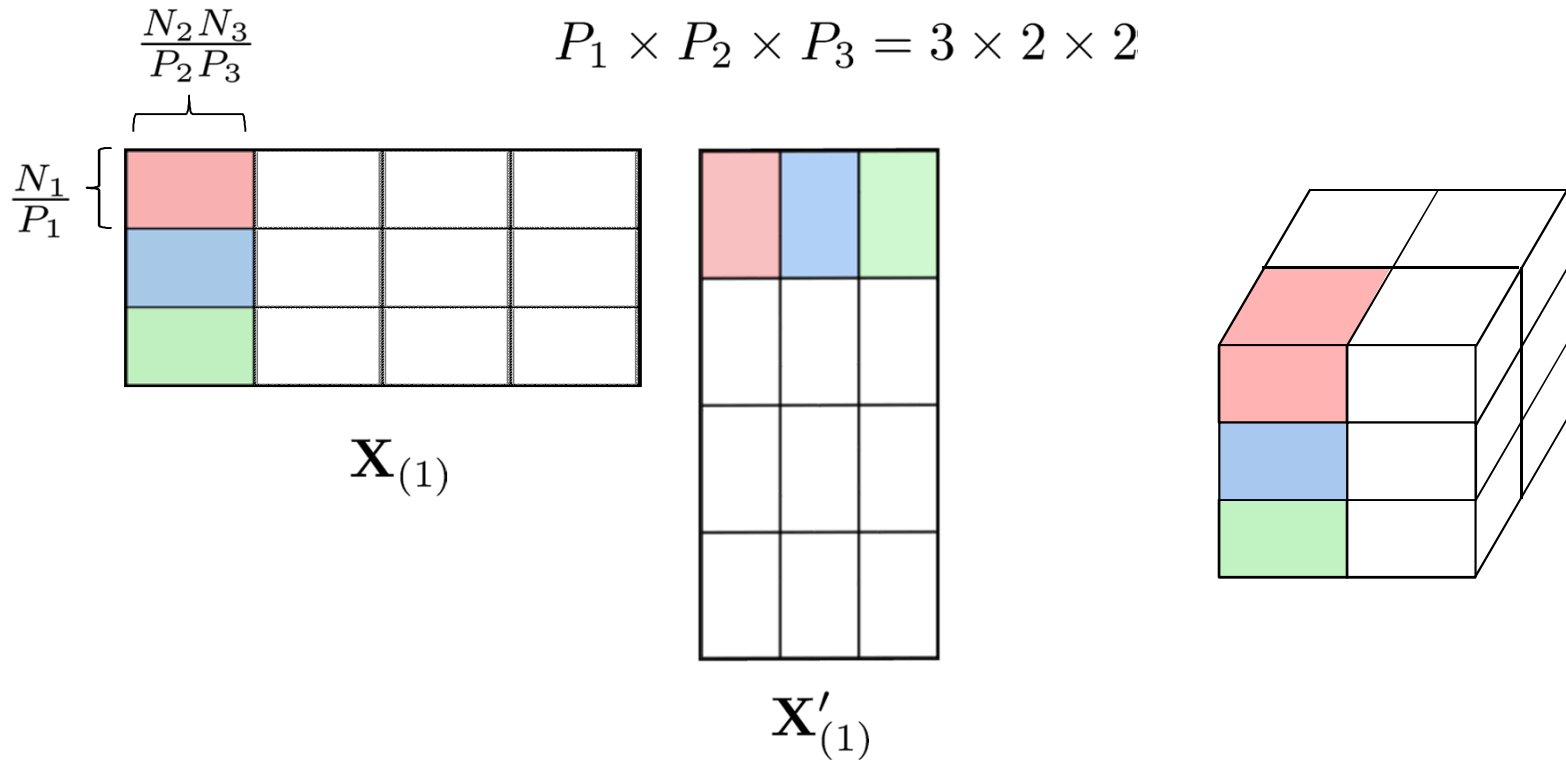
For N-way tensor, Cartesian block distribution on N-way processor grid

