

# An Input-Adaptive and In-Place Approach to Dense Tensor-Times-Matrix Multiply

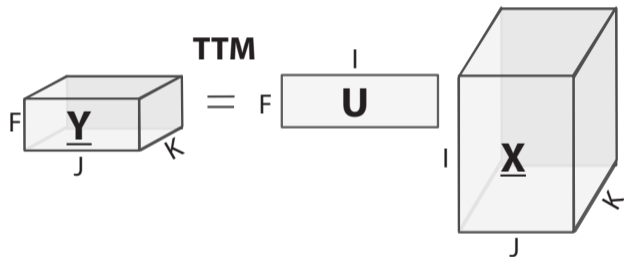
Jiajia Li, Casey Battaglini, Ioakeim Perros,  
Jimeng Sun, Richard Vuduc

Computational Science & Engineering,  
Georgia Institute of Technology

14<sup>th</sup> April 2016

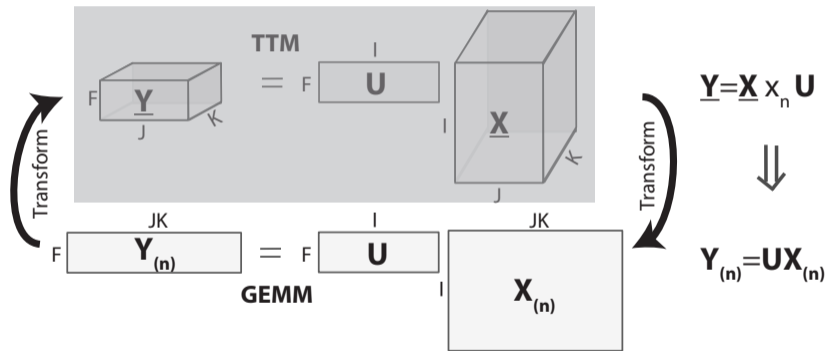


# The problem

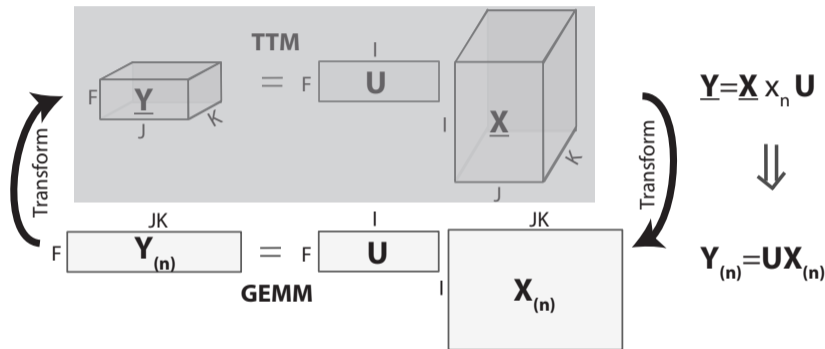


$$\underline{\mathbf{Y}} = \underline{\mathbf{X}} \times_n \mathbf{U}$$

# The problem



# The problem



Transform:

70% running time.  
50% space.

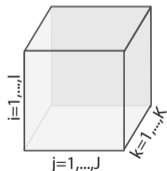
- We proposed an in-place TTM algorithm and employed auto-tuning method to adapt its parameters.

# Outline

- Background
- Motivation
- InTensLi Framework
- Experiments and Analysis
- Conclusion

# Tensor and Applications

- Tensor: interpreted as a multi-dimensional array, e.g.  $\underline{\mathbf{X}} \in \mathbb{R}^{I \times J \times K}$ .
  - Special cases: vectors ( $\mathbf{x}$ ) are 1D tensors, and matrices ( $\mathbf{A}$ ) are 2D tensors.
  - Tensor dimension ( $N$ ): also called mode or order.
  - We focus on dense tensors in this work.
- Applications
  - Quantum chemistry, quantum physics, signal and image processing, neuroscience, and data analytics.



