

# *Parallel Tensor Compression for Large-Scale Scientific Data*

**Woody Austin**  
Univ. Texas, Austin, TX

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

DMML Workshop  
Berkeley CA  
October 23, 2015

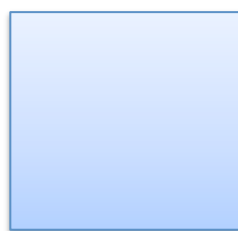
# A tensor is an N-way array

Vector  
N = 1



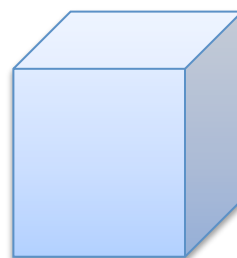
$\mathcal{X}$

Matrix  
N = 2



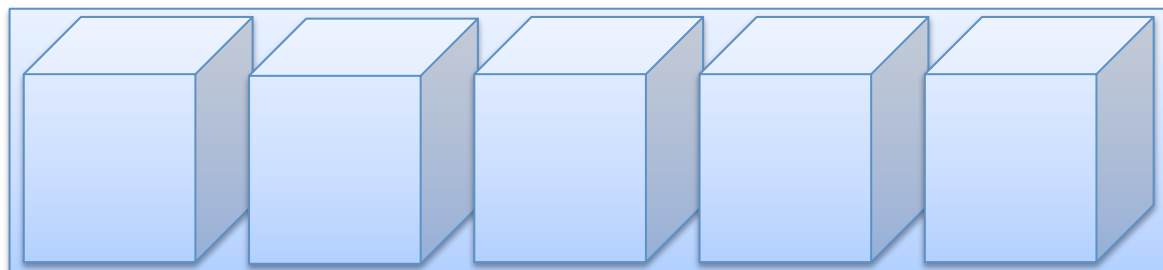
$\mathcal{X}$

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



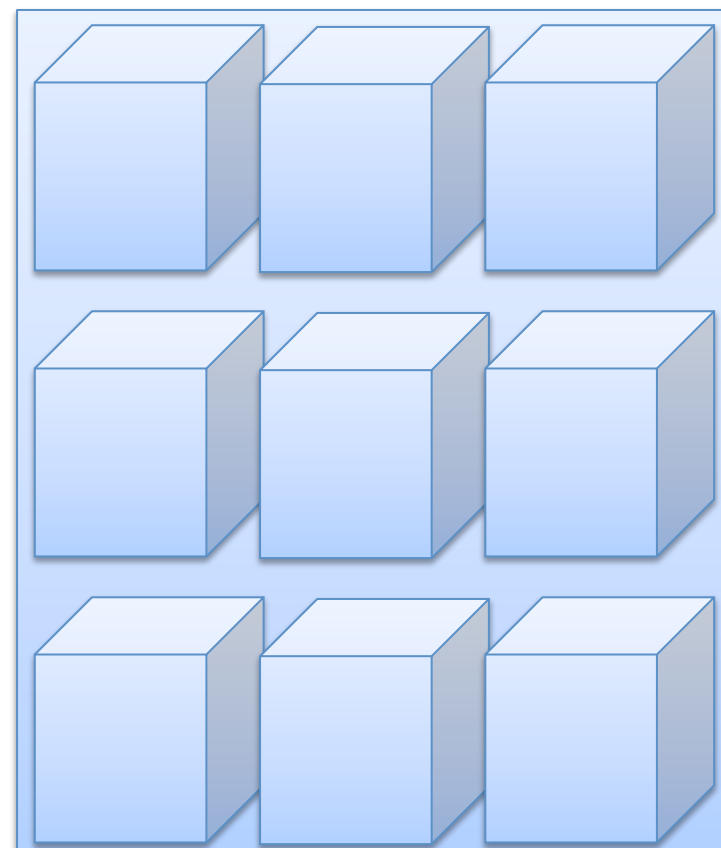
$\mathcal{X}$

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



$\mathcal{X}$

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

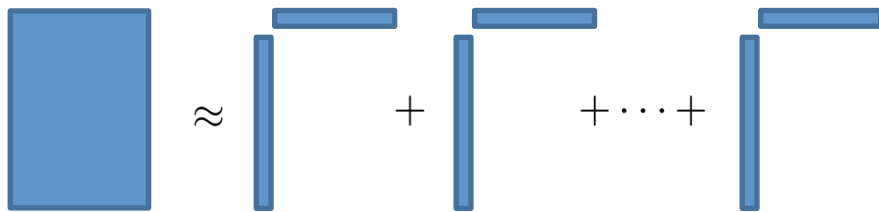


$\mathcal{X}$

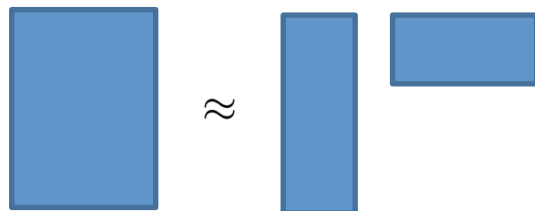
# Tensor decompositions are the new matrix decompositions

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

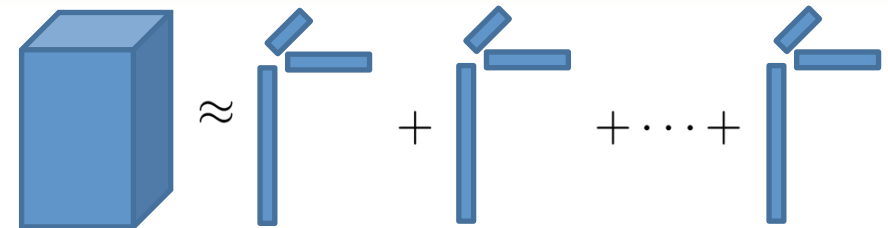
**Viewpoint 1:** Sum of outer products,  
useful for interpretation



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

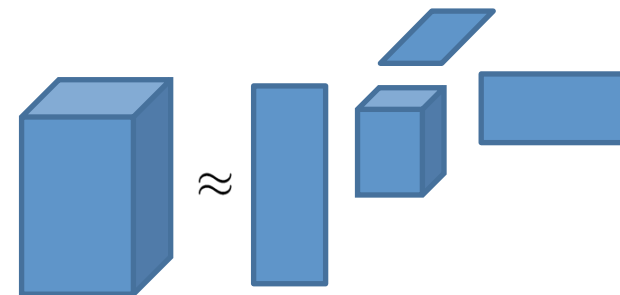


**CP Model:** Sum of d-way outer products,  
useful for interpretation



**CANDECOMP, PARAFAC, Canonical Polyadic, CP**

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

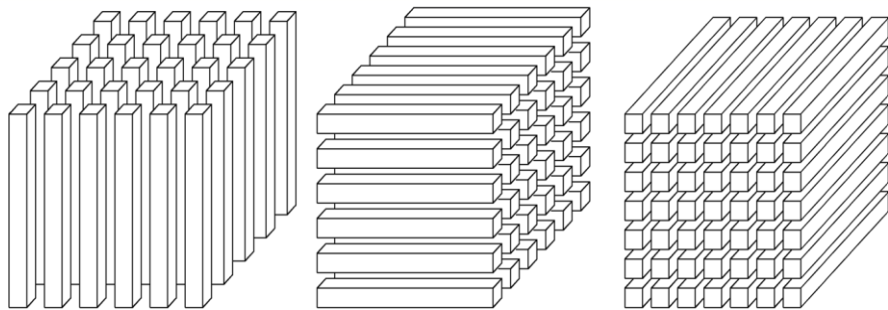


**HOSVD, Best Rank-(R1,R2,...,RN) decomposition**

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

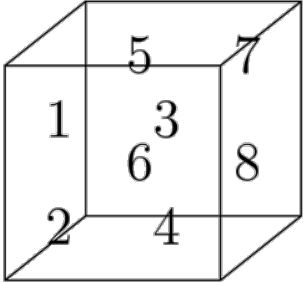
# Tensor fibers, mode- $n$ unfolding, and mode- $n$ Multiplication

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

$\mathbf{x} =$  

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

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

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

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

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

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

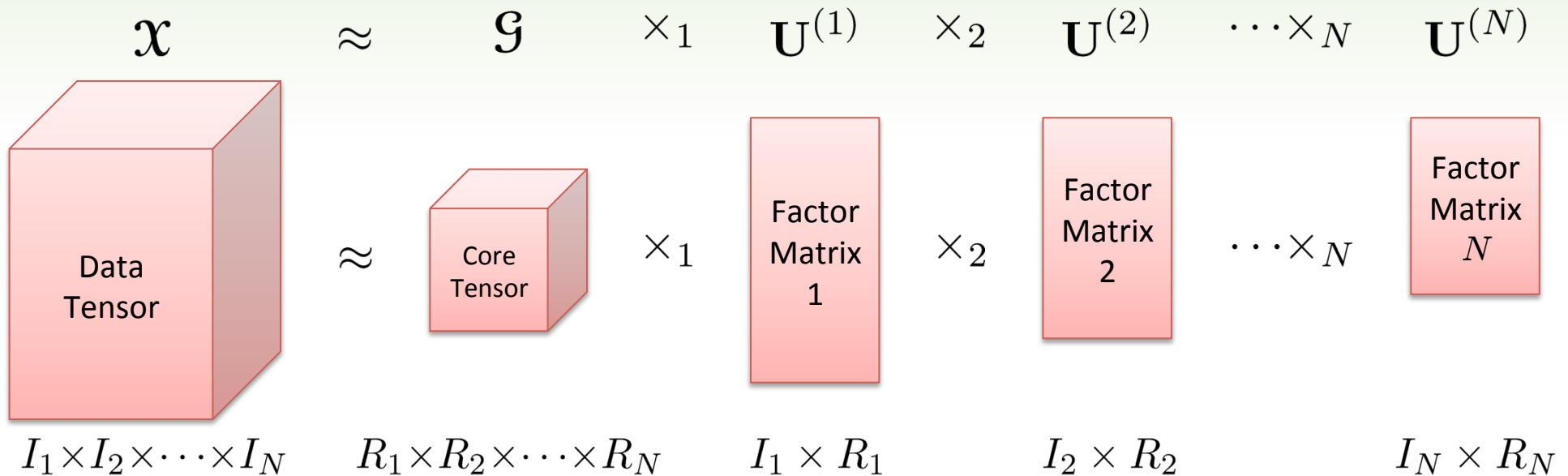
$$I_1 \times \cdots \times K \times \cdots \times I_N$$

Equivalent to matrix operation:

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

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

# Tucker decomposition



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

WLOG, assume  $\mathbf{U}^{(n)}$  has orthogonal columns for all  $n$ .

