



# *Parallel Tucker Compression for Large-scale Scientific Data*

Woody Austin  
Univ. Texas, Austin, TX

Grey Ballard and Tamara G. Kolda\*  
Sandia National Laboratories, Livermore, CA

SIAM Conference on Applied Linear Algebra (AL15)  
Atlanta, Georgia, USA  
October 29, 2015

Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.

# A Celebration of Dianne P. O'Leary



Sandia  
National  
Laboratories

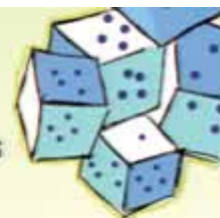


*PhD Advisor, Mentor, Role Model, Colleague, Friend*





# Dianne, Parallel Computing, & Tensors!



Unhappy student in Harland Glaz's scientific computing class. Didn't want Fortran and vector parallelization. Wanted C and MPI. He sent me to Dianne...

Reading class on parallel computing with this book:

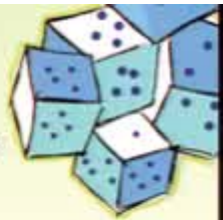


As a postdoc, the email that started me on tensors!

```
From: oleary@cs.umd.edu (Dianne O'Leary)
To: kolda@msr.epm.ornl.gov
CC: oleary@cs.umd.edu
Subject: reference
Received: from prost.cs.umd.edu (prost.cs.umd.edu [128.8.128.57]) by
msr.EPM.ORNL.GOV (8.8.3/8.8.3) with ESMTP id KAA14047 for
<kolda@msr.epm.ornl.gov>; Thu, 19 Feb 1998 10:57:21 -0500 (EST)
Date: Thu, 19 Feb 1998 15:57:19 GMT
Content-Type: text/plain; charset=utf-8
```

```
Leibovici has an article in LAA 269 p. 307 on one way to
extend the svd to 3-d and higher. I wonder if the sdd has
a similar extension?
```

Dianne



# A Tensor is an N-Way Array

Vector  
 $N = 1$



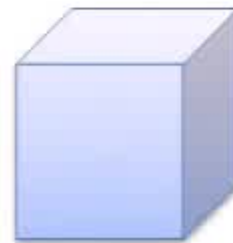
$\mathbf{x}$

Matrix  
 $N = 2$



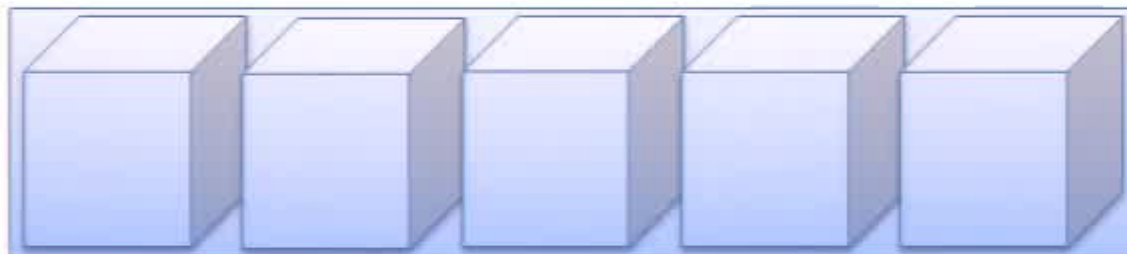
$\mathbf{X}$

3<sup>rd</sup>-Order Tensor  
 $N = 3$



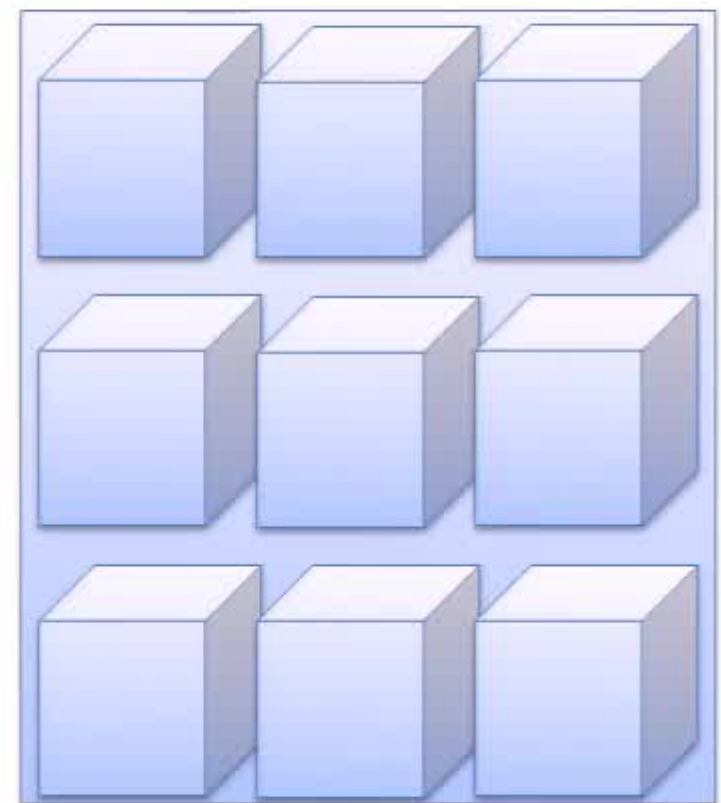
$\mathcal{X}$

4<sup>th</sup>-Order Tensor  
 $N = 4$



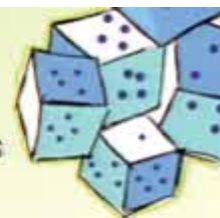
$\mathcal{X}$

5<sup>th</sup>-Order Tensor  
 $N = 5$



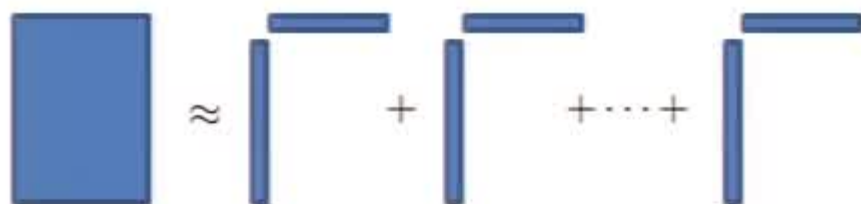
$\mathcal{X}$

# Tensor Decompositions are the New Matrix Decompositions

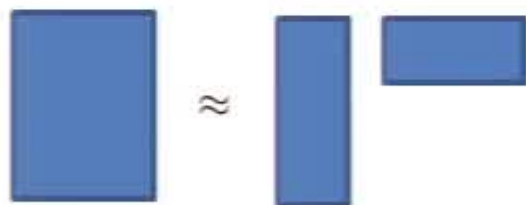


*Singular value decomposition (SVD), eigendecomposition (EVD), nonnegative matrix factorization (NMF), sparse SVD, etc.*

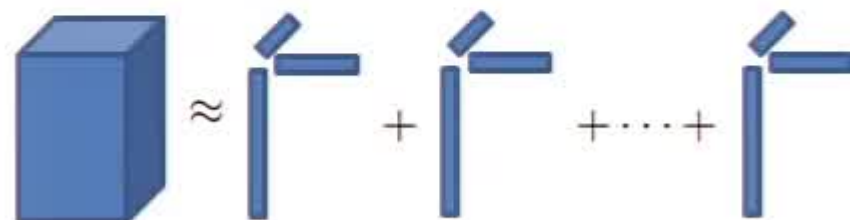
**Viewpoint 1:** Sum of outer products, useful for interpretation