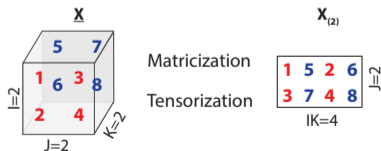
A third-order (or three-dimensional)  $I \times J \times K$  tensor.

# Tensor Representations

- Sub-tensor



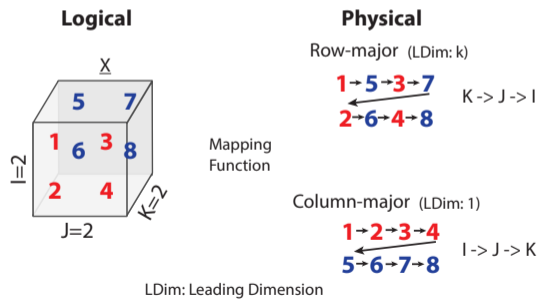
- Whole tensor



- Diff representations  $\rightarrow$  Diff algorithms  $\rightarrow$  Diff performance.

# Memory Mapping

- Tensor organization
  - Multi-dimensional array – logically
  - Linear storage – physically
- Memory mapping<sup>1</sup>.

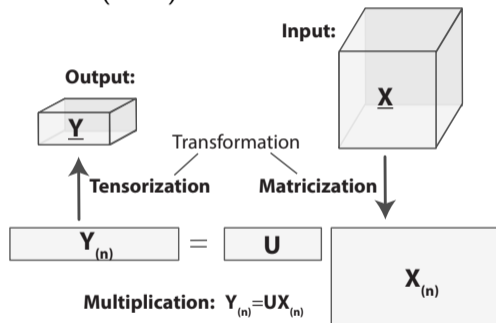


<sup>1</sup>GARCIA, R., and LUMSDAINE, A. Multirray: A c++ library for generic programming with arrays. Software Practive Experience 35 (2004), 159–188.



# T<sub>TM</sub> Algorithm

- Baseline tensor-times-matrix multiply (T<sub>TM</sub>) algorithm in TENSOR TOOLBOX and CYCLOPS Tensor Framework (CTF).



- T<sub>TM</sub> Applications

- Low-rank tensor decomposition.
- Tucker decomposition, e.g. TUCKER-HOOI algorithm.

$$\underline{\mathbf{Y}} = \underline{\mathbf{X}} \times_1 \mathbf{A}^{(1)T} \cdots \times_{n-1} \mathbf{A}^{(n-1)T} \times_{n+1} \mathbf{A}^{(n+1)T} \cdots \times_N \mathbf{A}^{(N)T}.$$

# Main Contributions

- Proposed an in-place tensor-times-matrix multiply ( $\text{INTTM}$ ) algorithm, by avoiding physical reorganization of tensors.
- Built an input-adaptive framework  $\text{INTENSLI}$  to automatically adapt parameters and generate the code.
- Achieved  $4\times$  and  $13\times$  speedups compared to the state-of-the-art  $\text{TENSOR TOOLBOX}$  and  $\text{CTF}$  tools.

## Observation 1: Transformation is expensive.

Notation: the number of words ( $Q$ ), floating-point operations ( $W$ ), last-level cache size ( $Z$ ).

The relation of them is  $Q \geq \frac{W}{8\sqrt{Z}} - Z^2$  for both general matrix-matrix multiply (GEMM) and  $T_{TM}$ .

- Suppose  $T_{TM}$  does the same number of flops as GEMM ( $\hat{W} = W$ ), the relation of Arithmetic Intensity of GEMM and  $T_{TM}$ :  $\hat{A} \approx A / (1 + \frac{A}{m})$  when counting transformation.  
 $(1 + \frac{A}{m})$  is the penalty.
- Assume cache size  $Z$  is 8MB, the penalty of a 3-D tensor is 33.

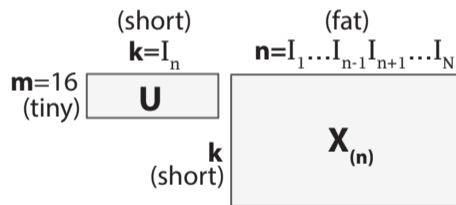
**Conclusion:** When  $T_{TM}$  and GEMM do the same number of flops, Arithmetic Intensity of  $T_{TM}$  is decreased by a penalty of 33 or more, as tensor dimension increases.