If  $R_n \geq \text{rank}(\mathbf{X}_{(n)})$  for all  $n$ , then decomposition is exact. Else, it's lossy.

Tucker (1966); Kapteyn, Neudecker, Wansbeek (1986)

# Optimization problem

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

Couple Facts: (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

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

$$\max_{\mathbf{U}^{(n)}} \| \mathbf{U}^{(n)\top} \mathbf{W}_{(n)} \|_F^2 \text{ subject to } \mathcal{W} = \{ \mathbf{U}^{(m)\top} \}_{m \neq n} \quad (*)$$

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

# Truncated Higher-Order SVD (HOSVD) is a sequence of truncated SVDs

```
1: procedure T-HOSVD( $\mathcal{X}$ ,  $\{R_n\}$ )
2:   for  $n = 1, \dots, N$  do
3:      $\mathbf{U}^{(n)} \leftarrow$  leading  $R_n$  left singular vectors of  $\mathbf{X}_{(n)}$ 
4:   end for
5:    $\mathcal{G} \leftarrow \mathcal{X} \times \{\mathbf{U}^{(n)\top}\}$ 
6:   return  $(\mathcal{G}, \{\mathbf{U}^{(n)}\})$ 
7: end procedure
```

Also known as “Tucker1” method.

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

# Truncated HOSVD is a sequence of truncated SVDs

```
1: procedure T-HOSVD( $\mathcal{X}$ ,  $\{R_n\}$ )
2:   for  $n = 1, \dots, N$  do
3:      $\mathbf{S}^{(n)} = \mathbf{X}_{(n)}\mathbf{X}_{(n)}^T$ 
4:      $\mathbf{U}^{(n)} \leftarrow$  leading  $R_n$  eigenvectors of  $\mathbf{S}^{(n)}$ 
5:   end for
6:    $\mathcal{G} \leftarrow \mathcal{X} \times \{\mathbf{U}^{(n)T}\}$ 
7:   return ( $\mathcal{G}$ ,  $\{\mathbf{U}^{(n)}\}$ )
8: end procedure
```

Also known as “Tucker1” method.

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



# Sequentially Truncated HOSVD improves further

```
procedure ST-HOSVD( $\mathcal{X}$ ,  $\{R_n\}$ )  
   $\mathcal{Y} \leftarrow \mathcal{X}$   
  for  $n = 1, \dots, N$  do  
     $\mathbf{S}^{(n)} \leftarrow \mathbf{Y}_{(n)} \mathbf{Y}_{(n)}^\top$   
     $\mathbf{U}^{(n)} \leftarrow$  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**



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

Vannieuwenhoven, Vandebril, and Meerbergen (2012)

# Higher-Order Orthogonal Iteration (HOOI) improves again

```
procedure HOOI( $\mathcal{X}$ ,  $\{R_n\}$ )
  ( $\mathcal{G}$ ,  $\{\mathbf{U}^{(n)}\}$ ) = ST-HOSVD( $\mathcal{X}$ ,  $\{R_n\}$ )
  repeat
    for  $n = 1, \dots, N$  do
       $\mathcal{Y} \leftarrow \mathcal{X} \times \{\mathbf{U}^{(m)\top}\}_{m \neq n}$ 
       $\mathbf{S}^{(n)} \leftarrow \mathbf{Y}_{(n)} \mathbf{Y}_{(n)}^\top$ 
       $\mathbf{U}^{(n)} \leftarrow$  leading  $R_n$  eigenvectors of  $\mathbf{S}^{(n)}$ 
    end for
     $\mathcal{G} \leftarrow \mathcal{Y} \times_N \mathbf{U}^{(N)\top}$ 
  until the quantity  $\|\mathcal{G}\|^2$  ceases to increase
  return ( $\mathcal{G}$ ,  $\{\mathbf{U}^{(n)}\}$ )
end procedure
```

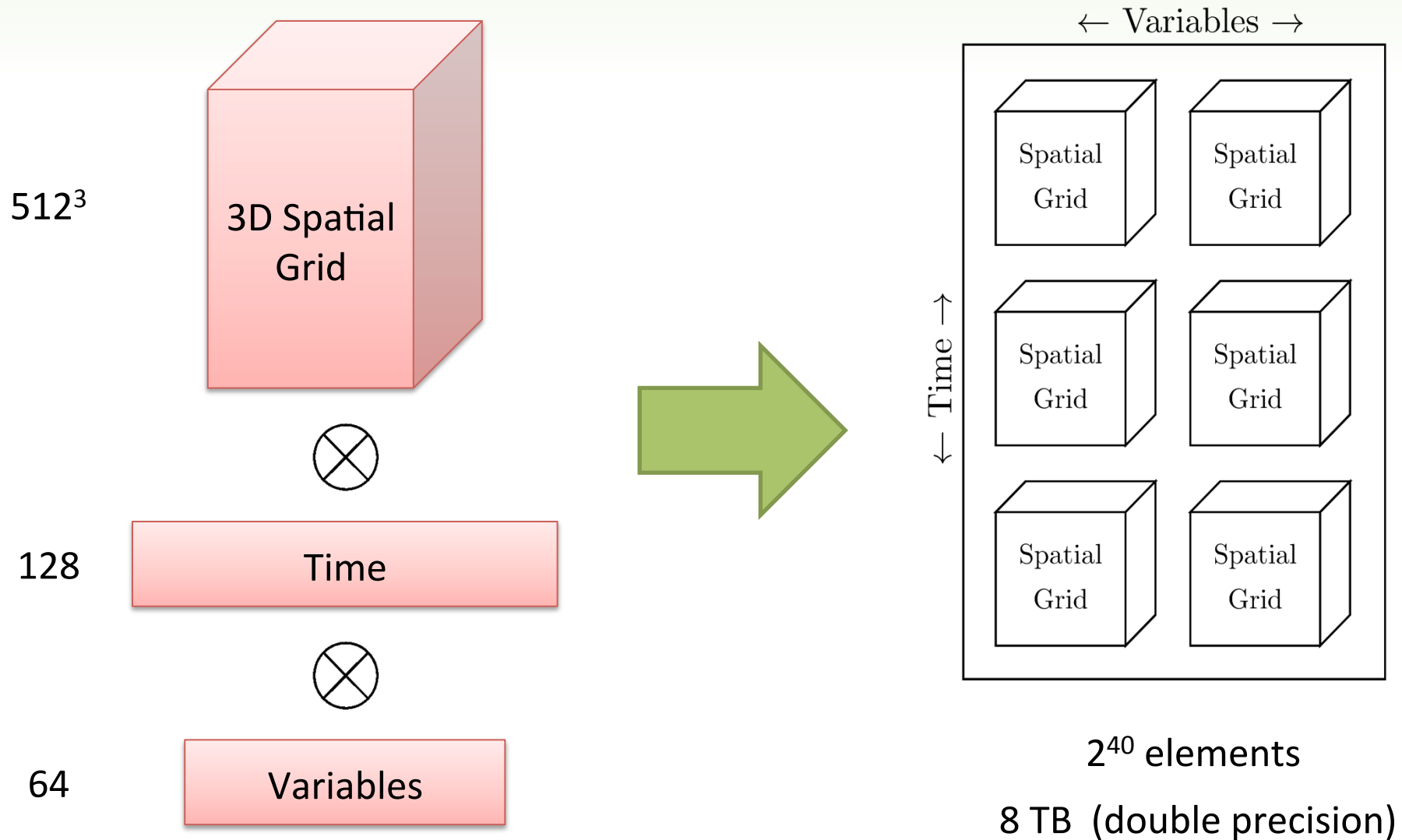
# Key kernels are TTM and Gram