**Viewpoint 2:** High-variance subspaces, useful for compression

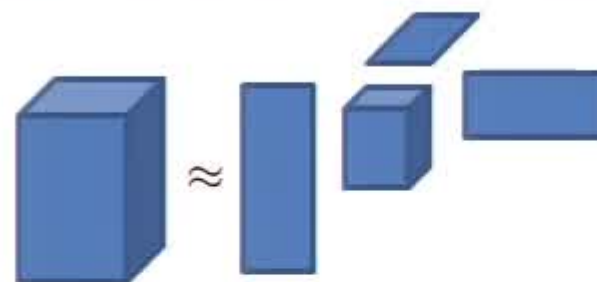


**CP Model:** Sum of d-way outer products, useful for interpretation



CANDECOMP, PARAFAC, Canonical Polyadic, CP

**Tucker Model:** Project onto high-variance subspaces to reduce dimensionality



HOSVD, Best Rank-(R1,R2,...,RN) decomposition

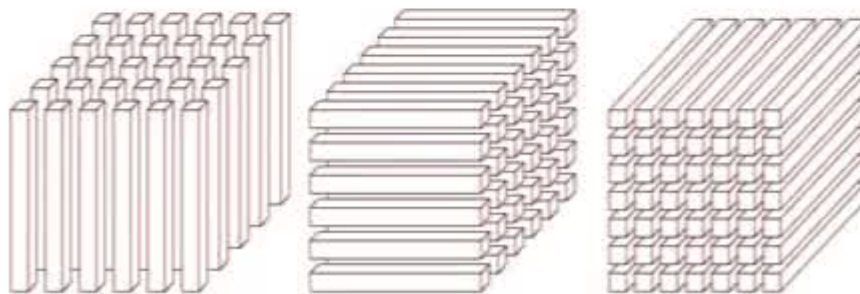
*Other models for compression include hierarchical Tucker and tensor train.*



# Tensor Fibers, Mode- $n$ Unfolding, and Mode- $n$ Multiplication

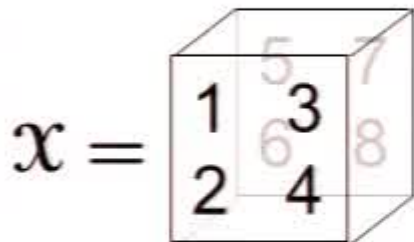


Tensor “mode- $n$  fibers” analogous to matrix rows and columns



Mode-1 Fibers    Mode-2 Fibers    Mode-3 Fibers

$\mathbf{X}_{(n)}$  denotes mode- $n$  unfolding, arranges mode- $n$  fibers as matrix columns



$$\mathbf{X}_{(1)} = \begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & 4 & 6 & 8 \end{bmatrix}$$

$$\mathbf{X}_{(2)} = \begin{bmatrix} 1 & 2 & 5 & 6 \\ 3 & 4 & 7 & 8 \end{bmatrix}$$

$$\mathbf{X}_{(3)} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$$

$$n\text{-rank}(\mathbf{X}) = \text{col-rank}(\mathbf{X}_{(n)})$$

Tensor-times-matrix (TTM) in mode- $n$  multiplies mode- $n$  fibers times matrix

$$I_1 \times \cdots \times I_n \times \cdots \times I_N \quad K \times I_n$$

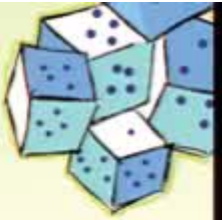
$$\mathbf{y} = \mathbf{X} \times_n \mathbf{U}$$

$$I_1 \times \cdots \times K \times \cdots \times I_N$$

Equivalent to matrix operation:

$$K \times \hat{I}_n \rightarrow \mathbf{Y}_{(n)} = \mathbf{U} \mathbf{X}_{(n)} \quad K \times I_n \quad I_n \times \hat{I}_n$$

$$I = \prod I_n, \quad \hat{I}_n = I / I_n$$



# Optimization Problem

$$\min_{\hat{\mathbf{x}}} \sum_{i_1 \dots i_N} (x_{i_1 \dots i_N} - \hat{x}_{i_1 \dots i_N})^2 \text{ subject to } \hat{\mathbf{x}} = \mathcal{G} \times \{ \mathbf{U}^{(n)} \}$$

Homework: (1) At an optimum, it must be the case that

$$\mathcal{G} = \mathbf{x} \times \{ \mathbf{U}^{(n)\top} \}$$

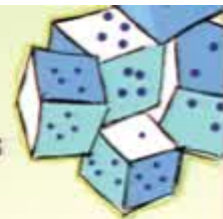
(2) The minimization problem above can be written as

$$\max_{\{ \mathbf{U}^{(n)} \}} \sum_{i_1 \dots i_N} g_{i_1 \dots i_N}^2 \text{ subject to } \mathcal{G} = \mathbf{x} \times \{ \mathbf{U}^{(n)\top} \}$$

$$\mathbf{U}^{(n)} = \arg \max_{\mathbf{U}} \| \mathbf{U}^\top \mathbf{Y}_{(n)} \|_F^2 \text{ subject to } \mathcal{Y} = \mathbf{x} \times \{ \mathbf{U}^{(m)\top} \}_{m \neq n}$$

Solution to (\*) is to choose  $\mathbf{U}^{(n)}$  to be the  $R_n$  leading left singular vectors of  $\mathbf{Y}_{(n)}$ .

# Sequentially Truncated HOSVD improves further



```
procedure ST-HOSVD( $\mathcal{X}$ ,  $\epsilon$ )
   $\mathcal{Y} \leftarrow \mathcal{X}$ 
  for  $n = 1, \dots, N$  do
     $\mathbf{S}^{(n)} \leftarrow \mathbf{Y}_{(n)} \mathbf{Y}_{(n)}^\top$ 
     $R_n \leftarrow \min R$  such that  $\sum_{r>R} \lambda_r(\mathbf{S}^{(n)}) \leq \epsilon^2 \|\mathcal{X}\|^2 / N$ 
     $\mathbf{U}^{(n)} \leftarrow$  leading  $R_n$  eigenvectors of  $\mathbf{S}^{(n)}$ 
     $\mathcal{Y} \leftarrow \mathcal{Y} \times_n \mathbf{U}^{(n)\top}$ 
  end for
   $\mathcal{G} \leftarrow \mathcal{Y}$ 
  return ( $\mathcal{G}$ ,  $\{\mathbf{U}^{(n)}\}$ )
end procedure
```