$$P_1 \times P_2 \times P_3 = 3 \times 5 \times 2$$



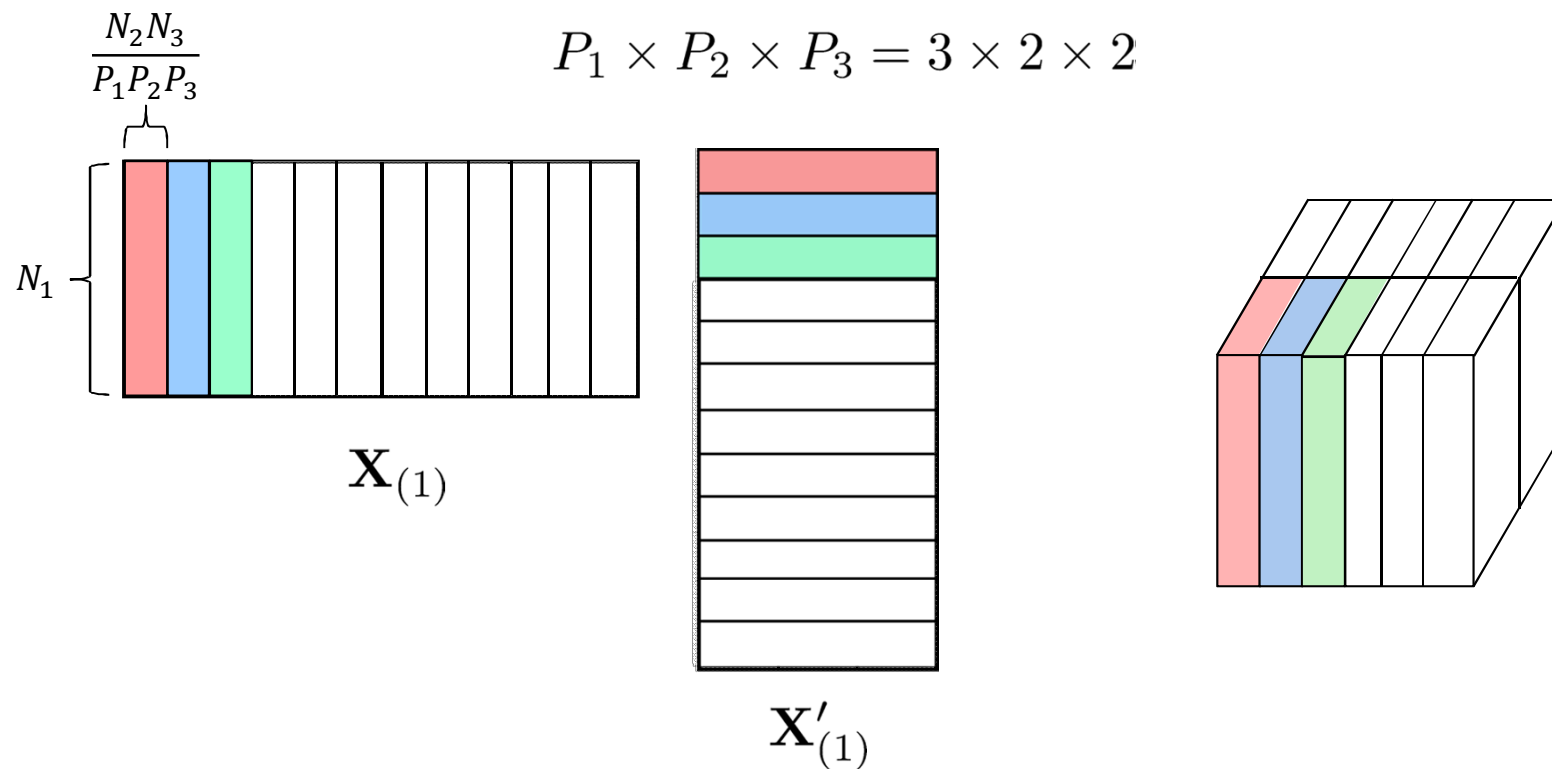
Local block size: $\frac{N_1}{P_1} \times \frac{N_2}{P_2} \times \frac{N_3}{P_3}$

Gram Matrix Kernel: Parallel Computation



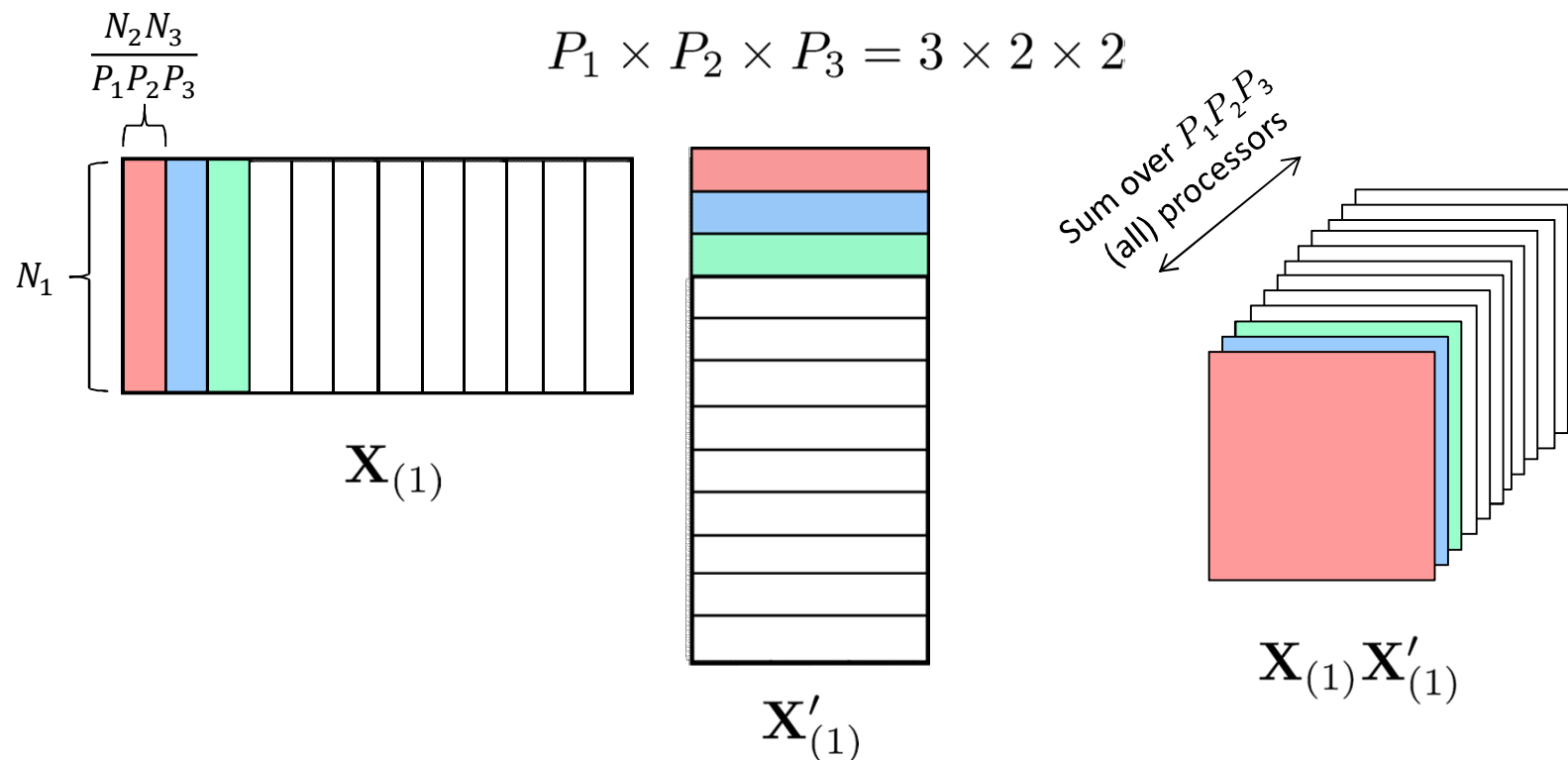
*Rearrange data to be block
column format*