```

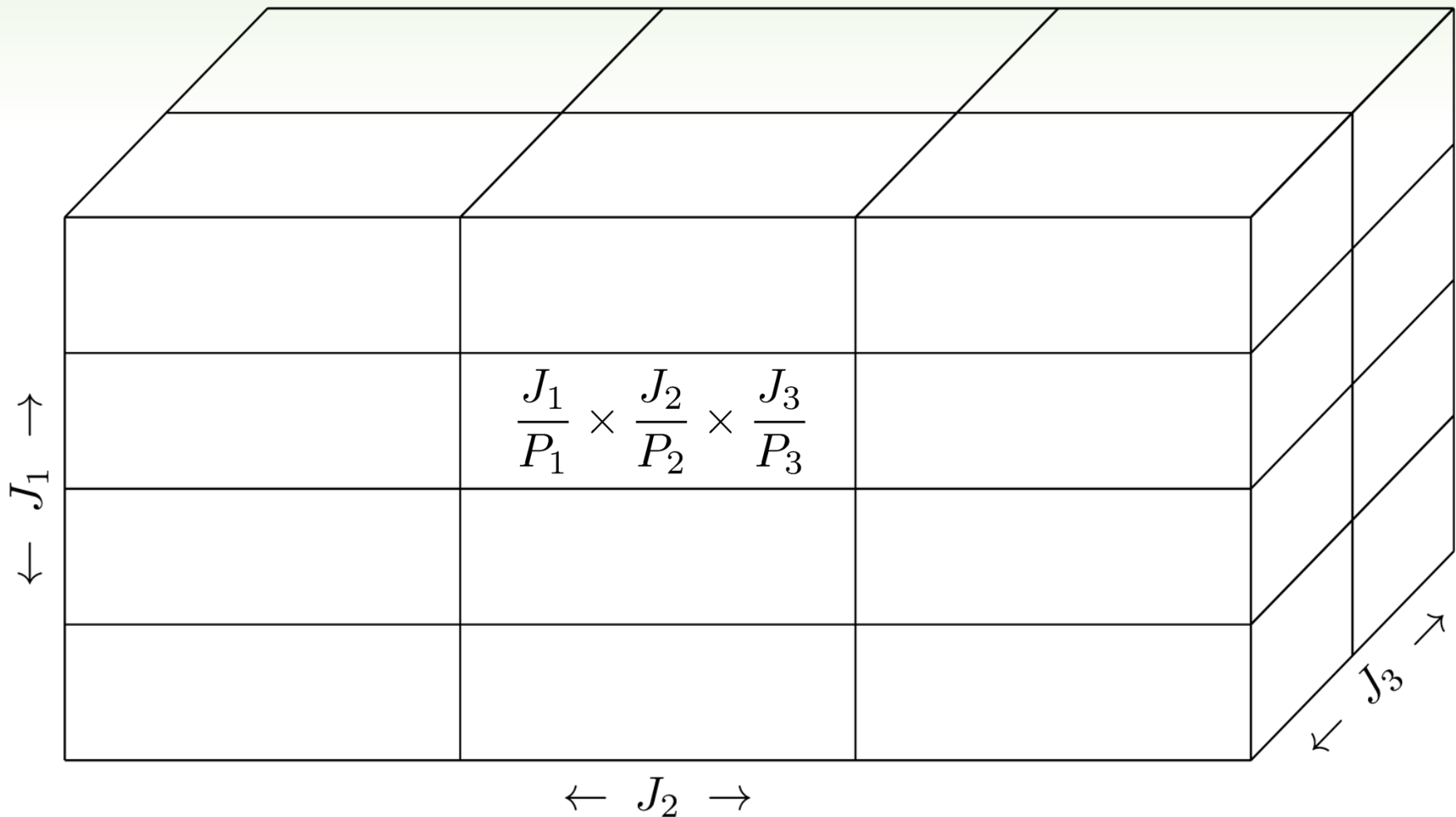
procedure ST-HOSVD( $\mathcal{X}$ ,  $\{R_n\}$ )
   $\mathcal{Y} \leftarrow \mathcal{X}$ 
  for  $n = 1, \dots, N$  do
    Gram  $\mathbf{S}^{(n)} \leftarrow \mathbf{Y}_{(n)} \mathbf{Y}_{(n)}^\top$ 
     $\mathbf{U}^{(n)} \leftarrow$  leading  $R_n$  eigenvectors of  $\mathbf{S}^{(n)}$ 
    TTM  $\mathcal{Y} \leftarrow \mathcal{Y} \times_n \mathbf{U}^{(n)\top}$ 
  end for
   $\mathcal{G} \leftarrow \mathcal{Y}$ 
  return ( $\mathcal{G}$ ,  $\{\mathbf{U}^{(n)}\}$ )
end procedure
  
```

Vannieuwenhoven, Vandebril, and Meerbergen (2012)

# Tensors in scientific applications are huge, need parallel methods



# Tensor distribution: Cartesian



Processor Grid:  $P_1 \times P_2 \times P_3 = 4 \times 3 \times 2$

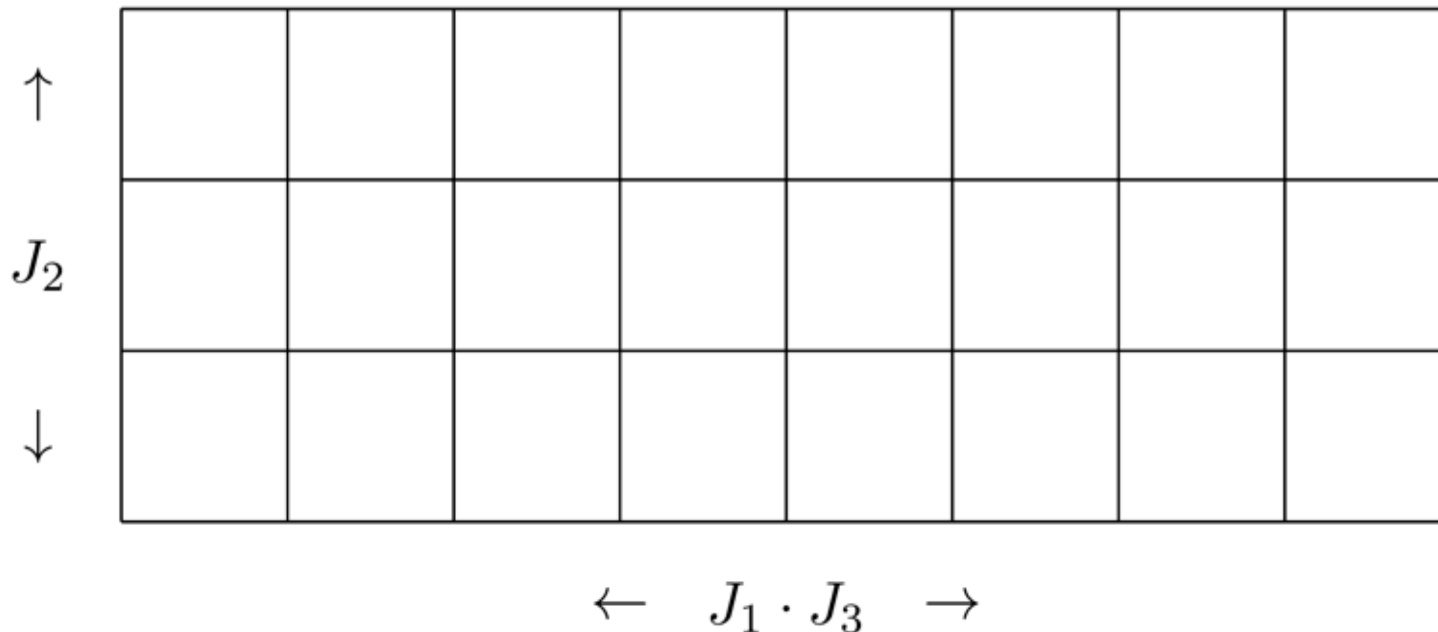
# Unfolded tensor distribution

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

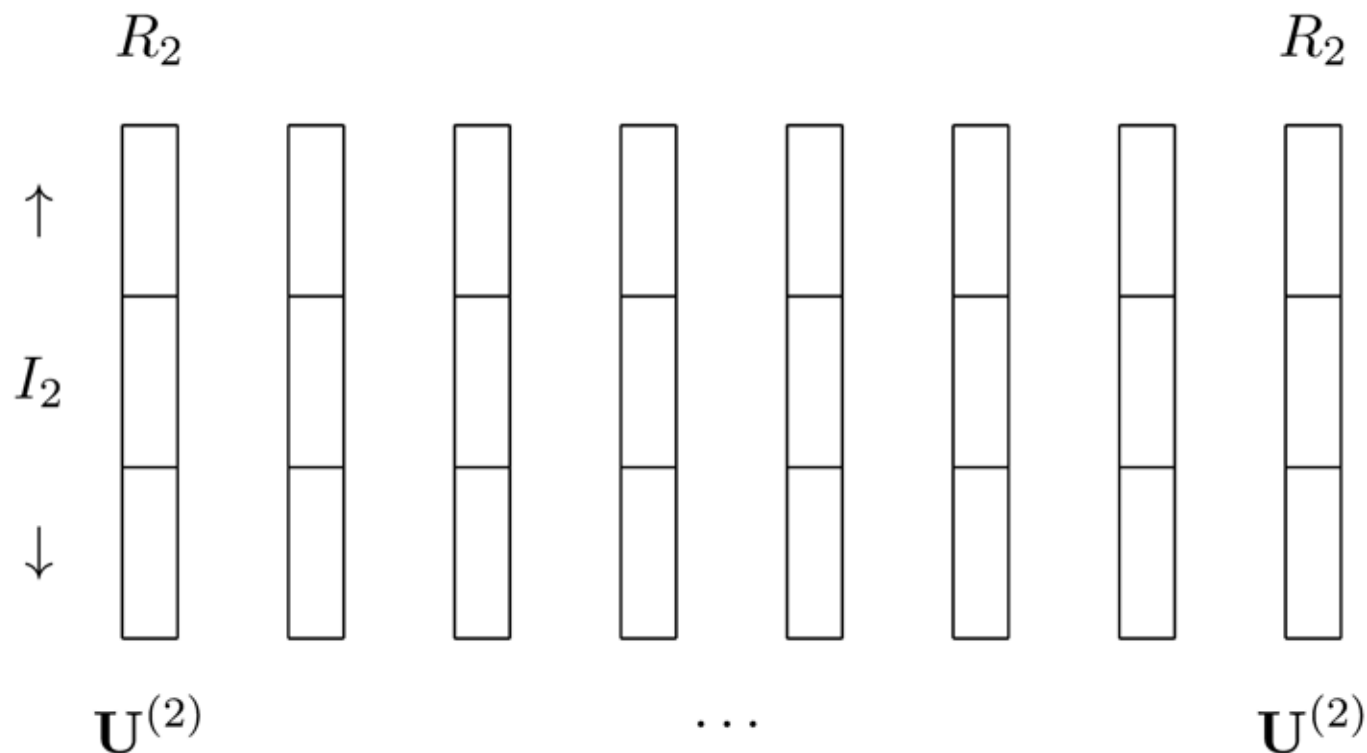
Global Unfolded Tensor:  $J_n \times \hat{J}_n$

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



# Redundant factor matrix distribution

Factor matrices are replicated on each processor fiber  
and 1D row-distributed on each fiber