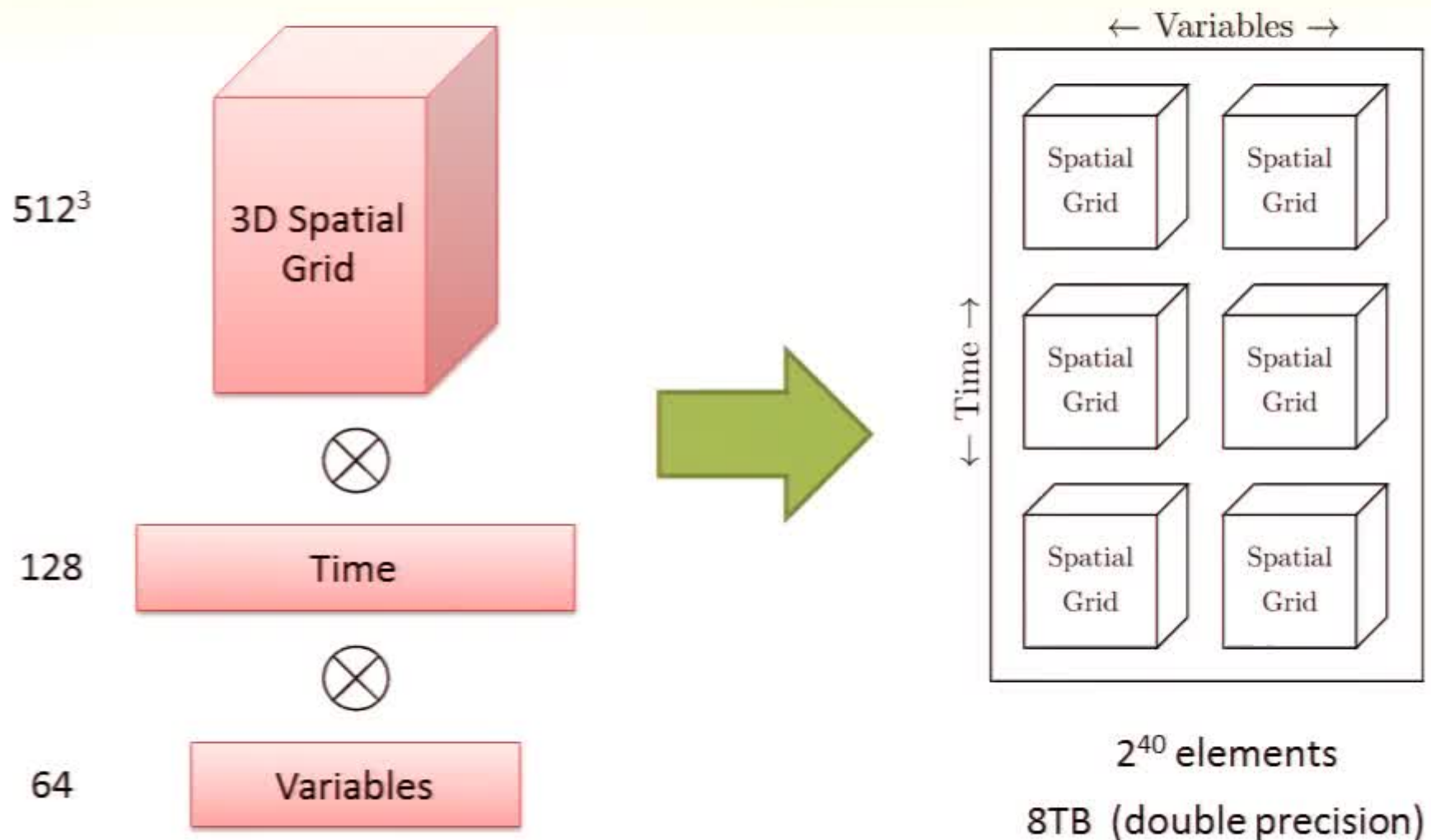
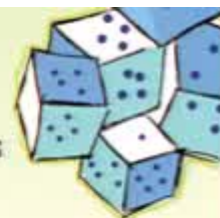
*Smaller at each step.*

$$\|\mathcal{X} - \hat{\mathcal{X}}\|^2 = \sum_{n=1}^N \left( \sum_{i=R_n+1}^{I_n} \lambda_i(\mathbf{S}^{(n)}) \right) \leq \epsilon^2 \|\mathcal{X}\|^2$$

Vannieuwenhoven, Vandebril, and Meerbergen (2012)



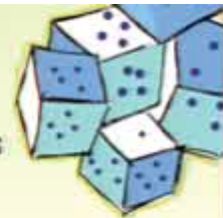
# Tensors in Scientific Applications are Huge, Need Parallel Methods



# Key Kernels in ST-HOSVD are TTM and Gram



```
procedure ST-HOSVD( $\mathcal{X}$ ,  $\epsilon$ )
   $\mathcal{Y} \leftarrow \mathcal{X}$ 
  for  $n = 1, \dots, N$  do
    Gram  $\mathbf{S}^{(n)} \leftarrow \mathbf{Y}_{(n)} \mathbf{Y}_{(n)}^\top$ 
     $R_n \leftarrow \min R$  such that  $\sum_{r>R} \lambda_r(\mathbf{S}^{(n)}) \leq \epsilon^2 \|\mathcal{X}\|^2 / N$ 
     $\mathbf{U}^{(n)} \leftarrow$  leading  $R_n$  eigenvectors of  $\mathbf{S}^{(n)}$ 
    TTM  $\mathcal{Y} \leftarrow \mathcal{Y} \times_n \mathbf{U}^{(n)\top}$ 
  end for
   $\mathcal{G} \leftarrow \mathcal{Y}$ 
  return ( $\mathcal{G}$ ,  $\{ \mathbf{U}^{(n)} \}$ )
end procedure
```



# Unfolded Tensor Distribution

Global Tensor Size:  $J_1 \times J_2 \times \cdots \times J_N$ ,  $J = \prod J_n$ ,  $\hat{J}_n = J/J_n$

Processor Grid Size:  $P_1 \times P_2 \times \cdots \times P_N$ ,  $P = \prod P_n$ ,  $\hat{P}_n = P/P_n$

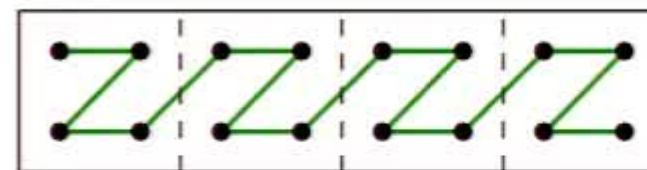
Global Unfolded Tensor:  $J_n \times \hat{J}_n$

Processor Grid:  $P_n \times \hat{P}_n$

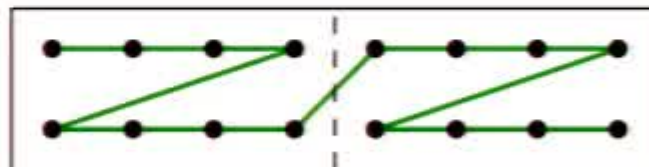
Local Layout: 2 x 2 x 2 x 2



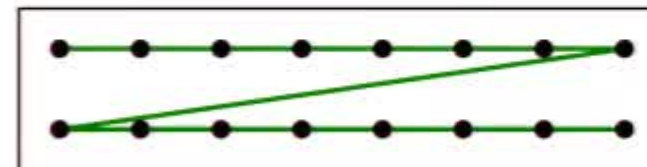
$n = 1$



$n = 2$

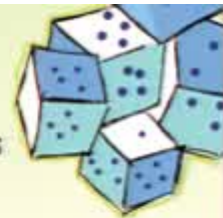


$n = 3$



$n = 4$





# Parallel Gram

$$\mathbf{S} = \mathbf{Y}_{(n)} \mathbf{Y}_{(n)}^T$$

procedure GRAM( $\mathbf{y}, n$ )

myProcID  $\leftarrow (p_1, p_2, \dots, p_N)$

myProcCol  $\leftarrow (p_1, \dots, p_{n-1}, *, p_{n+1}, \dots, p_N)$

myProcRow  $\leftarrow (*, \dots, *, p_k, *, \dots, *)$

$\mathbf{V}^{[p_n]} \leftarrow \bar{\mathbf{Y}}_{(n)} \bar{\mathbf{Y}}_{(n)}^T$

for  $i = 1$  to  $P_n - 1$  do

$j \leftarrow (p_n - i) \bmod P_n$

$k \leftarrow (p_n + i) \bmod P_n$

Send  $\bar{\mathbf{y}}$  to process  $(p_1, \dots, p_{n-1}, j, \dots, p_N)$

Receive  $\mathbf{W}$  from process  $(p_1, \dots, p_{n-1}, k, \dots, p_N)$

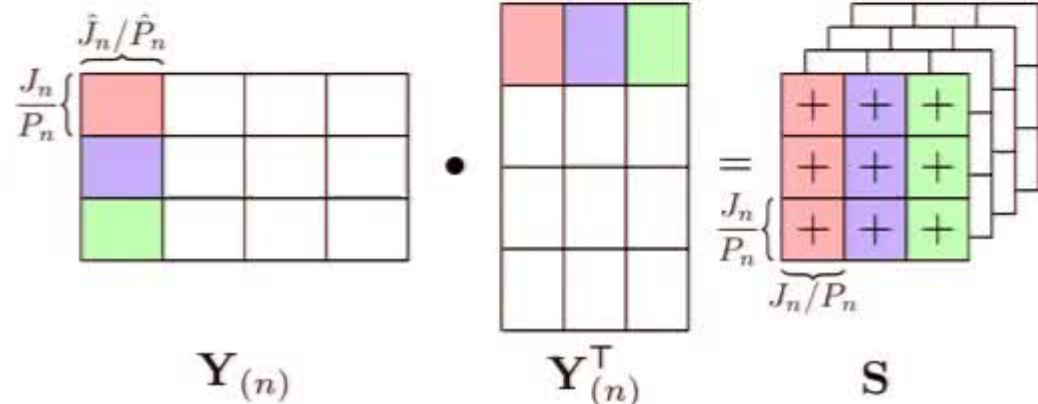
$\mathbf{V}^{[k]} \leftarrow \bar{\mathbf{Y}}_{(n)} \mathbf{W}^T$