Gram Matrix Kernel: Parallel Computation



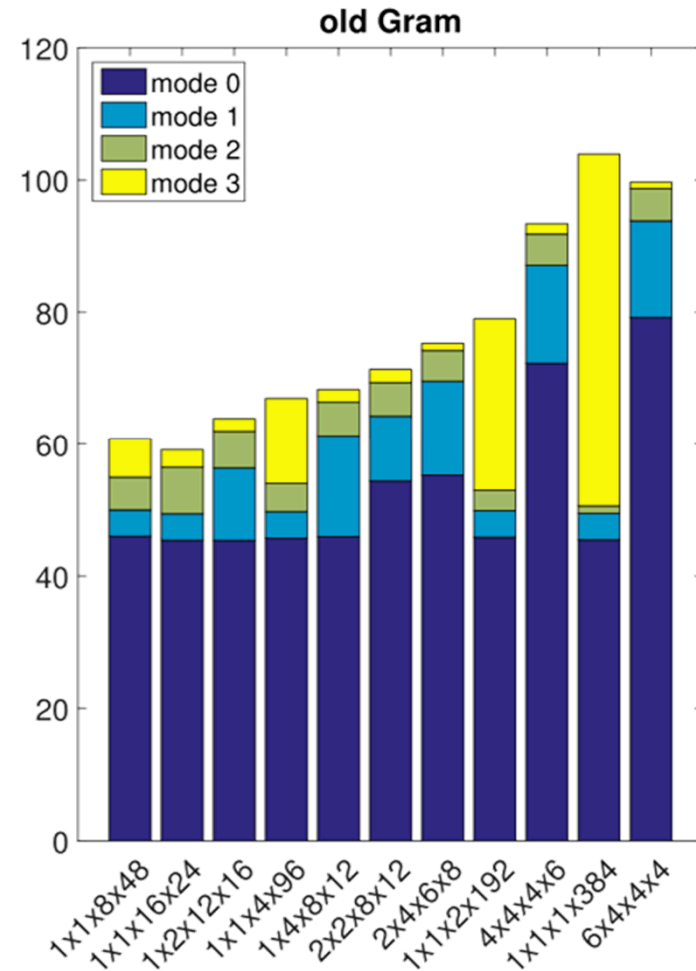
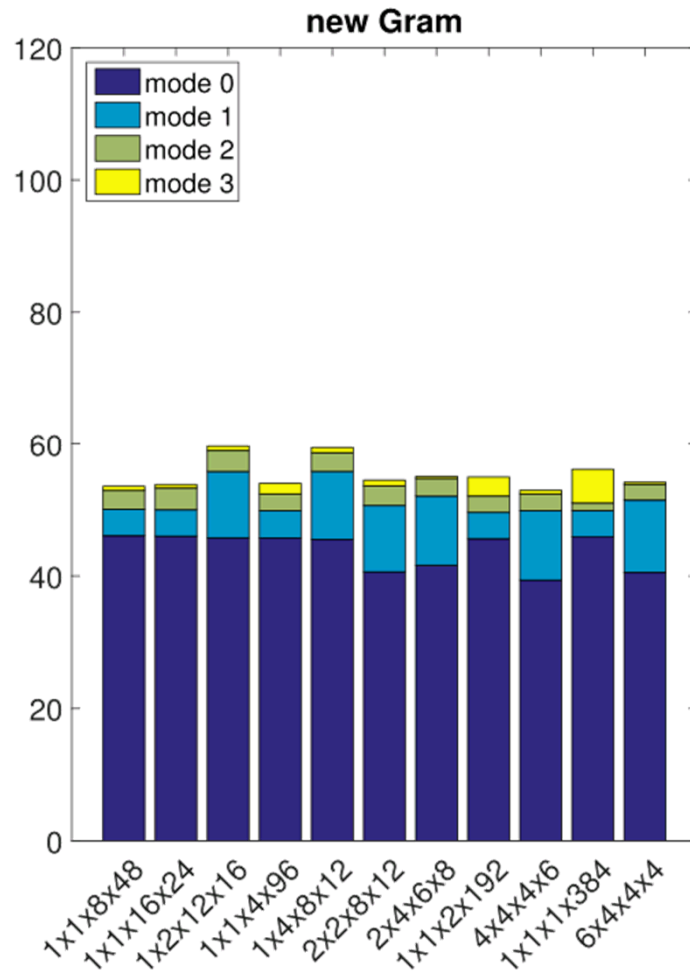
- Each processor column rearranges its data (with P_1 nodes)

Gram Matrix Kernel: Parallel Computation



- Each processor column rearranges its data (with P_1 nodes)
- Then computes local outer product
- Sum across all $P_1 P_2 P_3$ groups with all-reduce