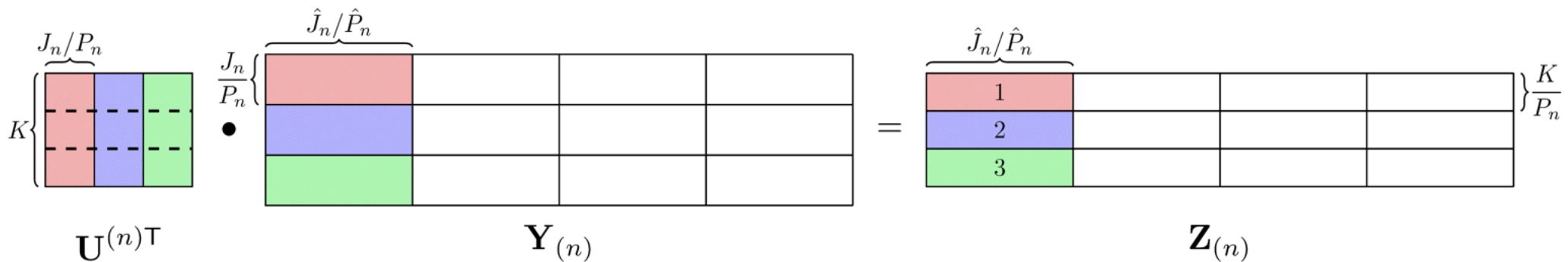
Processor Grid:  $4 \times 3 \times 2$

# Parallel TTM

$$\mathcal{Z} = \mathcal{Y} \times_n \mathbf{U}^{(n)\top} \Leftrightarrow \mathbf{Z}_{(n)} = \mathbf{U}^{(n)\top} \mathbf{Y}_{(n)}$$

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

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



**procedure** TTM( $\mathcal{Y}, \mathbf{V}, n$ )

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

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

**for**  $\ell = 1, \dots, P_n$  **do**

$\mathcal{W} \leftarrow \bar{\mathcal{Y}} \times_n \bar{\mathbf{V}}^{[\ell]}$

$\bar{\mathcal{Z}} \leftarrow \text{REDUCE}(\mathcal{W}, \text{myProcCol}, \ell)$

**end for**

**return**  $\bar{\mathcal{Z}}$

**end procedure**

$$C_{\text{TTM}} = 2\gamma \frac{JK}{P} + \alpha P_n \log P_n + \beta (P_n - 1) \frac{\hat{J}_n K}{P}$$

$$M_{\text{TTM}} = \underbrace{J/P}_{\bar{\mathcal{Y}}} + \underbrace{J_n K/P_n}_{\bar{\mathbf{V}}} + \underbrace{\hat{J}_n K/P}_{\bar{\mathcal{Z}}} + \underbrace{\hat{J}_n K/P}_{\mathcal{W}}$$

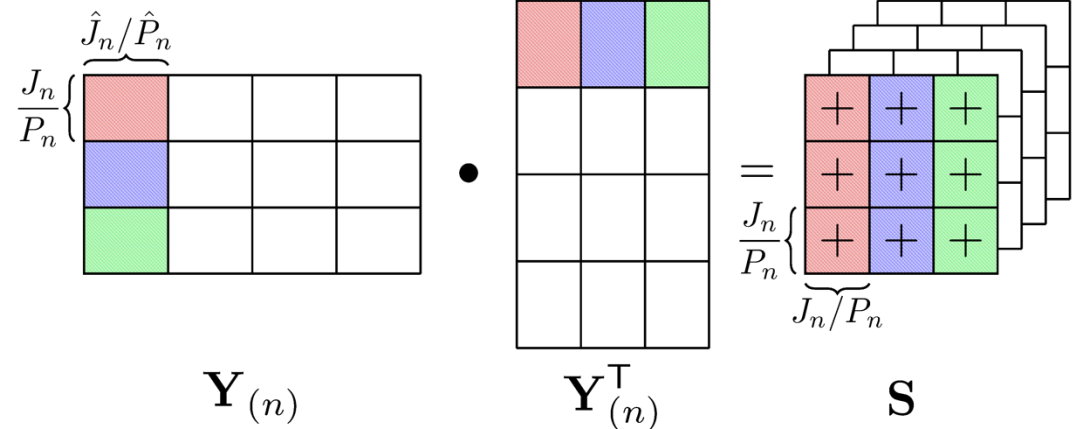


# Parallel Gram

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