---

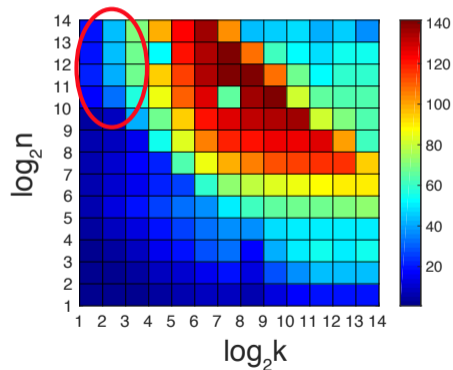
<sup>2</sup>G. Ballard, E. Carson, J. Demmel, M. Hoemmen, N. Knight, and O. Schwartz. Communication lower bounds and optimal algorithms for numerical linear algebra. *Acta Numerica*, 23:pp. 1–155, 2014.

## Observation 2: Performance of the multiplication in $T_{TM}$ is far below peak.

- $T_{TM}$  algorithm involves a variety of rectangular problem sizes.



(a) TTM's multiplication.



(b) GEMM performance in Intel MKL with 4 threads.

## Observation 3: T<sub>TM</sub> organization is critical to data locality.

- There are many ways to organize data accesses.
- Choose slice representation.

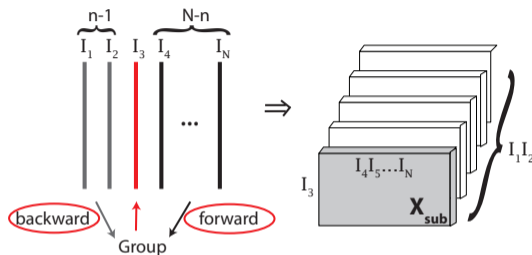
Table 1 : Different representation forms of mode-1 T<sub>TM</sub> on a  $I \times J \times K$  tensor.

Mode-1 Product Representation Forms		BLAS Level	Transformation
Full reorganization	<i>Tensor representation</i> $\underline{\mathbf{Y}} = \underline{\mathbf{X}} \times_1 \mathbf{U}$	—	—
	<i>Matrix representation</i> $\mathbf{Y}_{(1)} = \mathbf{U}\mathbf{X}_{(1)}$	L3	Yes
Sub-tensor extraction	<i>Fiber representation</i> $\mathbf{y}(f, :, k) = \mathbf{X}(:, :, k)\mathbf{u}(f, :)$ , <i>Loops</i> : $k = 1, \dots, K, f = 1, \dots, F$	L2	No
	<i>Slice representation</i> $\mathbf{Y}(:, :, k) = \mathbf{U}\mathbf{X}(:, :, k)$ , <i>Loops</i> : $k = 1, \dots, K$	L3	No

# Layout

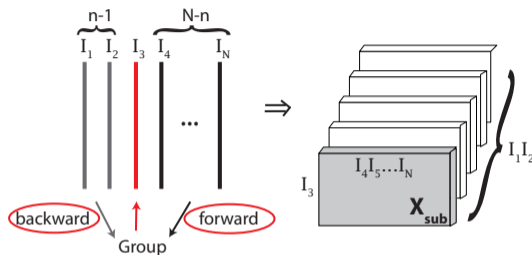
- 1 Background
- 2 Motivation
- 3 InTensLi Framework**
  - Algorithmic Strategy
  - InTensLi Framework
- 4 Experiments and Analysis
- 5 Conclusion
- 6 References

# Algorithmic Strategy



- To avoid data copy,
  - Rules: 1) compress only contiguous dimensions; 2) always include the leading dimension.
  - Lemma: TTM can be performed on up to  $\max\{n-1, N-n\}$  contiguous dimensions without physical reorganization.