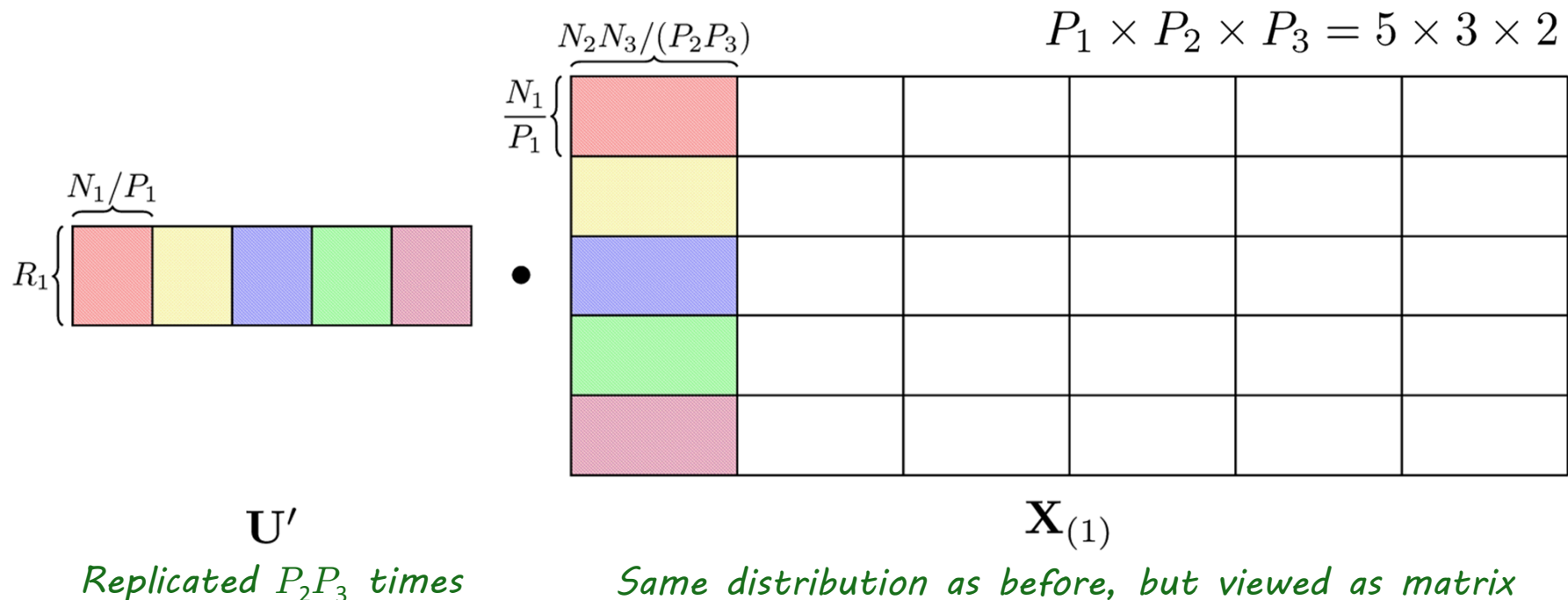
New Gram Improves Worst-Case Grid Choice Performance



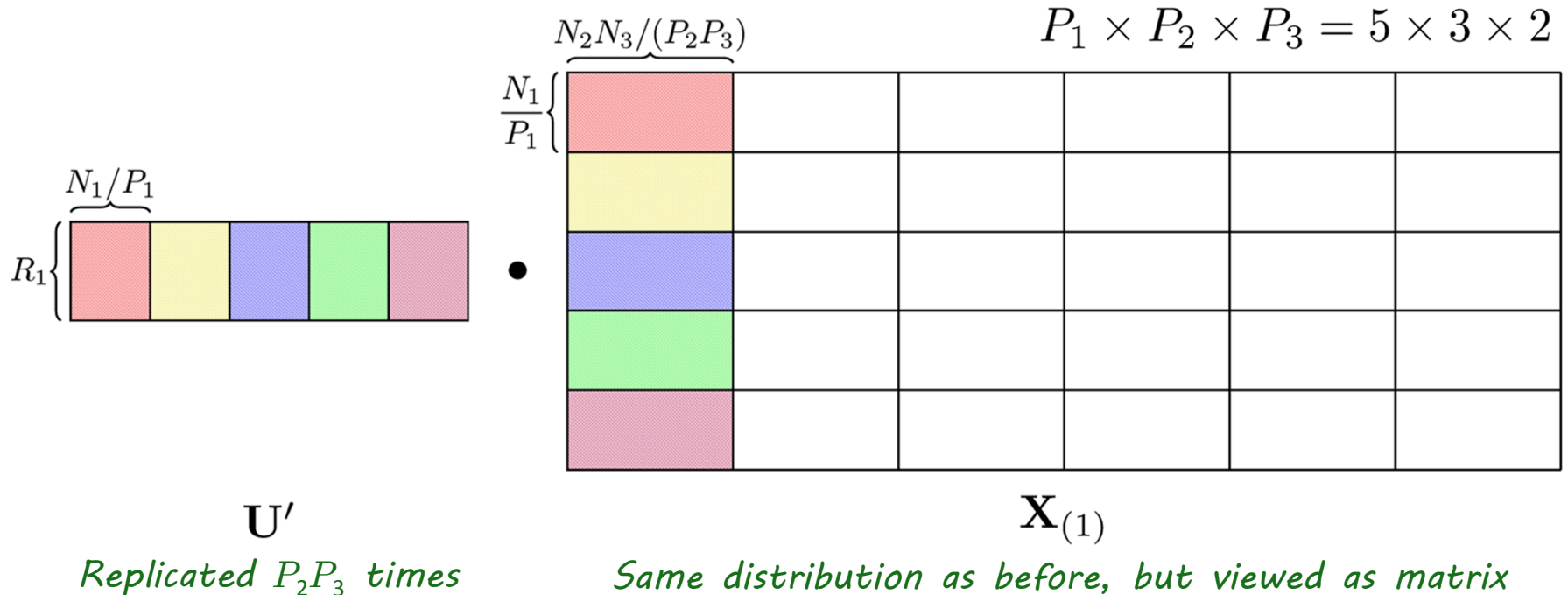
4-way Tensor of size 384^4 reduced to size 96^4 .

Comparisons on Kahuna @ Sandia: 120 Nodes x 2 Intel Haswell 14-core x 256 GB

Tensor-Times-Matrix Kernel: Parallel Computation



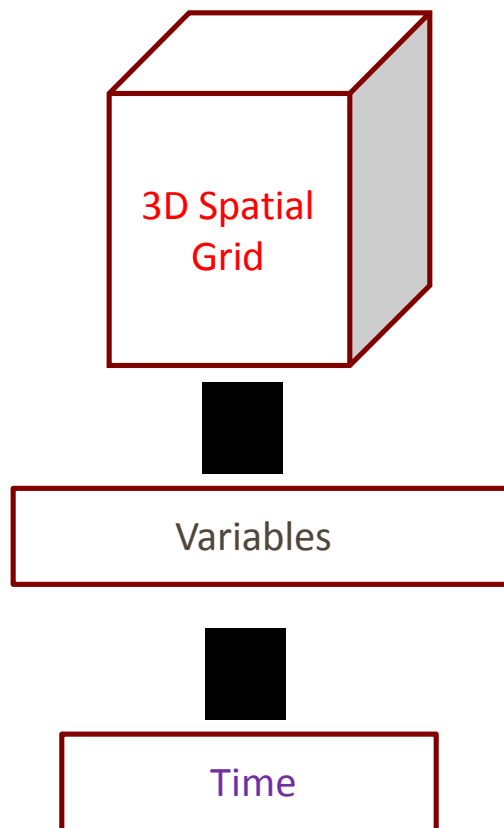
Tensor-Times-Matrix Kernel: Parallel Computation