```

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  $\leftarrow (*, \dots, *, p_k, *, \dots, *)$ 
   $\mathbf{V}^{[p_n]} \leftarrow \bar{\mathbf{Y}}_{(n)} \bar{\mathbf{Y}}_{(n)}^\top$ 
  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  $(p_1, \dots, p_{n-1}, k, \dots, p_N)$ 
     $\mathbf{V}^{[k]} \leftarrow \bar{\mathbf{Y}}_{(n)} \mathbf{W}_{(n)}^\top$ 
  end for
   $\bar{\mathbf{S}} = \text{All-Reduce}(\mathbf{V}, \text{myProcRow})$ 
  return  $\bar{\mathbf{S}}$ 
end procedure
  
```



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

$$M_{\text{GRAM}} = \underbrace{J/P}_{\bar{\mathbf{y}}} + \underbrace{J/P}_{\mathbf{W}} + \underbrace{J_n^2/P_n}_{\mathbf{V}} + \underbrace{J_n^2/P_n}_{\bar{\mathbf{S}}}$$

# Parallel eigenvector computation

```
1: procedure EIGENVECTORS( $\bar{\mathbf{S}}, R_n, n$ )
2:   myProcID = ( $p_1, p_2, \dots, p_N$ )
3:   myProcCol  $\leftarrow$  ( $p_1, \dots, p_{n-1}, *, p_{n+1}, \dots, p_N$ )
4:    $\mathbf{S} = \text{ALL-GATHER}(\bar{\mathbf{S}}, \text{myProcCol})$ 
5:    $\mathbf{U}^{(n)} = \text{LOCAL-EIGENVECTORS}(\mathbf{S}, R_n)$ 
6:    $\bar{\mathbf{U}}^{(n)} = \text{ROW-SUBSET}(\mathbf{U}^{(n)}, P_n, p_n)$  ▷ Extract  $p_n$ -th block
7:   return  $\bar{\mathbf{U}}^{(n)}$ 
8: end procedure
```

Every processor redundantly computes the leading eigenvectors of the Gram matrix