end for

$\bar{\mathbf{S}} = \text{All-Reduce}(\mathbf{V}, \text{myProcRow})$

return  $\bar{\mathbf{S}}$

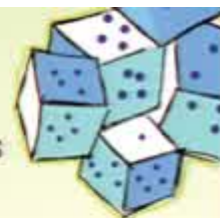
end procedure



$$C_{\text{GRAM}} = \gamma 2J_n J/P + 2(P_n - 1)(\alpha + \beta J/P) + 2\alpha \log \hat{P}_n + 2\beta(\hat{P}_n - 1)J_n^2/P$$

$$M_{\text{GRAM}} = J/P + J/P + J_n^2/P_n + J_n^2/P_n$$

# Application Results: Compression versus Accuracy



Ranks depend on error:

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \sqrt{\sum_{n=1}^N \left( \sum_{i=R_n+1}^{I_n} \lambda_i(\mathbf{S}^{(n)}) \right)} \leq \epsilon \|\mathbf{x}\|$$

Compression ratio:

$$C = \prod_{k=1}^N I_n / \left( \prod_{k=1}^N R_n + \sum_{k=1}^N I_n R_n \right).$$

*Simulation of an autoignitive premixture of air and ethanol in Homogeneous Charge Compression Ignition*

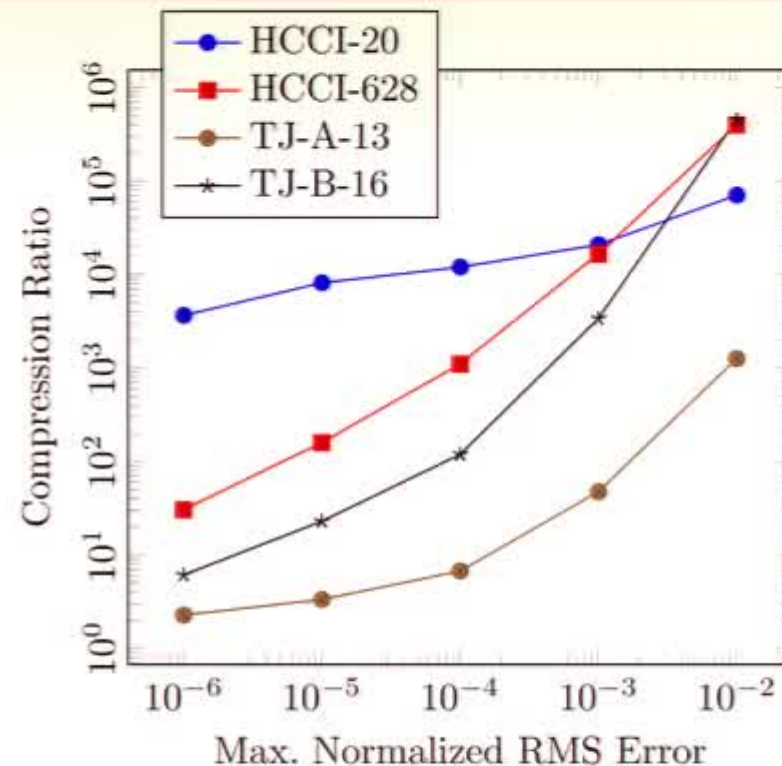
HCCI-628: 672 x 672 x 33 x 628, 72 GB

*Temporally-evolving planar slot jet flame with DME (dimethyl ether) as the fuel*

TJ-A-13: 300 x 500 x 240 x 35 x 13, 122 GB

TJ-B-16: 460 x 700 x 360 x 35 x 16, 512 GB

Thanks to Hemanth Kolla and Ankit Bhagatwala for combustion application data, from Sandia's S3D direct numerical simulation code

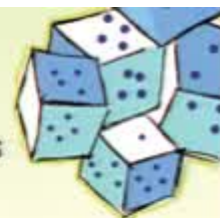


$\epsilon = 10^{-5}$

Dataset	Reduced Size	Max. Elem. Error	Comp. Ratio
HCCI-1	(16, 16, 4, 1)	3.6e-5	573
HCCI-20	(20, 18, 6, 5)	2.0e-4	7083
HCCI-628	(192, 183, 16, 104)	1.2e-3	139
TJ-A-1	(257, 139, 186, 20, 1)	1.7e-3	9
TJ-A-13	(300, 209, 240, 25, 13)	3.2e-3	3