- Each node locally computes product of its part of \mathbf{U} with its part of $\mathbf{X}_{(1)}$
- Communication is a reduce-scatter on the result within set of P_1 nodes
- Output is distributed same as $\mathbf{X}_{(1)}$ but with a smaller first dimension, i.e., size R_1/P_1
- Works the same for every mode (with different views)

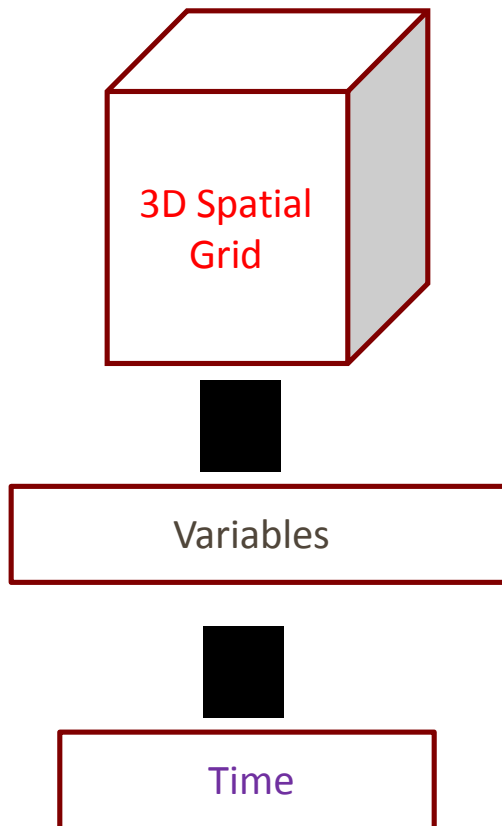
Results

Combustion results



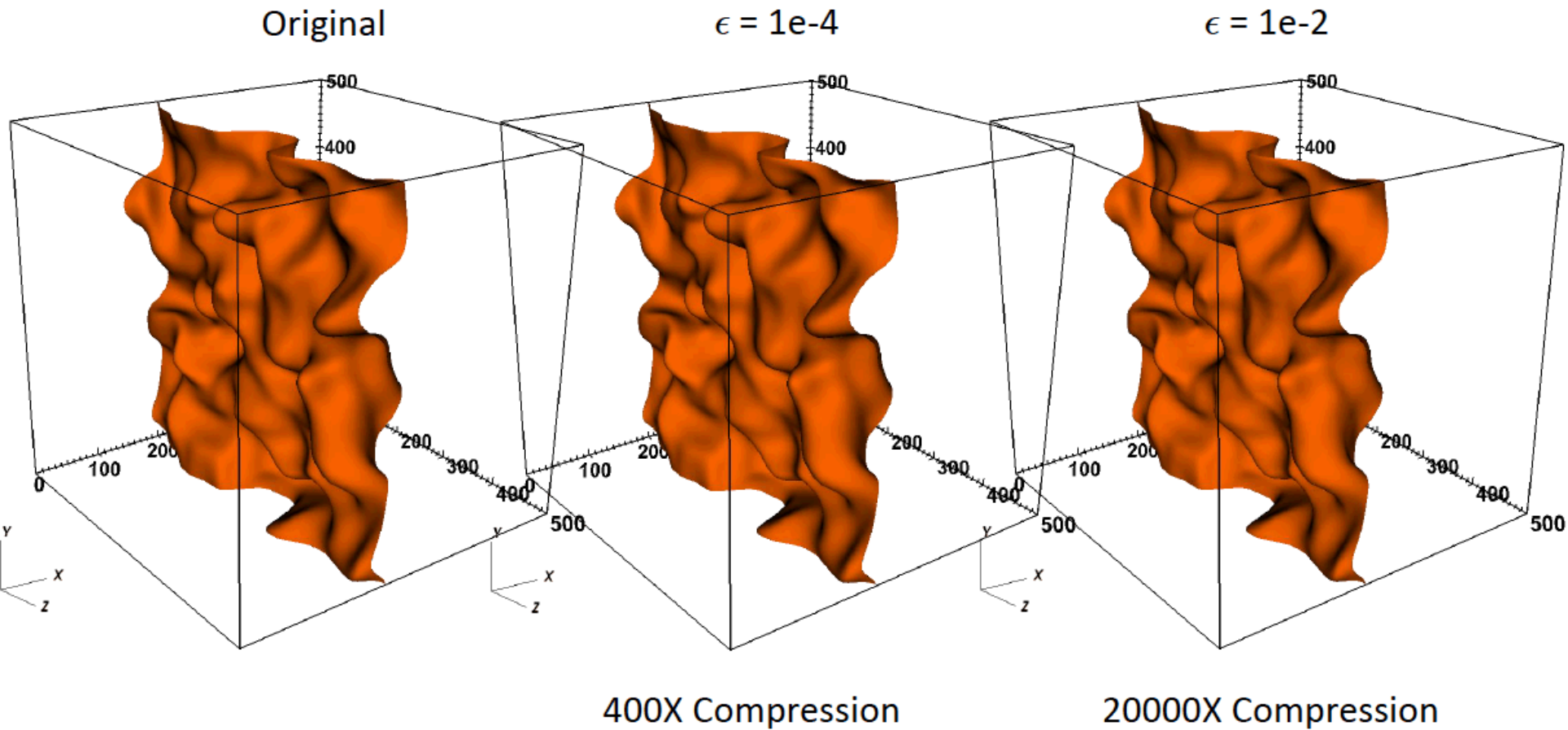
	Original	$\epsilon = 10^{-4}$	$\epsilon = 10^{-2}$
HCCI	672 x 672 x 32 x 626 (67 GB)	330 x 310 x 31 x 199 (5 GB)	111 x 105 x 22 x 46 (90 MB)
SP	500 x 500 x 500 x 11 x 400 4 TB	95 x 129 x 125 x 7 x 125 (10 GB)	30 x 38 x 35 x 6 x 11 (20 MB)
JICF	1500 x 2080 x 1500 x 18 x 10 6 TB	424 x 387 x 261 x 18 x 10 (57 GB)	90 x 61 x 48 x 13 x 6 (156 MB)