# Application results: compression versus accuracy

Ranks depend on error:

$$\|\mathbf{x} - \mathcal{M}\|^2 \leq \sum_{n=1}^N \left( \sum_{i=R_n+1}^{I_n} \sigma_i(\mathbf{X}_{(n)})^2 \right)$$

Compression ratio:

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

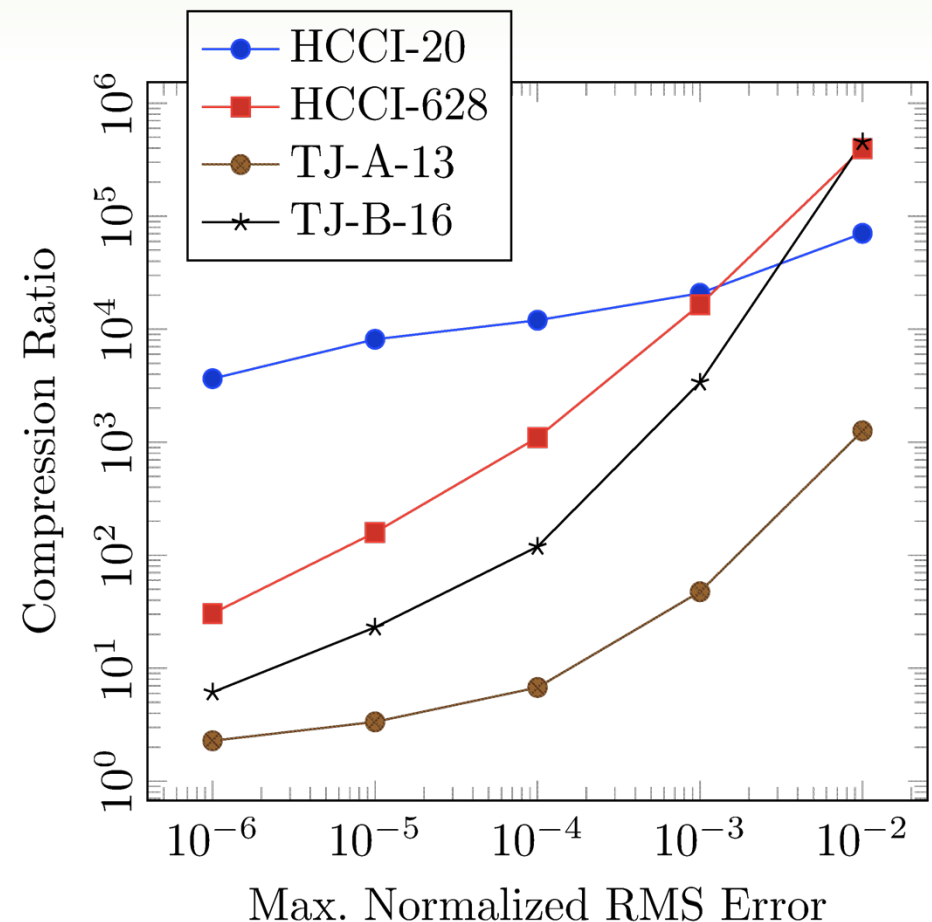
*Simulation of an autoignitive premixture of air and ethanol in Homogeneous Charge Compression Ignition*

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

# Sample results for one species in HCCI: error is negligible

Original  $\mathcal{X}$

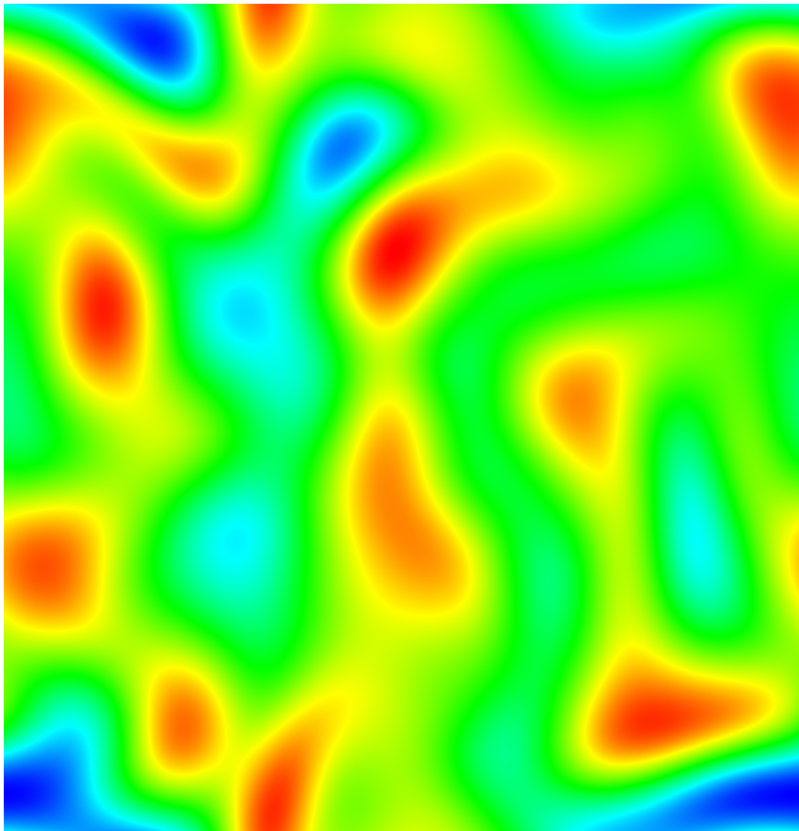
672 x 672 x 33 x 8

Compression: 586

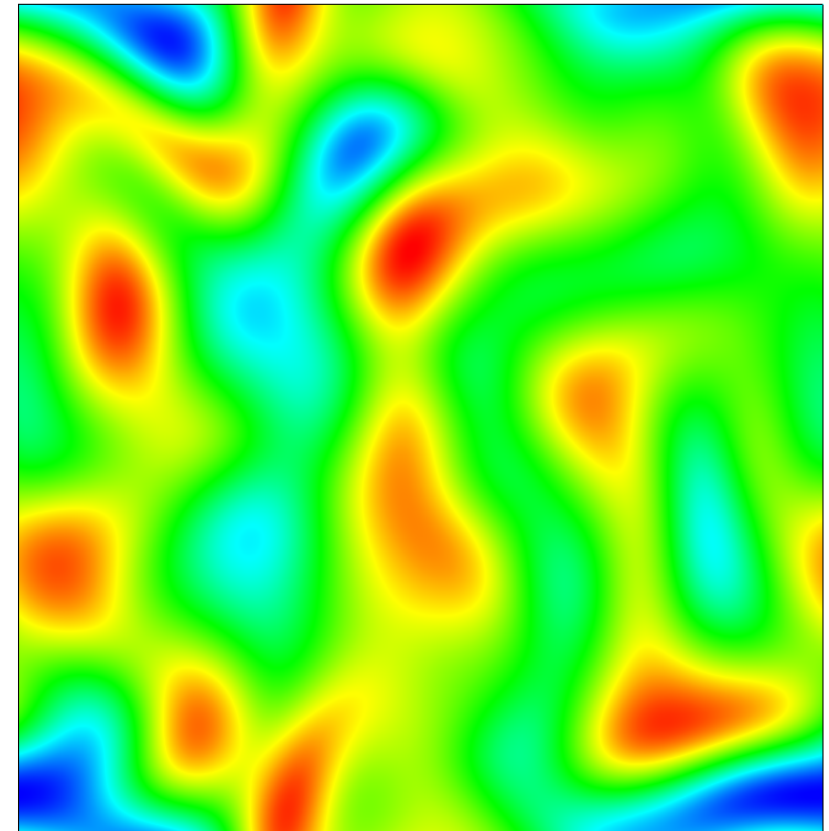
Recovered  $\hat{\mathcal{X}}$

48 x 48 x 20 x 3

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



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



910MB compressed to 1.5MB

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

# Sample results for “derived” quantity: error is negligible

Original  $\mathcal{X}$

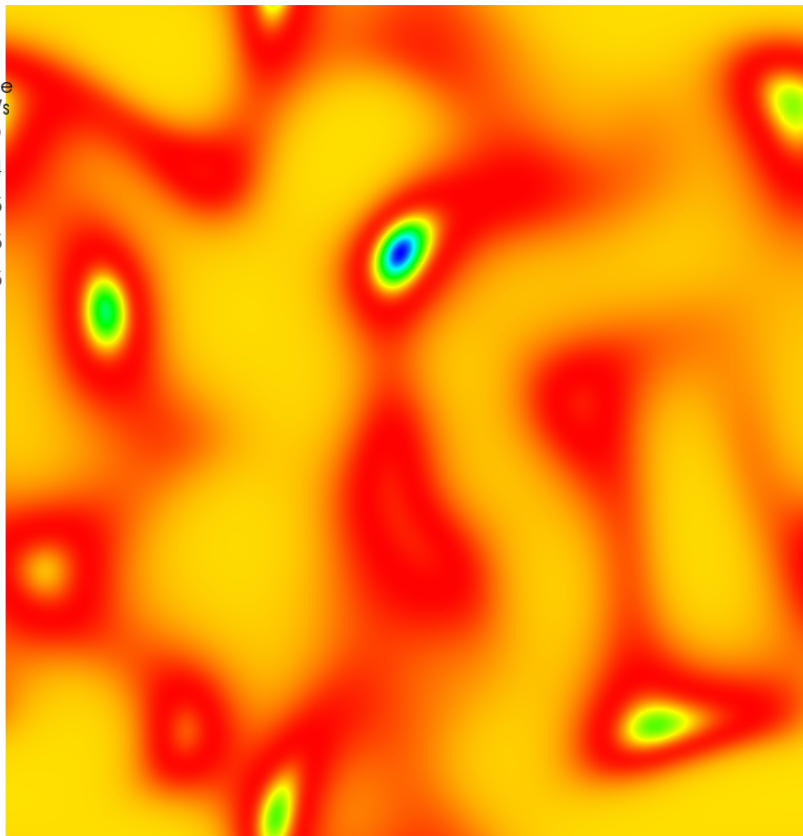
672 x 672 x 33 x 8

Compression: 586

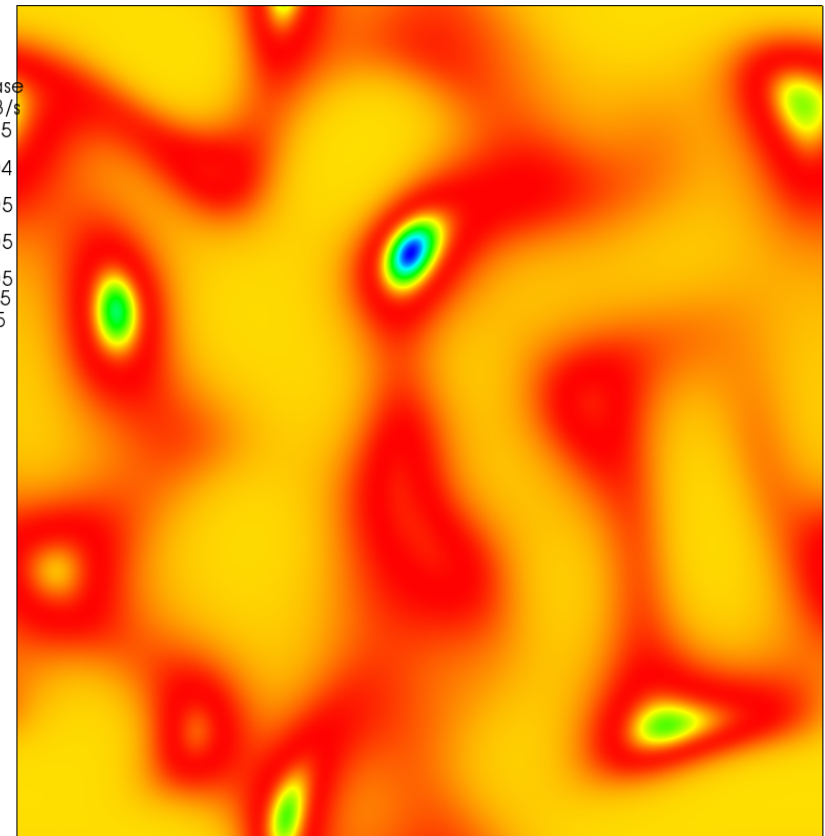
Recovered  $\hat{\mathcal{X}}$

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 compressed to 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

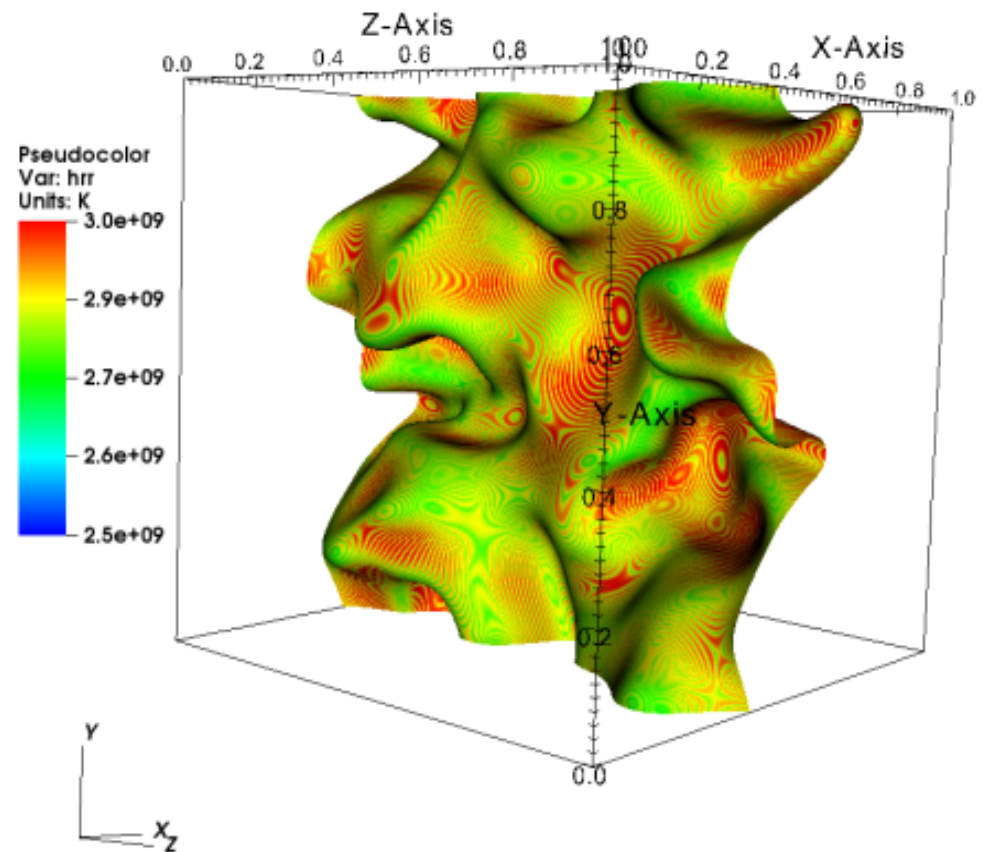
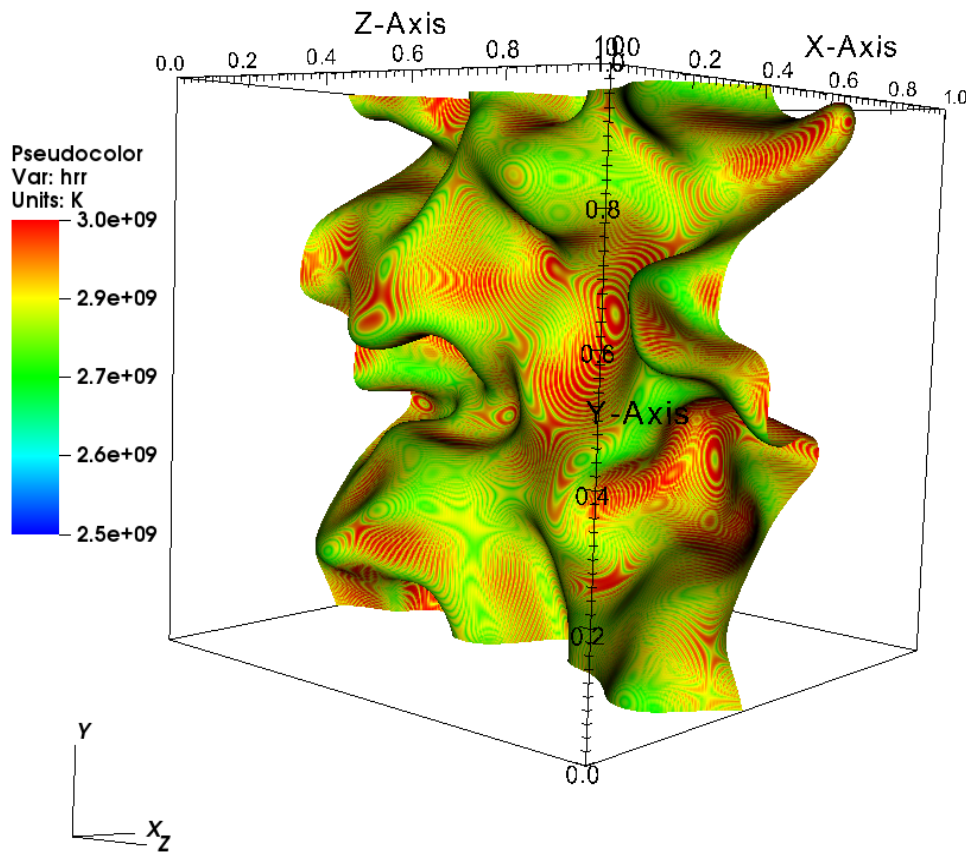
Original  $\mathcal{X}$