# Algorithmic Strategy



- To avoid data copy,
  - Rules: 1) compress only contiguous dimensions; 2) always include the leading dimension.
  - Lemma: TTM can be performed on up to  $\max\{n - 1, N - n\}$  contiguous dimensions without physical reorganization.
- To get high performance of GEMM,
  - Find an approximate matrix size according to computer architecture.
  - Use auto-tuning method in INTENSLI framework.



# INT<sub>TM</sub> Algorithm and Comparison

- INT<sub>TM</sub>'s AI:  $\tilde{A} \lesssim \frac{\hat{Q}}{8\sqrt{Z}} = 8\sqrt{Z} \approx A.$
- Traditional T<sub>TM</sub>'s AI:  $\hat{A} \approx \frac{A}{1 + \frac{A}{m}}.$
- INT<sub>TM</sub> eliminates the AI by a factor  $1 + \frac{A}{m}.$

**Input:** A dense tensor  $\underline{\mathbf{X}} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , a dense matrix  $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ , and an integer  $n$ ;  
**Output:** A dense tensor  $\underline{\mathbf{Y}} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N};$

```

// Nested loops, using  $P_L$  threads
1: parfor  $i_l = 1$  to  $I_l$ , all  $i_l \in M_L$  do
2:   if  $M_C$  are on the left of  $i_n$  then
3:      $\mathbf{X}_{\text{sub}} = \text{inplace-mat}(\underline{\mathbf{X}}, M_C, i_n);$ 
4:      $\mathbf{Y}_{\text{sub}} = \text{inplace-mat}(\underline{\mathbf{Y}}, M_C, j);$ 

// Matrix-matrix multiplication, using  $P_C$  threads
5:      $\mathbf{Y}_{\text{sub}} = \mathbf{X}_{\text{sub}} \mathbf{U}'$ ,  $\mathbf{U}'$  is the transpose of  $\mathbf{U}$ .
6:   else
7:      $\mathbf{X}_{\text{sub}} = \text{inplace-mat}(\underline{\mathbf{X}}, i_n, M_C);$ 
8:      $\mathbf{Y}_{\text{sub}} = \text{inplace-mat}(\underline{\mathbf{Y}}, j, M_C);$ 

// Matrix-matrix multiplication, using  $P_C$  threads
9:      $\mathbf{Y}_{\text{sub}} = \mathbf{U} \mathbf{X}_{\text{sub}}$ 
10:   end if
11: end parfor
12: return  $\underline{\mathbf{Y}}$ ;

```

In-place Tensor-Times-Matrix Multiply (INT<sub>TM</sub>) algorithm.

# INT<sub>TM</sub> Algorithm and Comparison

- INT<sub>TM</sub>'s AI:  $\tilde{A} \lesssim \frac{\hat{Q}}{8\sqrt{Z}} = 8\sqrt{Z} \approx A$ .
- Traditional T<sub>TM</sub>'s AI:  $\hat{A} \approx \frac{A}{1 + \frac{A}{m}}$ .
- INT<sub>TM</sub> eliminates the AI by a factor  $1 + \frac{A}{m}$ .

**Input:** A dense tensor  $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , a dense matrix  $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ , and an integer  $n$ ;  
**Output:** A dense tensor  $\mathbf{Y} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N}$ ;

```

// Nested loops, using  $P_L$  threads
1: parfor  $i_l = 1$  to  $I_l$ , all  $i_l \in M_L$  do
2:   if  $M_C$  are on the left of  $i_n$  then
3:      $\mathbf{X}_{\text{sub}} = \text{inplace-mat}(\mathbf{X}, M_C, i_n)$ ;
4:      $\mathbf{Y}_{\text{sub}} = \text{inplace-mat}(\mathbf{Y}, M_C, j)$ ;

// Matrix-matrix multiplication, using  $P_C$  threads
5:      $\mathbf{Y}_{\text{sub}} = \mathbf{X}_{\text{sub}} \mathbf{U}'$ ,  $\mathbf{U}'$  is the transpose of  $\mathbf{U}$ .
6:   else
7:      $\mathbf{X}_{\text{sub}} = \text{inplace-mat}(\mathbf{X}, i_n, M_C)$ ;
8:      $\mathbf{Y}_{\text{sub}} = \text{inplace-mat}(\mathbf{Y}, j, M_C)$ ;

// Matrix-matrix multiplication, using  $P_C$  threads
9:      $\mathbf{Y}_{\text{sub}} = \mathbf{U} \mathbf{X}_{\text{sub}}$ 
10:   end if
11: end parfor
12: return  $\mathbf{Y}$ ;

