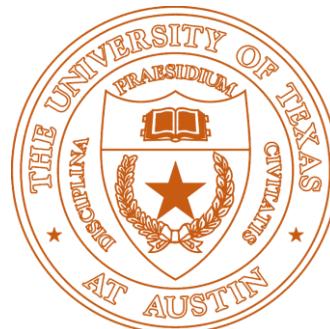


# Strassen's Algorithm for Tensor Contraction



Jianyu Huang, Devin A. Matthews, Robert A. van de Geijn

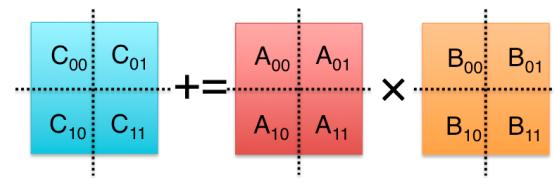
The University of Texas at Austin

September 14-15, 2017

Tensor Computation Workshop

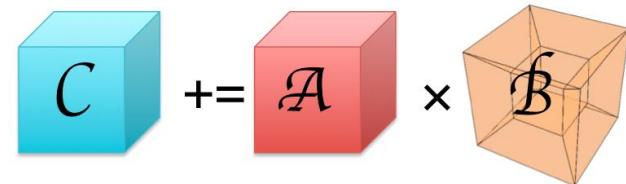
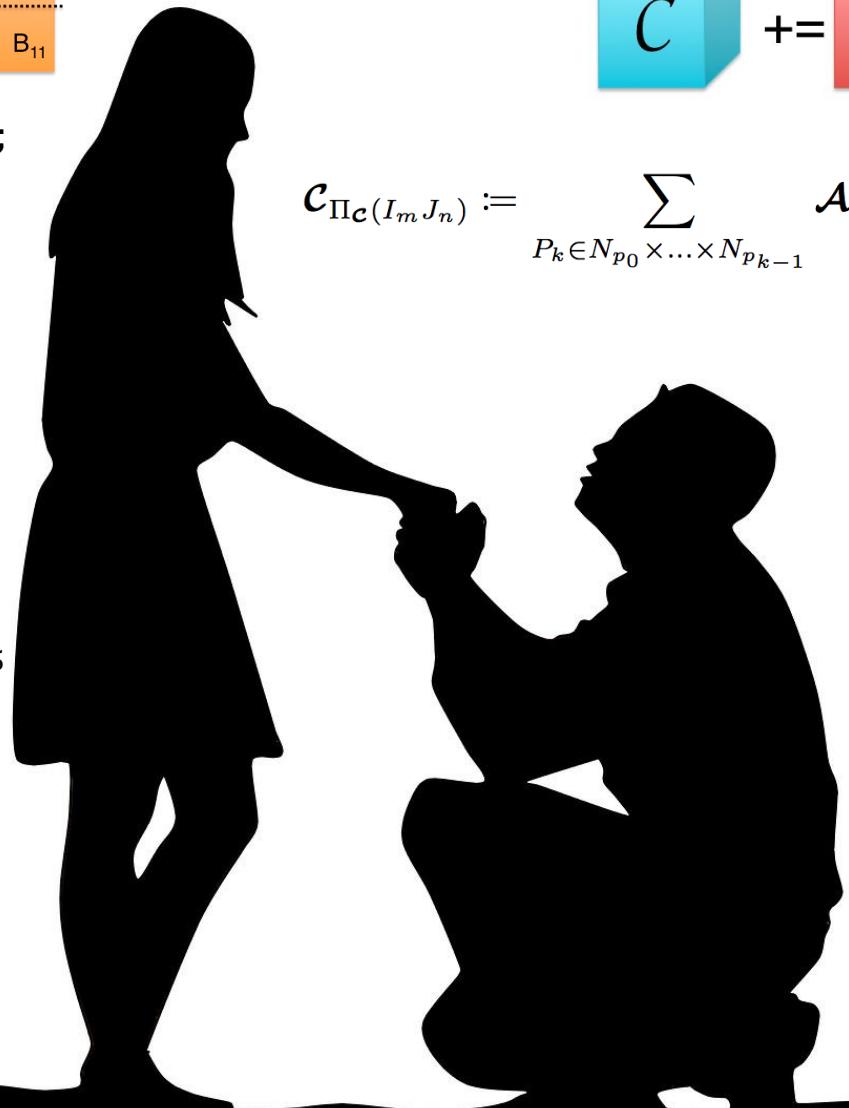
Flatiron Institute, 162 Fifth Ave, New York City.