Compression: 1000

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

500 x 500 x 500 x 11 x 5

50 x 50 x 50 x 11 x 5



52 GB compressed to 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!

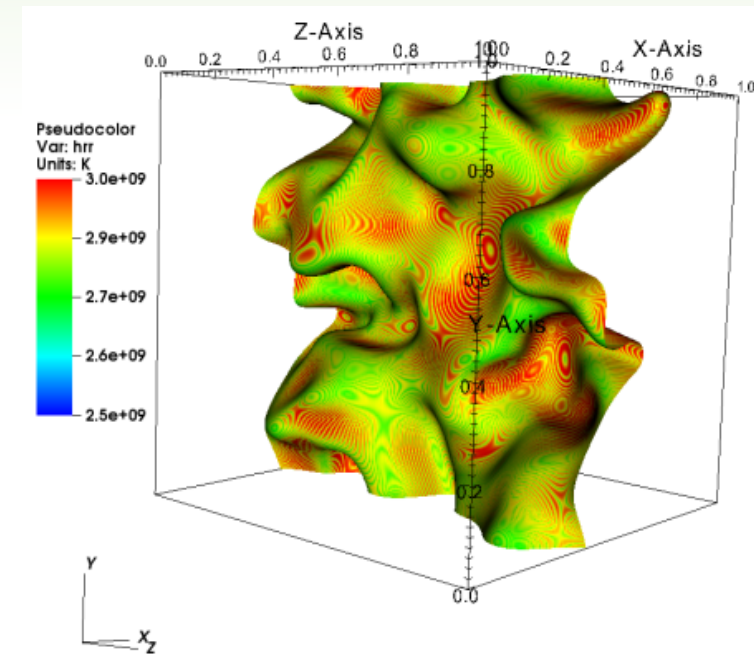
$$\hat{\mathcal{X}} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \times_4 \mathbf{U}^{(4)} \times_5 \mathbf{U}^{(5)}$$

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

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

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

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

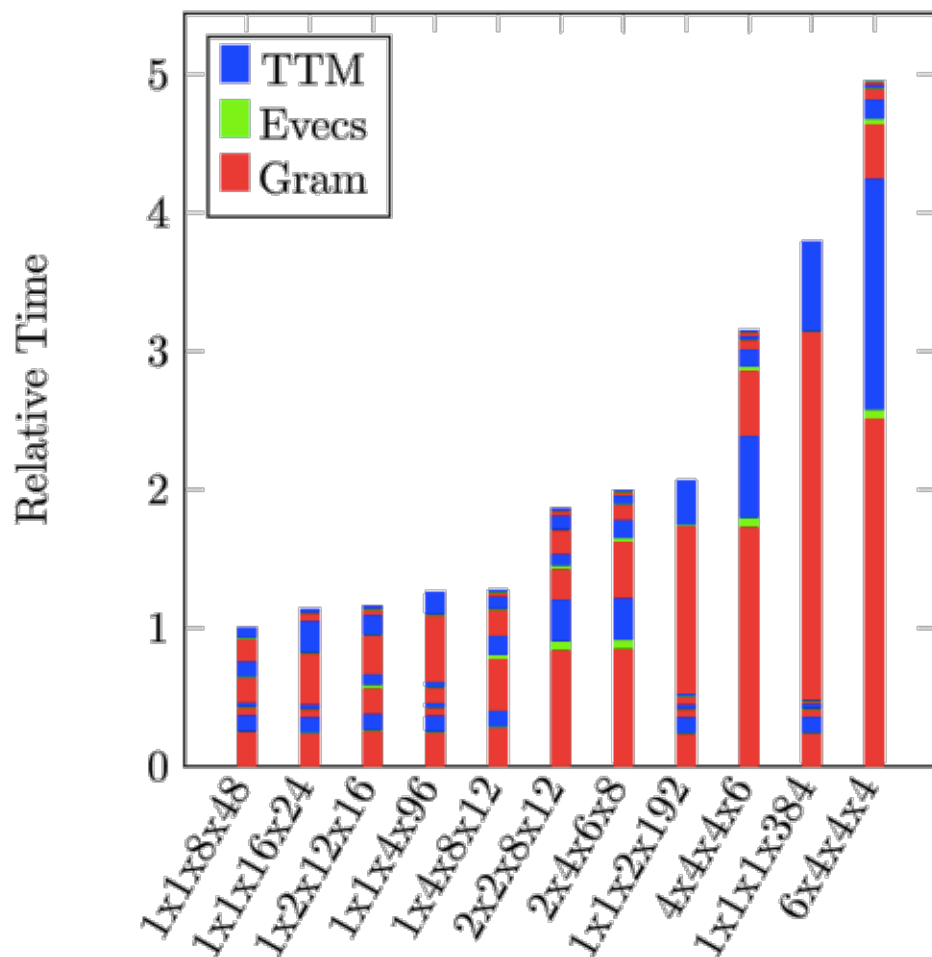


# Parameter choices: processor grid configuration & mode order

## Processor Grid Configuration

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

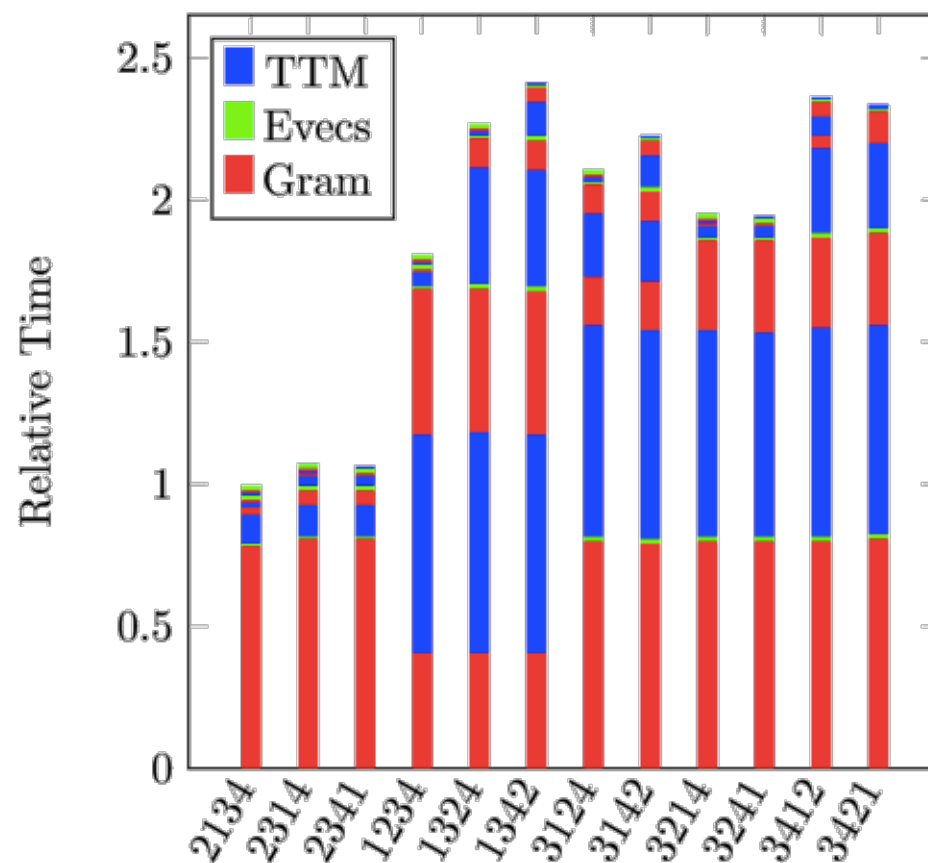
$R: 96 \times 96 \times 96 \times 96$



## Mode Order

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

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



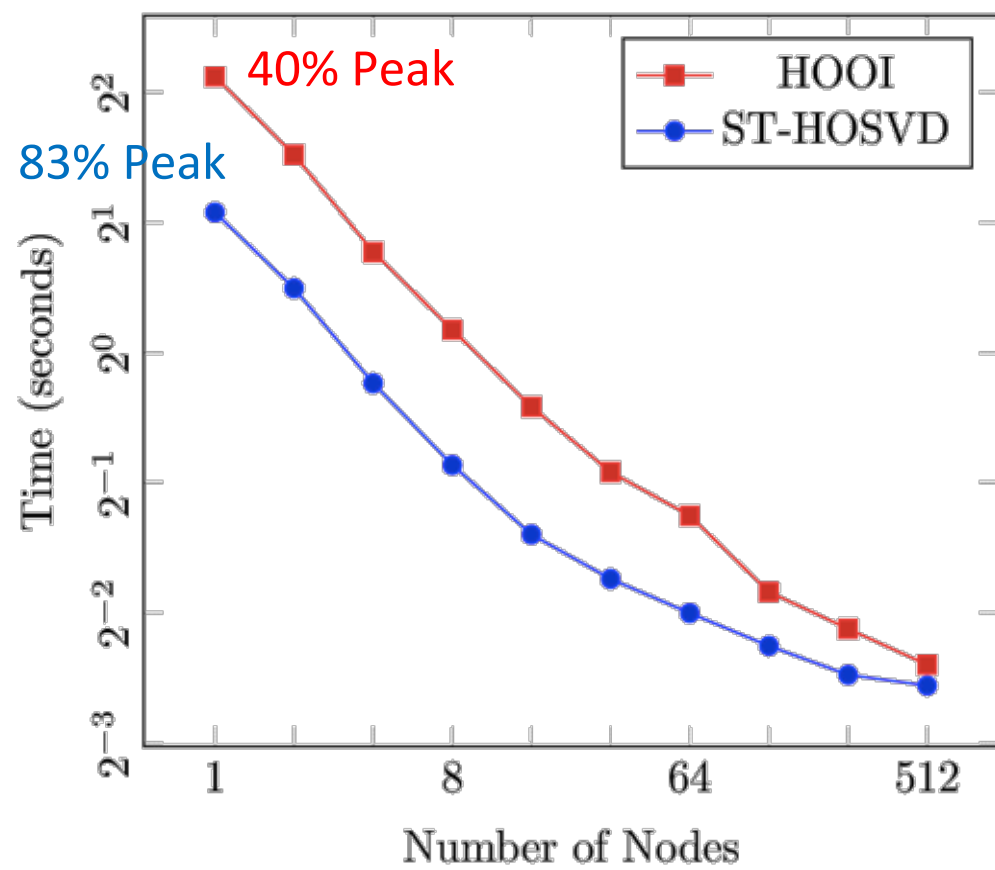


# Strong & weak scaling on Edison

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

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

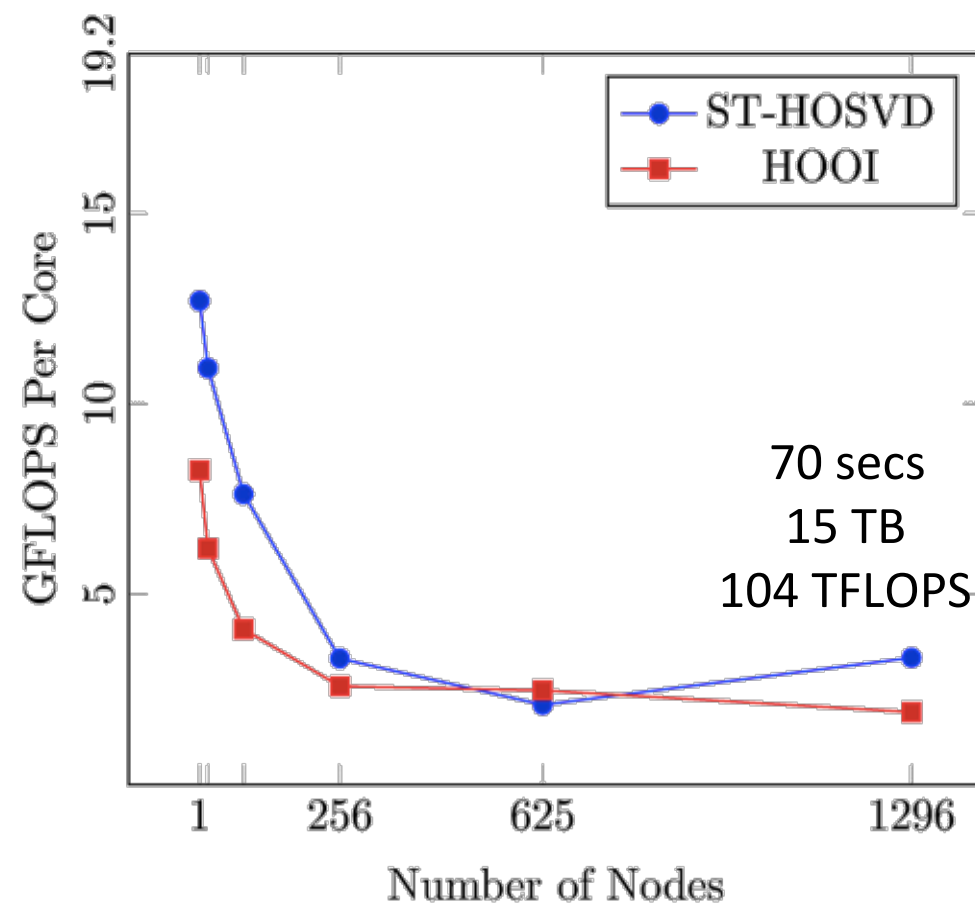
$R$ :  $42 \times 115 \times 81 \times 19$



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

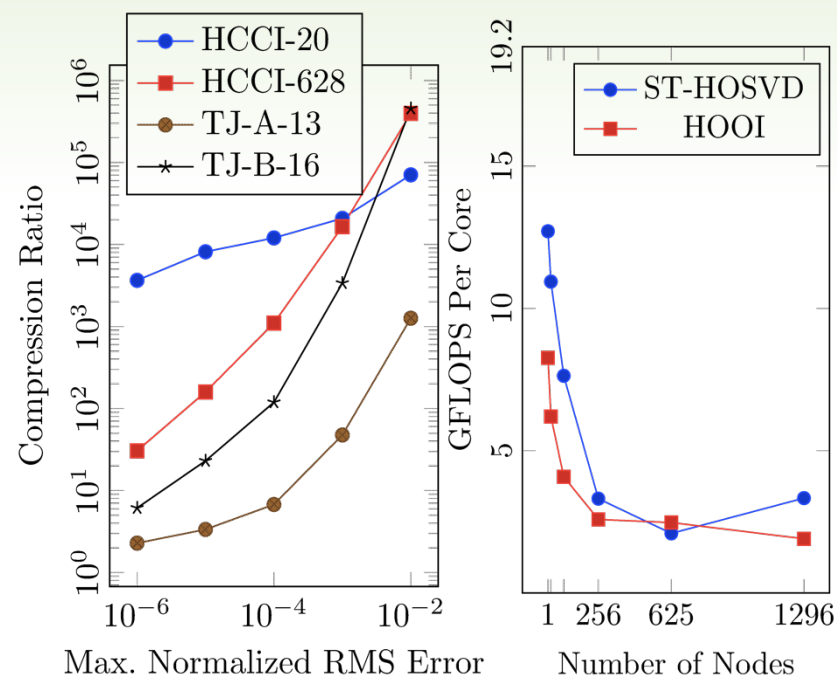