Combustion results

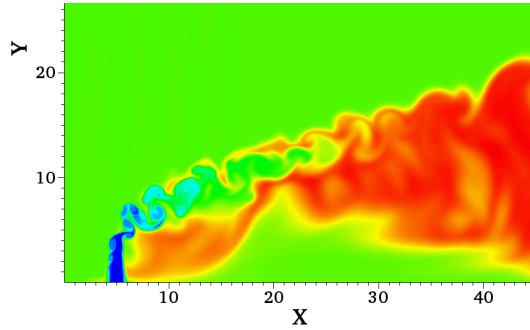


	Original	$\epsilon = 10^{-4}$	$\epsilon = 10^{-2}$
HCCI	672 x 672 x 32 x 626 (67 GB)	330 x 310 x 31 x 199 (14 X)	111 x 105 x 22 x 46 (760 X)
SP	500 x 500 x 500 x 11 x 400 4 TB	95 x 129 x 125 x 7 x 125 (410 X)	30 x 38 x 35 x 6 x 11 (20000 X)
JICF	1500 x 2080 x 1500 x 18 x 10 6 TB	424 x 387 x 261 x 18 x 10 (110 X)	90 x 61 x 48 x 13 x 6 (40000 X)

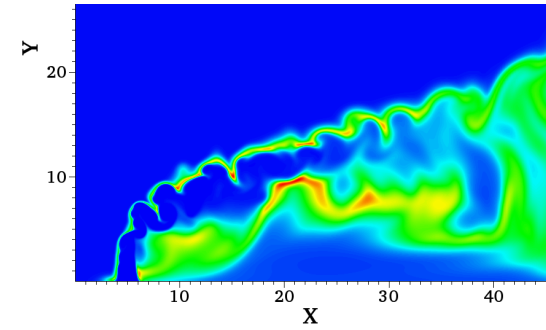
Original vs reconstructed



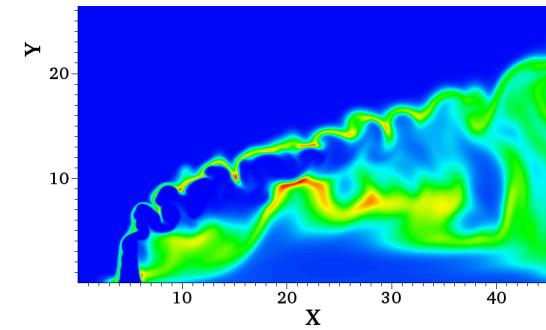
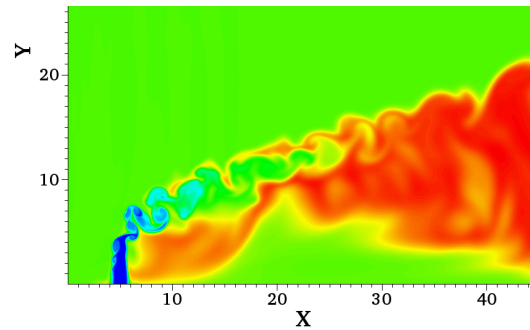
Original vs reconstructed



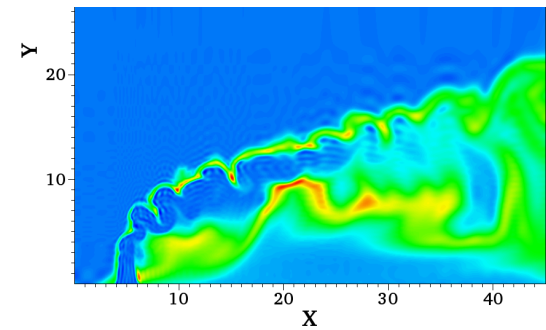
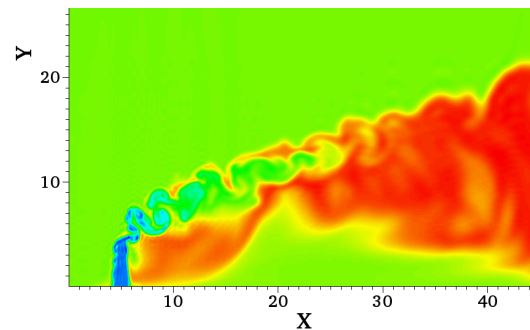
Original



$\epsilon = 10^{-4}$
(110X)