# Sample Results for one Species in HCCI: Error is Negligible

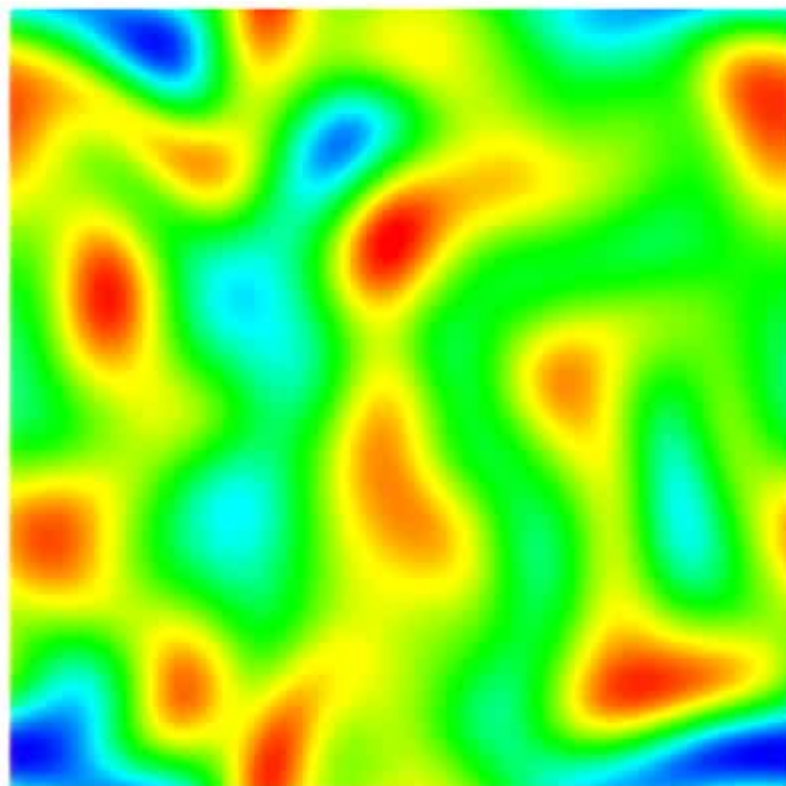


Compression: 586

Original  $\mathcal{X}$

672 x 672 x 33 x 8

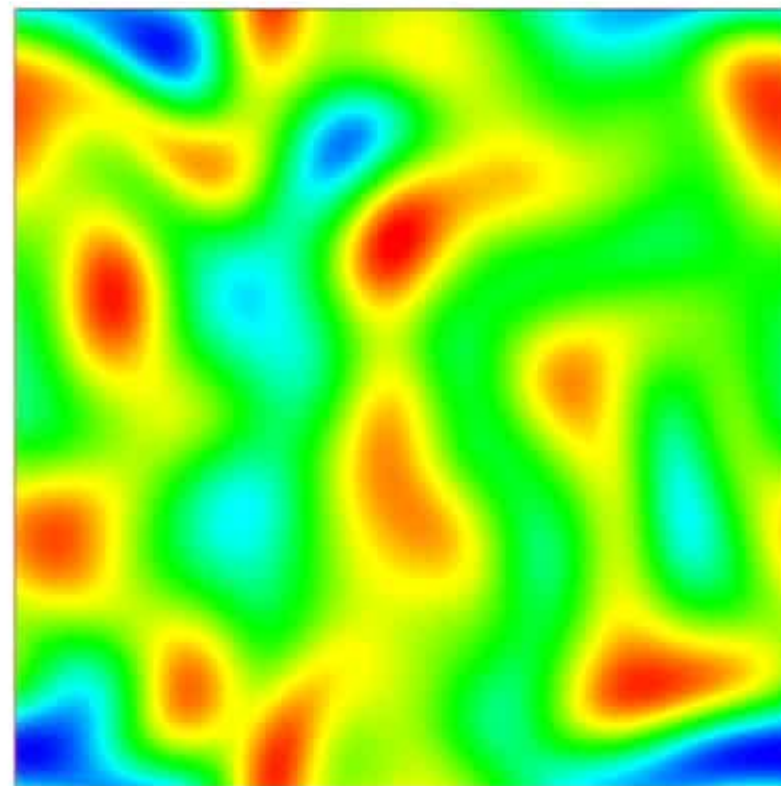
Pseudocolor  
Var: T  
Units: K  
1019.  
966.0  
913.3  
860.6  
808.0  
Max: 1019.  
Min: 808.0



Recovered  $\hat{\mathcal{X}}$

48 x 48 x 20 x 3

Pseudocolor  
Var: T  
Units: K  
1019.  
966.0  
913.3  
860.7  
808.0  
Max: 1019.  
Min: 808.0

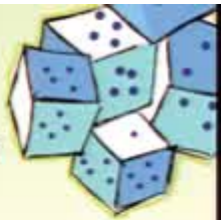


910MB  $\Rightarrow$  1.5MB

$$\frac{\|\mathcal{X} - \hat{\mathcal{X}}\|}{\|\mathcal{X}\|} = 3.08 \times 10^{-8}$$



# Sample Results for “Derived” Quantity: Error is Negligible

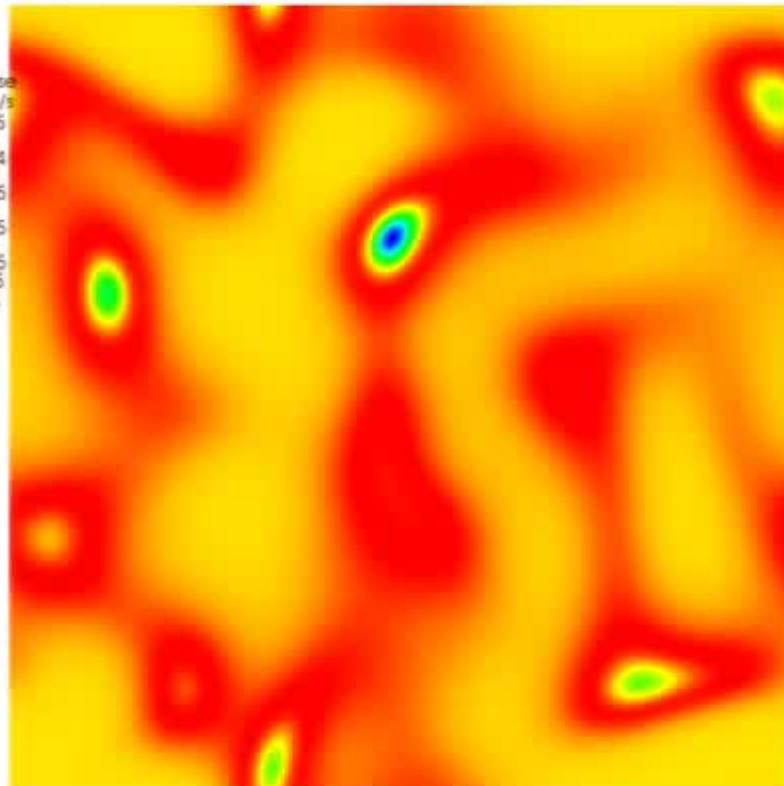


Compression: 586

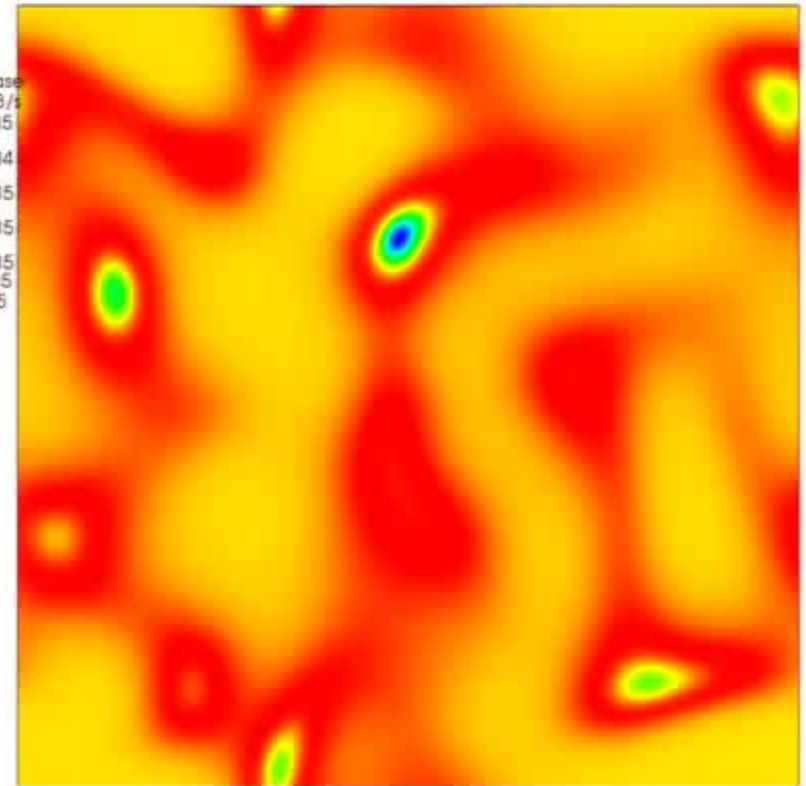
Original  $\mathcal{X}$   
672 x 672 x 33 x 8

Recovered  $\hat{\mathcal{X}}$   
48 x 48 x 20 x 3

Pseudocolor  
Var: heat\_release  
Units: erg/cm<sup>3</sup>/s  
1.192e+05  
-1.596e+04  
-1.511e+05  
-2.862e+05  
-4.214e+05  
Max: 1.192e+05  
Min: -4.214e+05



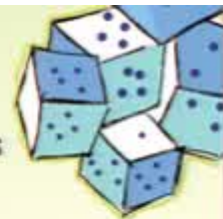
Pseudocolor  
Var: heat\_release  
Units: erg/cm<sup>3</sup>/s  
1.188e+05  
-1.668e+04  
-1.522e+05  
-2.877e+05  
-4.232e+05  
Max: 1.188e+05  
Min: -4.232e+05