```

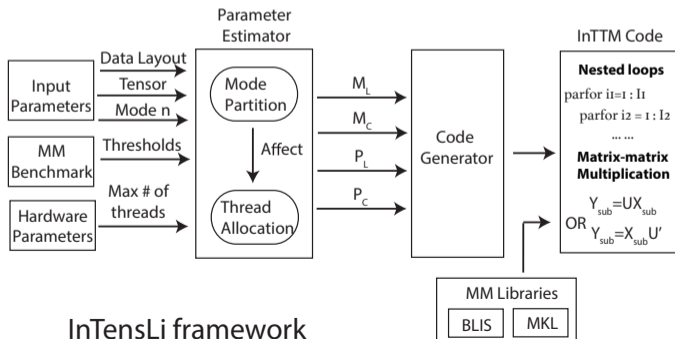
In-place Tensor-Times-Matrix Multiply (INT<sub>TM</sub>) algorithm.

# Layout

- 1 Background
- 2 Motivation
- 3 InTensLi Framework**
  - Algorithmic Strategy
  - InTensLi Framework**
- 4 Experiments and Analysis
- 5 Conclusion
- 6 References

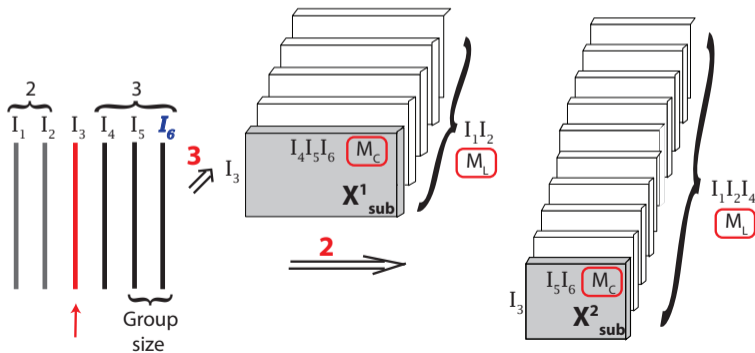
# INTENSLI Framework

- Input: tensor features, hardware configuration, and MM benchmark.
- Parameter estimation
  - Mode partitioning:  $M_L$  and  $M_C$ .
  - Thread allocation:  $P_L$  and  $P_C$ .
- Code generation

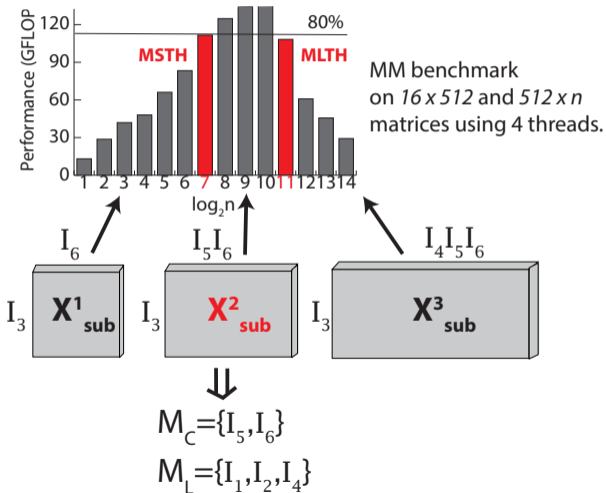


# Parameter Estimation – Mode Partitioning

- Determine grouping direction
  - Row-major  $\leftrightarrow$  forward
  - Column-major  $\leftrightarrow$  backward
- Group size decides INTM algorithm.