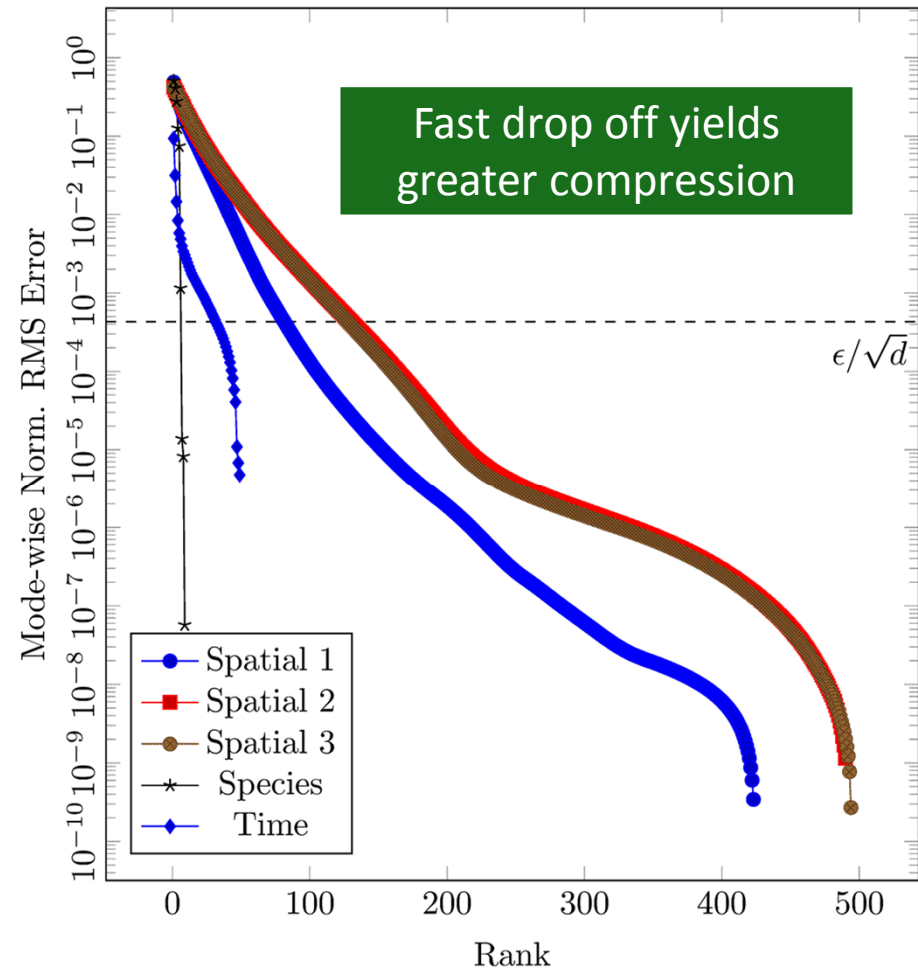
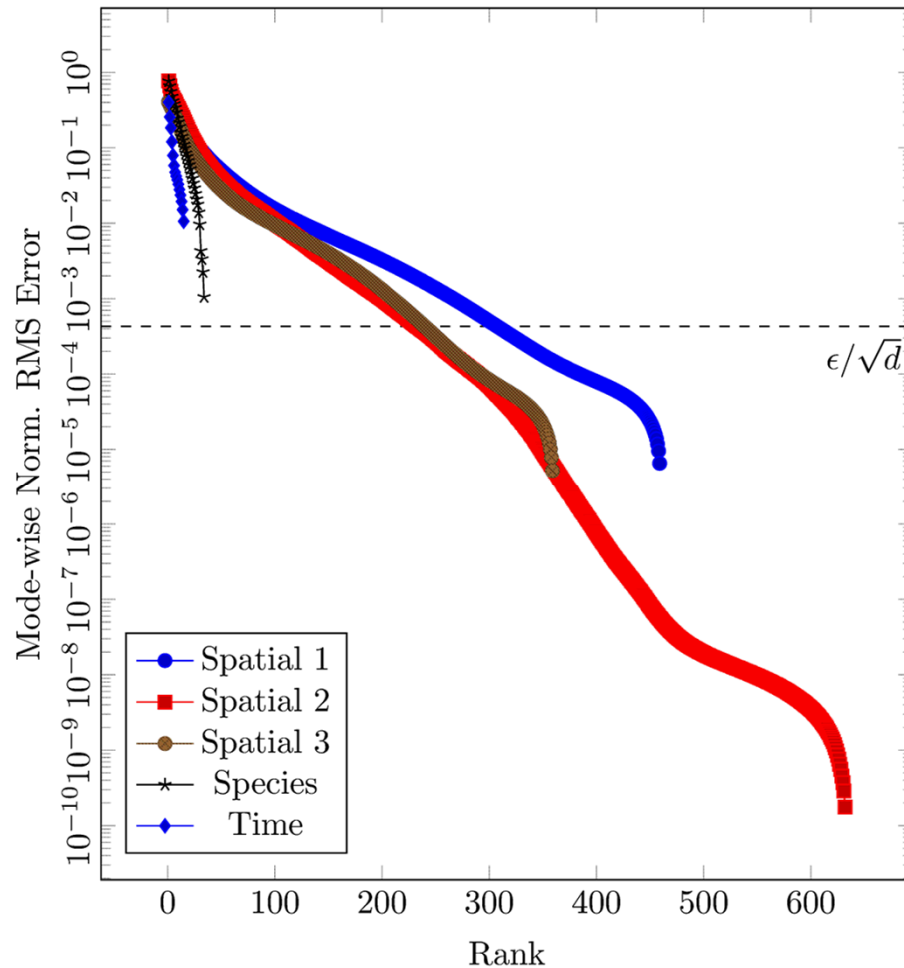
$\epsilon = 10^{-2}$
(40000X)



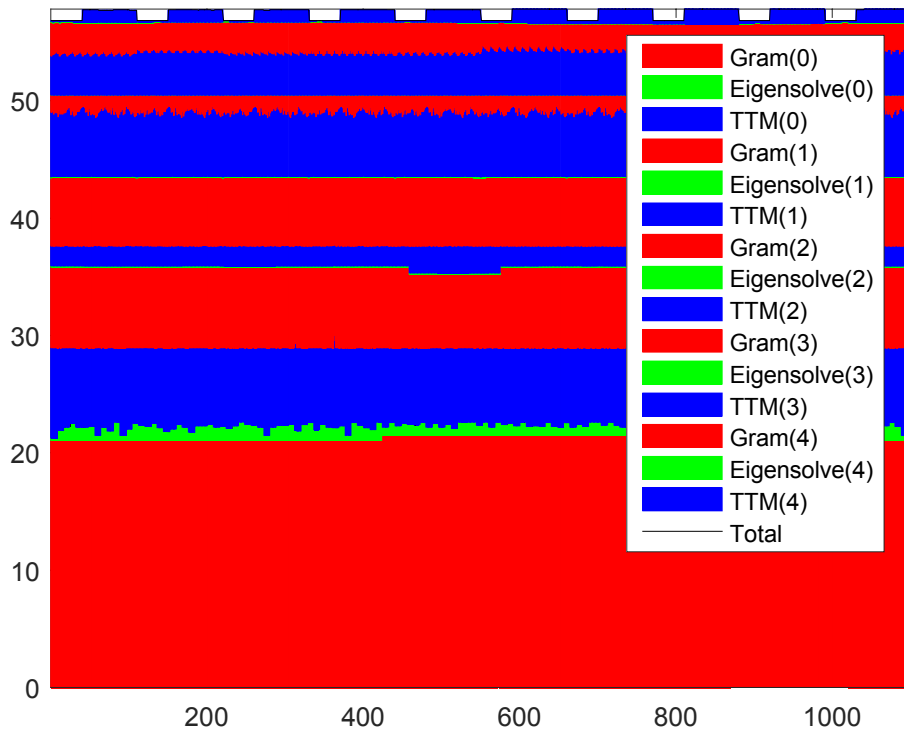
Compression Depends on Data Redundancy (Eigenvalues Don't Lie)