910MB  $\Rightarrow$  1.5MB

$$\frac{\|\mathcal{X} - \hat{\mathcal{X}}\|}{\|\mathcal{X}\|} = 3.08 \times 10^{-8}$$

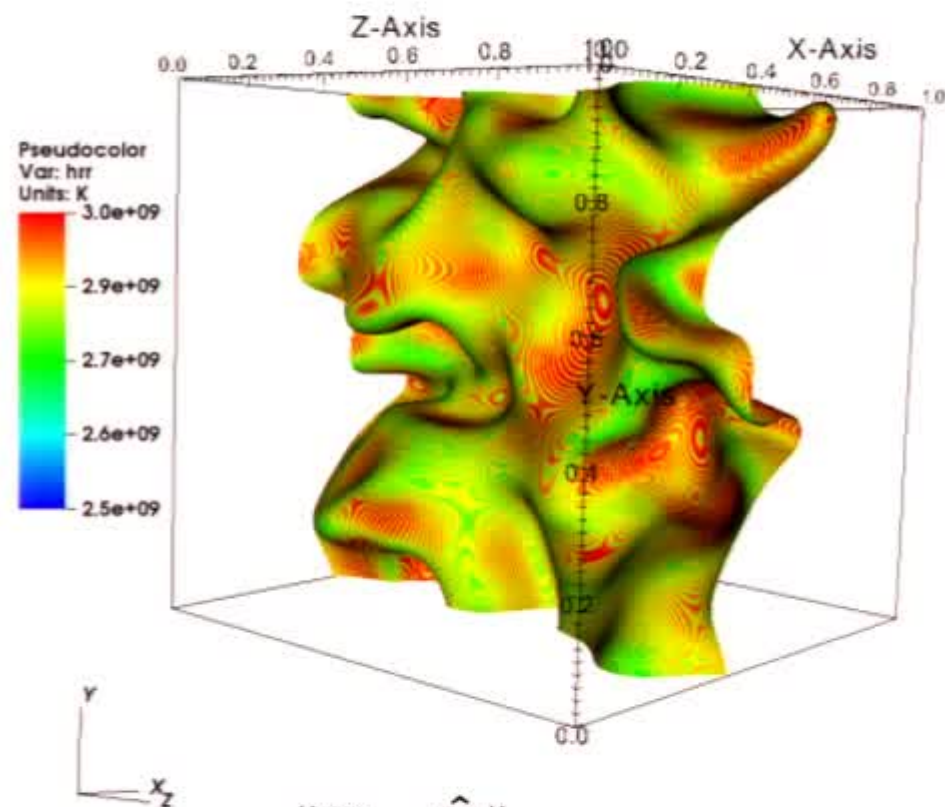
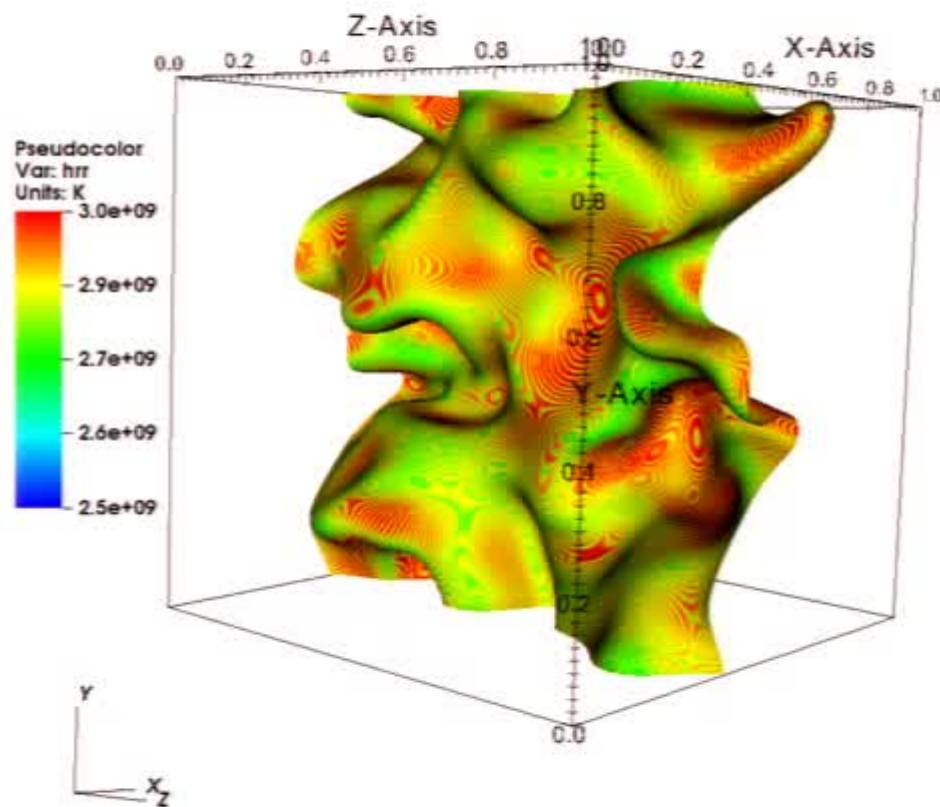
# Sample Results for one Species in 3D HCCI: Error is Negligible



Original  $\mathcal{X}$   
500 x 500 x 500 x 11 x 5

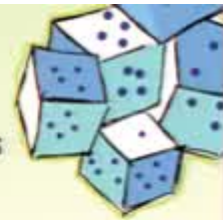
Compression: 1000

Recovered  $\hat{\mathcal{X}}$   
50 x 50 x 50 x 11 x 5



52 GB  $\Rightarrow$  52 MB

$$\frac{\|\mathcal{X} - \hat{\mathcal{X}}\|}{\|\mathcal{X}\|} = 3.3 \times 10^{-9}$$



# Partial Reconstruction

Reconstruction requires as much space  
as the original data!

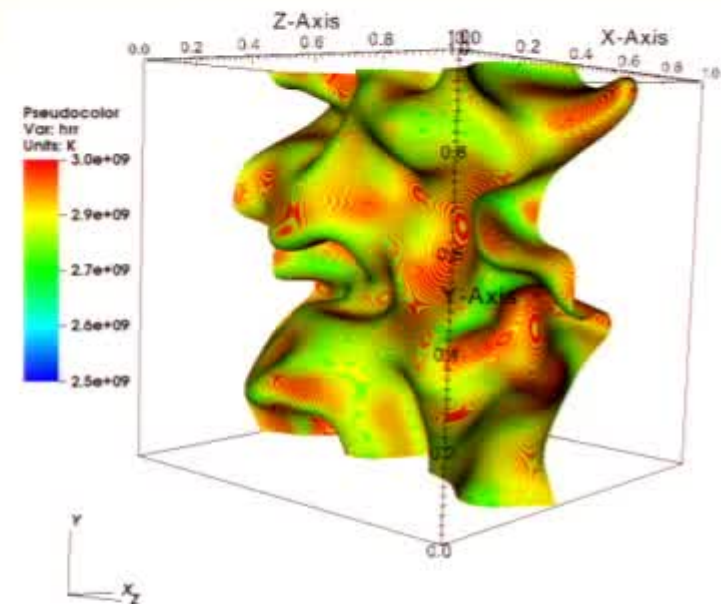
$$\hat{\mathcal{X}} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \times_4 \mathbf{U}^{(4)} \times_5 \mathbf{U}^{(5)}$$

$$I_1 \times I_2 \times I_3 \times I_4 \times I_5$$

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

$$\bar{\mathcal{X}} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \times_4 \underbrace{\mathbf{U}^{(4)} \mathbf{e}_k}_{\text{Pick out } k\text{th species}} \times_5 \underbrace{\mathbf{U}^{(5)} \mathbf{e}_l}_{\text{Pick out } l\text{th time}}$$