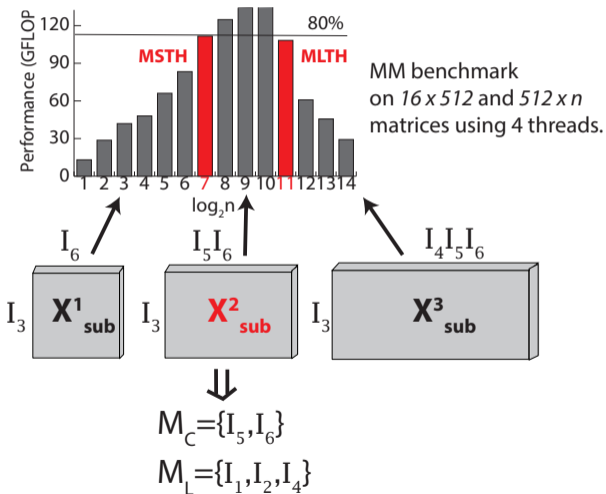


# Choosing Group Size



- *MSTH* and *MLTH*: Thresholds of GEMM's size, the size of all the three operating matrices.
- *MSTH* = 1.04MB and *MLTH* = 7.04MB in our experiments.

# Choosing Group Size



- **MSTH** and **MLTH**: Thresholds of GEMM's size, the size of all the three operating matrices.
- **MSTH** = 1.04MB and **MLTH** = 7.04MB in our experiments.
- **Decide  $M_C$** : Use **MSTH** and **MLTH** to decide group size, then decide  $M_C$ .
- **Decide  $M_L$** : The rest modes of  $M_C$ , except mode- $n$ .

# Thread Allocation and Code Generation

- Thread allocation
  - In most cases, maximum performance is obtained by only two configurations:
    - Small matrices: all threads are allocated to nested loops.
    - Large matrices: all threads are allocated to GEMM operation.
  - A threshold  $PTH$  is set to distinguish the GEMM size, which is 800 KB in our tests.
- Code generation
  - Generate nested loops and wrappers for the GEMM kernel.
  - Code generated in C++, using OpenMP with the collapse directive.



# Experimental Platforms

- Double-precision is adopted in our experiments.
- We employ 8 and 32 threads on the two platforms respectively, considering hyper-threading.
- Xeon E7-4820 has a relatively large memory (512 GiB), allowing us to test a larger range of (dense) tensor sizes than has been common in prior single-node studies.

Table 2 : Experimental Platforms Configuration

Parameters	Intel Core i7-4770K	Intel Xeon E7-4820
Microarchitecture	Haswell	Westmere
Frequency	3.5 GHz	2.0 GHz
# of physical cores	4	16
Hyper-threading	On	On
Peak GFLOP/s	224	128
Last-level cache	8 GiB	18 GiB
Memory size	32 GiB	512 GiB
Memory bandwidth	25.6 GB/s	34.2 GB/s
# of memory channels	2	4
Compiler	icc 15.0.2	icc 15.0.0

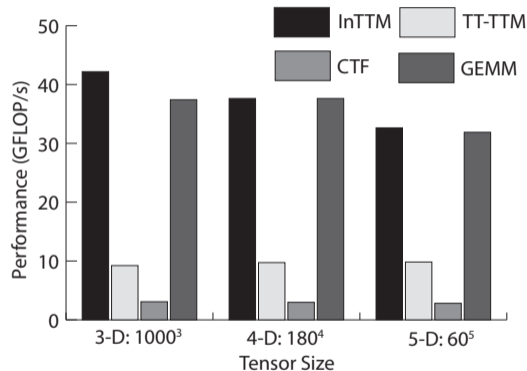
# Performance Comparison

## Implementations

- **INTTM**: INTENSLI generated C++ code with OpenMP.
- **TT-TTM**: TENSOR TOOLBOX library in MATLAB.
- **CTF**: C++ code, supporting MPI+OpenMP parallelization.
- **GEMM**: C++ code, baseline TTM algorithm without transformation.