TJLR: 460 x 700 x 360 x 35 x 16

SP-50: 500 x 500 x 500 x 11 x 50



Parallel TuckerMPI performance



- Total of 55s; 1100 cores.
- 4.4TB -> 10GB (410X).
- Bulk of time is in first mode (GRAM computation).
- Fast BLAS for Eigensolve.
- Time for I/O is order of magnitude greater (~450s)

Key Feature: Need Only Do Partial Reconstruction on Laptops, etc.

Reconstruction requires as much space as the original data!

$$\hat{\mathbf{X}} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \times_4 \mathbf{U}^{(4)} \times_5 \mathbf{U}^{(5)}$$

$$N_1 \times N_2 \times N_3 \times N_4 \times N_5$$

But we can just reconstruct the portion that we need at the moment:

$$\bar{\mathbf{X}} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{C}^{(3)} \mathbf{U}^{(3)} \times_4 \mathbf{C}^{(4)} \mathbf{U}^{(4)} \times_5 \mathbf{C}^{(5)} \mathbf{U}^{(5)}$$

$$N_1 \times N_2 \times \frac{N_3}{2} \times 1 \times 1$$

$$\mathbf{C}^{(3)} = \begin{bmatrix} 1/2 & 0 & \cdots & 0 \\ 1/2 & 0 & \cdots & 0 \\ 0 & 1/2 & \cdots & 0 \\ 0 & 1/2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

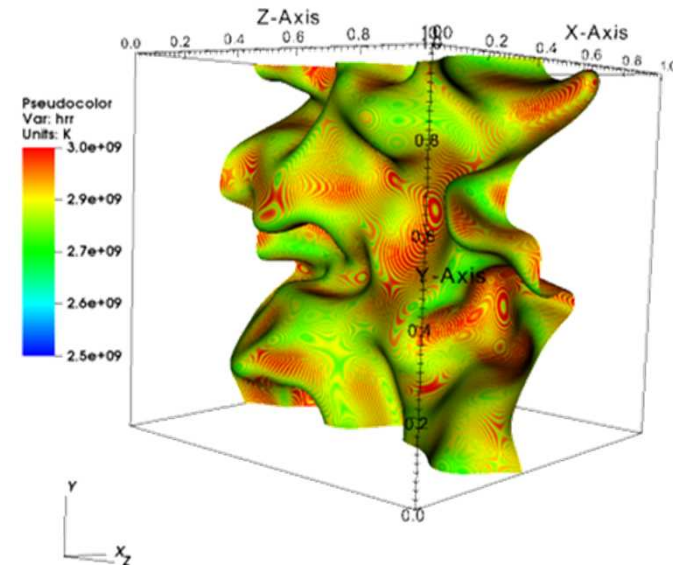
Downsample

$$\mathbf{C}^{(4)} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$

Pick single variable

$$\mathbf{C}^{(5)} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$

Pick single time step



Software: TuckerMPI



`git@gitlab.com:tensors/TuckerMPI.git`

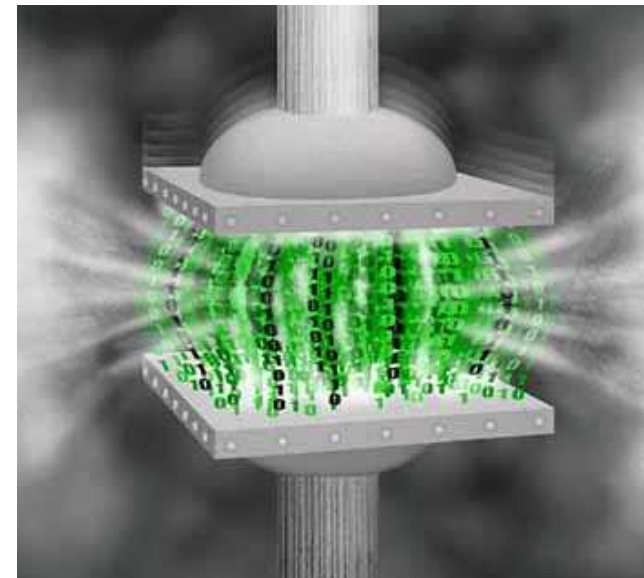
Alicia Klinvex, Woody Austin, Grey Ballard, Hemanth Kolla, Tammy Kolda

- Open source code for computing Tucker compression
- MPI/BLAS/LAPACK/C++11
- Still in development but available for testing
- Looking for new applications and users
- Interesting in partnering

Tensor Tucker Decomposition for Compression for Scientific Data



- First parallel implementation of Tucker decomposition
 - 5-way data using regular grid
 - Process 4TB data in < 1 min
- Up to 40000X compression on real-world data
 - Specify desired relative RMSE
 - Discovers latent multi-linear structure
 - Enables “smart” compression rather than discarding data that may be useful
- More work to do...
 - In situ computations
 - Real-time visualization for computational steering, etc.
 - Adaptive and non-uniform grids
 - Experimental data
 - Randomization
 - Extensive testing



<http://www.analyticbridge.com/profiles/blogs/how-much-is-big-data-compressible-an-interesting-theorem>

Tammy Kolda
tgkolda@sandia.gov

W. Austin, G. Ballard, and T. G. Kolda, *Parallel Tensor Compression for Large-Scale Scientific Data*, IPDPS'16 (arXiv:1510.06689)