$$I_1 \times I_2 \times I_3 \times 1 \times 1$$





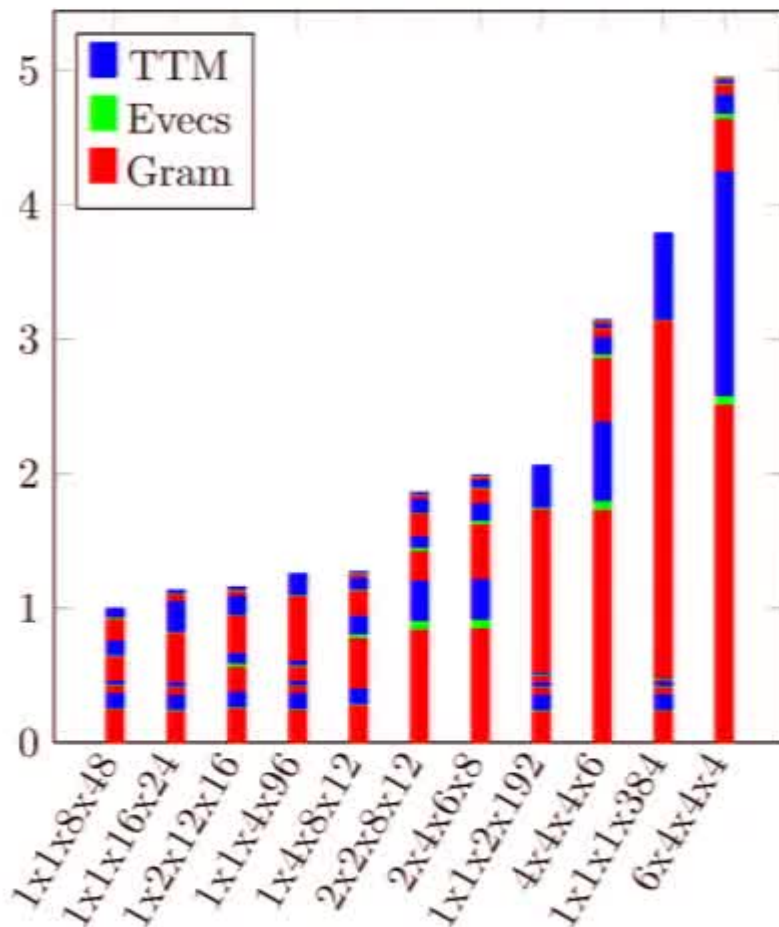
# Parameter Choices: Processor Grid Configuration & Mode Order