$I$ :  $200k \times 200k \times 200k \times 200k$

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



# Parallel Tucker Compression

- First-ever implementation of distributed-memory parallel Tucker decomposition
- Up to  $10^6$  compression on real-world data with minimal loss in accuracy
- Scales well – achieving 17% of peak on over 30,000 cores
- Future work
  - Detailed application studies
  - Use QR/SVD instead of Gram/EVD



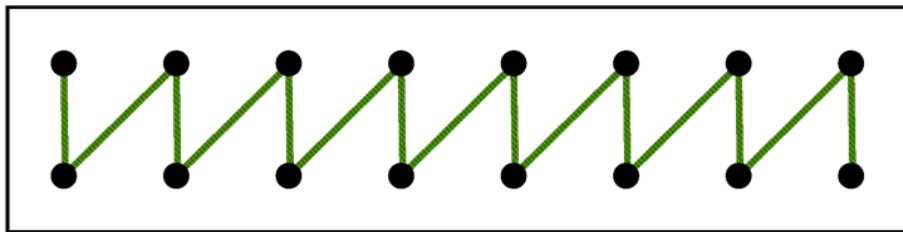
For more information:  
Grey Ballard,  
gmballa@sandia.gov

W. Austin, G. Ballard, and T. G. Kolda,  
***Parallel Tensor Compression for Large-Scale Scientific Data,***  
<http://arxiv.org/abs/1510.06689>, October 2015

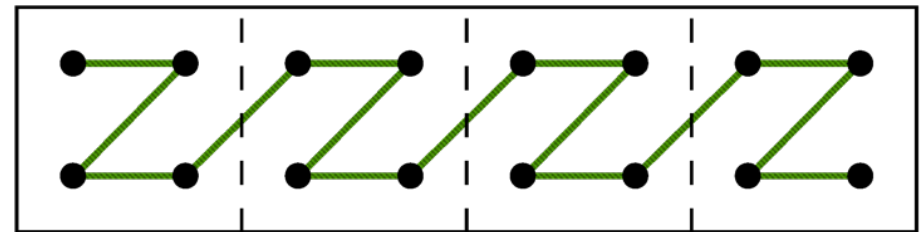
# Backup Slides

# Local unfolded tensor layout

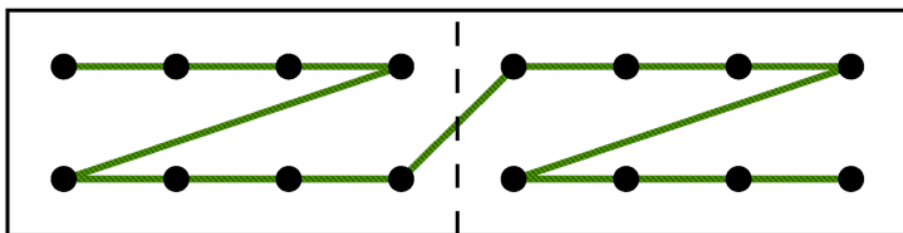
Local Layout:  $2 \times 2 \times 2 \times 2$



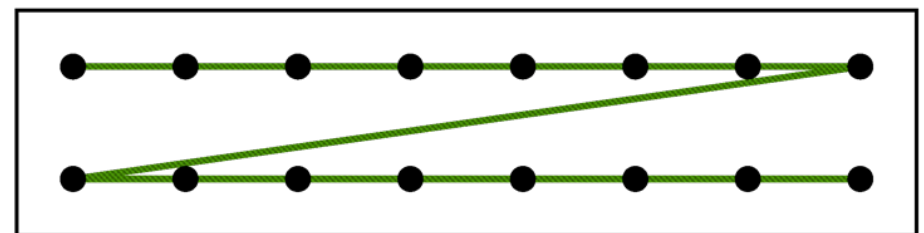
$n = 1$



$n = 2$



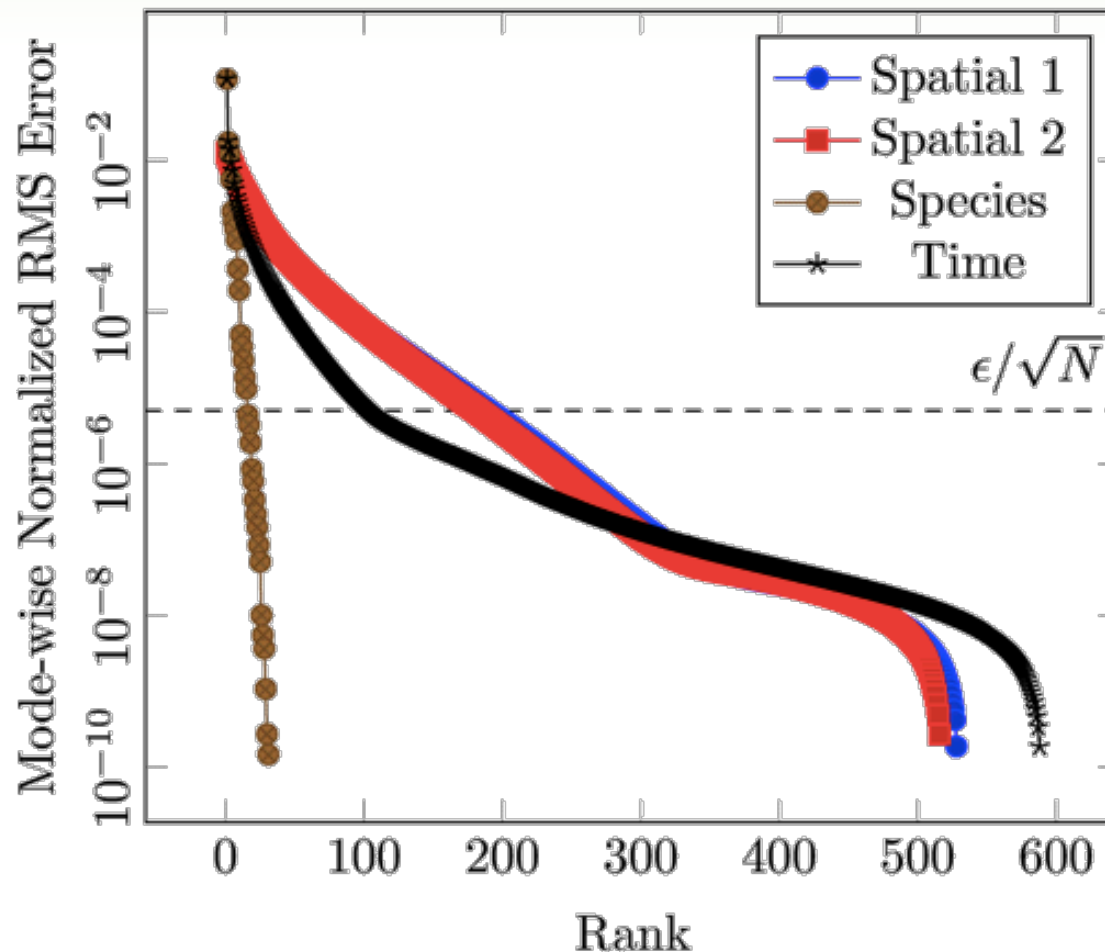
$n = 3$



$n = 4$

# Mode-wise contributions to approximation error bound

2D HCCI Data (628 time steps)



$$672 \times 672 \times 33 \times 628 \rightarrow 192 \times 183 \times 16 \times 104$$

# Elementwise errors

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

Table 1: Compression and maximum absolute elementwise error of centered data for normalized RMS error of  $1e-5$ .