## Speedup

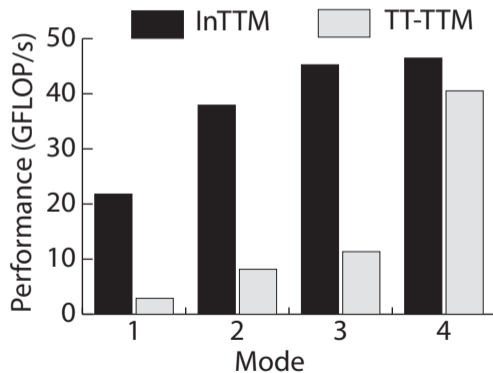
- Obtain  $4\times$  and  $13\times$  speedup compared to TENSOR TOOLBOX and CTF.
- Get close to GEMM-only's performance.



Performance comparison of TTM on mode-2 over diverse dimensional tensors.

# Analysis

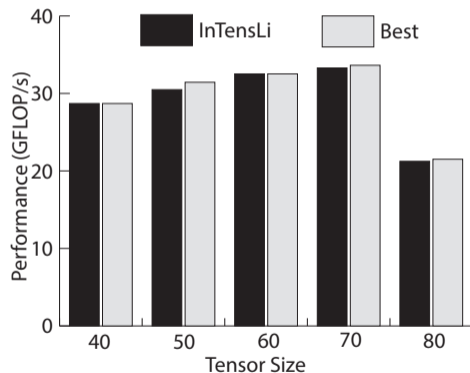
- Performance of different modes.
  - `INTENSLI` is stable for different mode- $n$  products, while `TENSOR TOOLBOX` is not.



Performance behavior of `INTTM` against `TENSOR TOOLBOX` (`TT-TTM`) for different mode products on a  $160 \times 160 \times 160 \times 160$  tensor.

# Analysis

- Parameter selection
  - Compare INTENSLI with exhaustive search, the performance is close to optimal.



Comparison between the performance of  $T_{TM}$  on mode-1 with predicted configuration and the actually highest performance on 5th-order tensors.

# Conclusion

## Summary

- Proposed an in-place tensor-times-matrix multiply (`INTTM`) algorithm, by avoiding physical reorganization of tensors.
- Built an input-adaptive framework `INTENSLI` to automatically do optimization and generate the code.
- Achieved  $4\times$  and  $13\times$  speedups compared to the state-of-the-art `TENSOR TOOLBOX` and `CTF` tools.

## Source code

- <https://github.com/hpcgarage/InTensLi>
- Contact: Jiajia Li (jiajiali@gatech.edu)

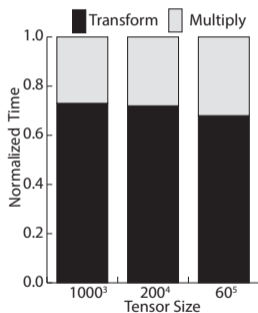
# References

- E. Solomonik, D. Matthews, J. Hammond, and J. Demmel. Cyclops tensor framework: reducing communication and eliminating load imbalance in massively parallel contractions. Technical Report UCB/EECS-2012-210, EECS Department, University of California, Berkeley, Nov 2012.
- B. W. Bader, T. G. Kolda, et al. Matlab tensor toolbox version 2.5. Available from <http://www.sandia.gov/~tgkolda/TensorToolbox/index-2.6.html>, January 2012
- T. Kolda and B. Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- ...

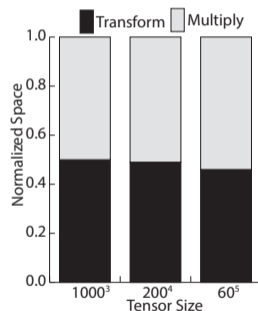
## Backup Slides

## Observation 1: Transformation is expensive.

- Transformation takes about 70% of the total run-time, and close to 50% of the total storage.



(a) Time Profiling



(b) Space Profiling

Profiling of T<sub>TM</sub> algorithm on mode-2 product on 3rd, 4th, and 5th-order tensors, where the output tensors are low-rank representations of corresponding input tensors.



## Observation 3: T<sub>TM</sub> organization is critical to data locality.

- There are many ways to organize data accesses.