## Processor Grid Configuration

$I: 384 \times 384 \times 384 \times 384$

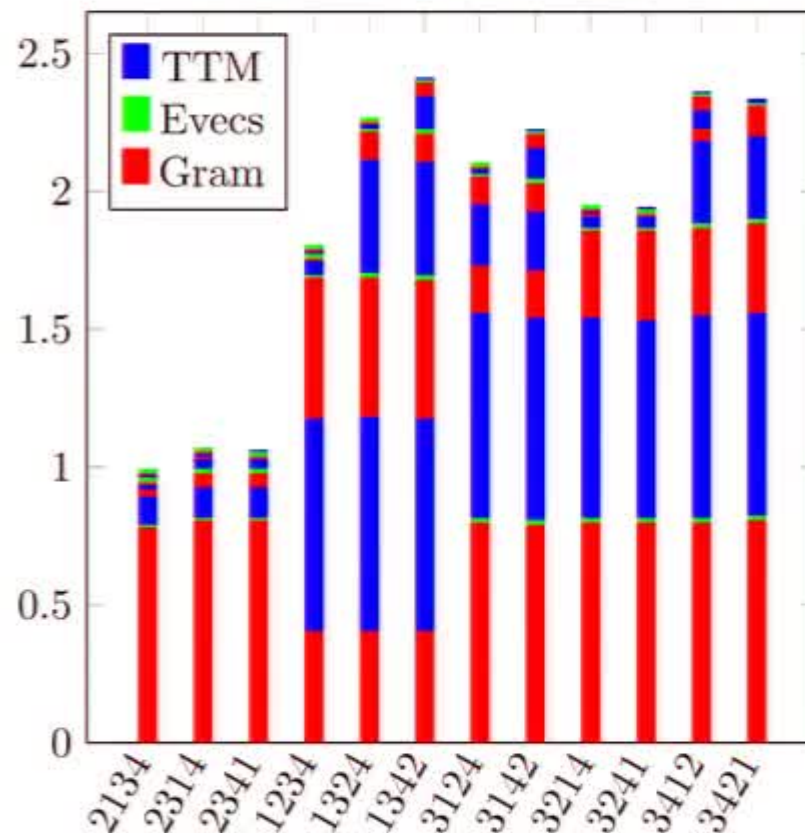
$R: 96 \times 96 \times 96 \times 96$



## Mode Order

$I: 25 \times 250 \times 250 \times 250$

$R: 10 \times 10 \times 100 \times 100$



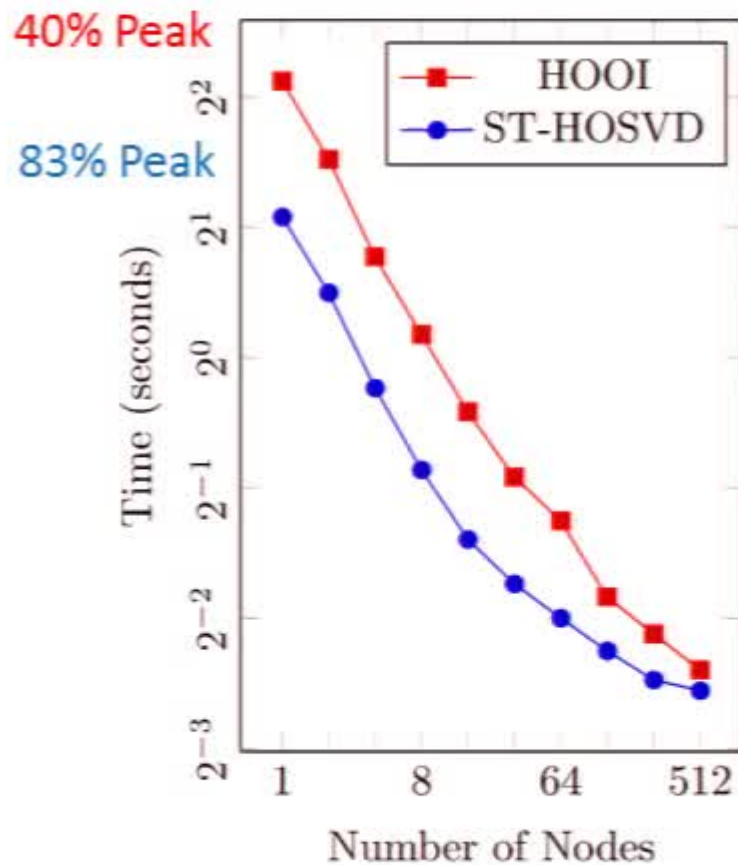


# Strong & Weak Scaling

## Strong Scaling with $24 \times 2^k$ processors

$I: 500 \times 300 \times 240 \times 35$

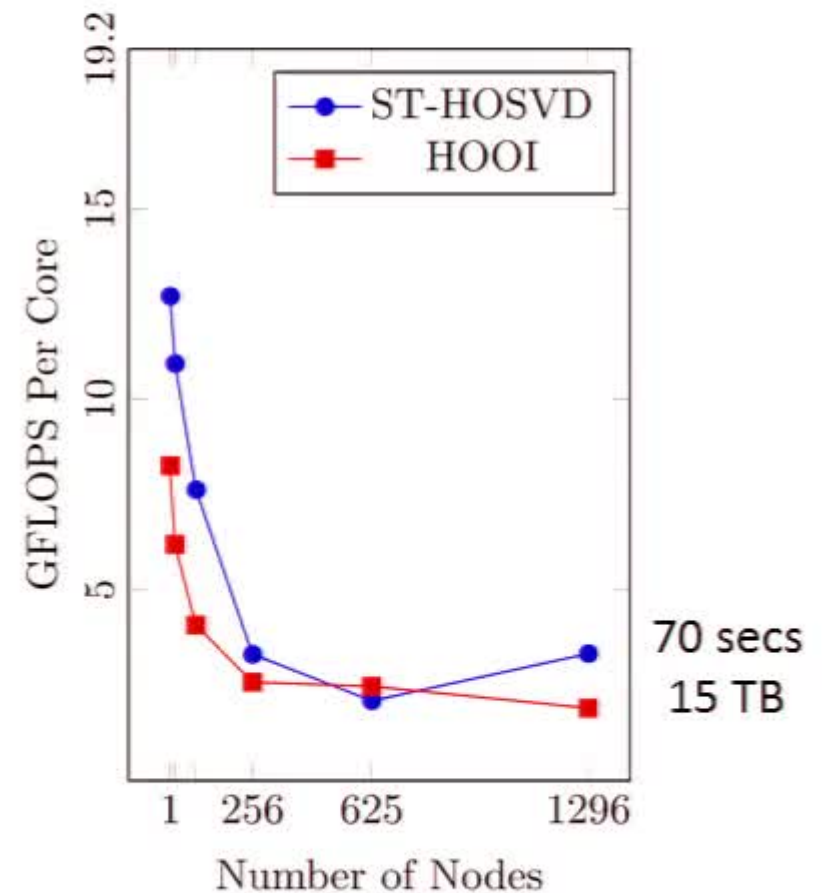
$R: 42 \times 115 \times 81 \times 19$



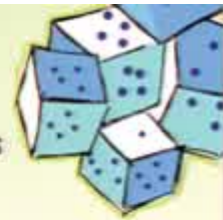
## Weak Scaling with $24 \times k^4$ processors

$I: 200k \times 200k \times 200k \times 200k$

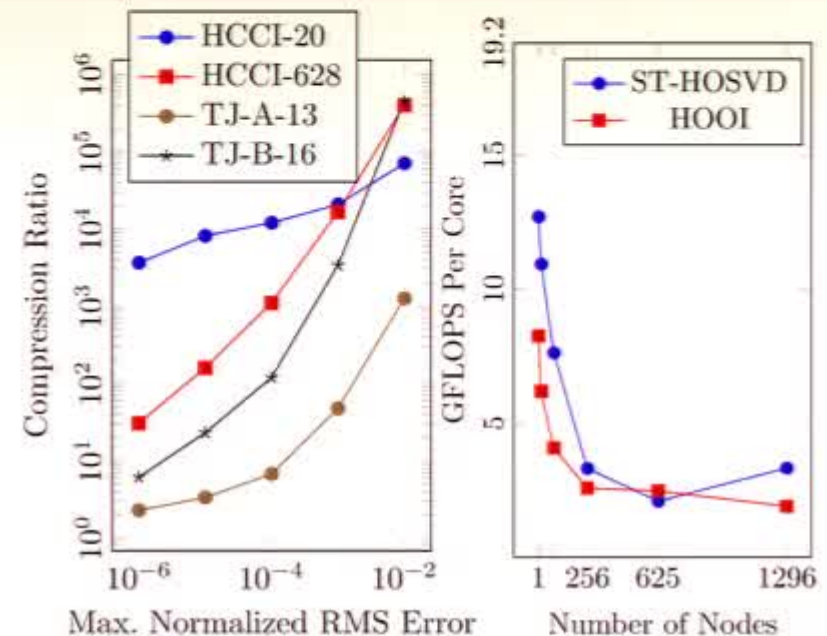
$R: 20k \times 20k \times 20k \times 20k$



# Parallel Tucker Compression



- First-ever implementation of distributed-memory parallel Tucker decomposition
  - Avoids unnecessary data permutations
- Up to  $10^5$  compression on real-world data with minimal loss in accuracy
- Scales well – achieving 17% of peak on 1000s of processors
- Future work
  - Detailed application studies
  - Use QR instead of Gram



For more information:  
Tammy Kolda,  
[tgkolda@sandia.gov](mailto:tgkolda@sandia.gov)

W. Austin, G. Ballard, and T. G. Kolda, *Parallel Tensor Compression for Large-Scale Scientific Data*, arXiv:1510.06689, October 2015, submitted for publication



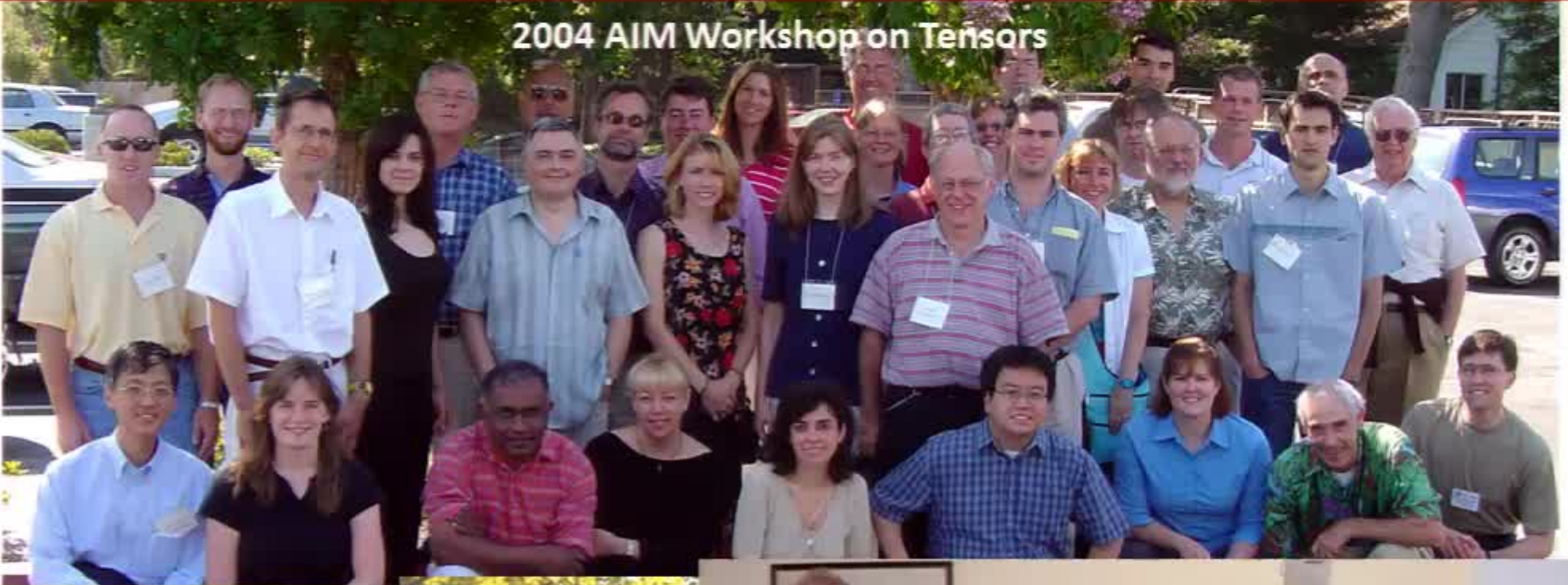


Sandia  
National  
Laboratories



# In Memory of... Carla Martin

2004 AIM Workshop on Tensors



Householder 2008