# Marry Strassen with Tensor Contraction ?



$$\begin{aligned}M_0 &:= (A_{00}+A_{11})(B_{00}+B_{11}); \\M_1 &:= (A_{10}+A_{11})B_{00}; \\M_2 &:= A_{00}(B_{01}-B_{11}); \\M_3 &:= A_{11}(B_{10}-B_{00}); \\M_4 &:= (A_{00}+A_{01})B_{11}; \\M_5 &:= (A_{10}-A_{00})(B_{00}+B_{01}); \\M_6 &:= (A_{01}-A_{11})(B_{10}+B_{11}); \\C_{00} &+= M_0 + M_3 - M_4 + M_6 \\C_{01} &+= M_2 + M_4 \\C_{10} &+= M_1 + M_3 \\C_{11} &+= M_0 - M_1 + M_2 + M_5\end{aligned}$$

$O(n^3) \rightarrow O(n^{2.8})$



$$C_{\Pi_C(I_m J_n)} := \sum_{P_k \in N_{p_0} \times \dots \times N_{p_{k-1}}} A_{\Pi_A(I_m P_k)} \cdot B_{\Pi_B(P_k J_n)} + C_{\Pi_C(I_m J_n)}$$

Practical Speedup?

# Outline

- Background
  - High-performance GEMM
  - High-performance Strassen
  - High-performance Tensor Contraction
- Strassen's Algorithm for Tensor Contraction
- Performance Model
- Experiments
- Conclusion

PROPOSAL

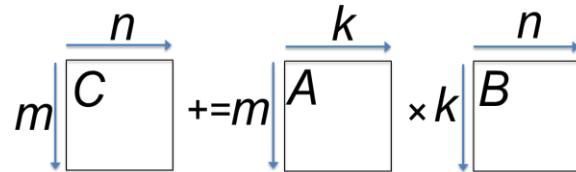


MARRIAGE



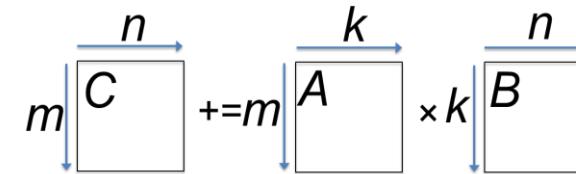
# High-performance matrix multiplication (GEMM)

# State-of-the-art GEMM in BLIS



- BLAS-like Library Instantiation Software (**BLIS**) is a portable framework for instantiating BLAS-like dense linear algebra libraries.
  - ❑ Field Van Zee, and Robert van de Geijn. “BLIS: A Framework for Rapidly Instantiating BLAS Functionality.” *ACM TOMS* 41.3 (2015): 14.
- BLIS provides a refactoring of **GotoBLAS** algorithm (best-known approach on CPU) to implement **GEMM**.
  - ❑ Kazushige Goto, and Robert van de Geijn. “High-performance implementation of the level-3 BLAS.” *ACM TOMS* 35.1 (2008): 4.
  - ❑ Kazushige Goto, and Robert van de Geijn. “Anatomy of high-performance matrix multiplication.” *ACM TOMS* 34.3 (2008): 12.
- GEMM implementation in BLIS has 6-layers of loops. The outer 5 loops are written in **C**. The inner-most loop (micro-kernel) is written in **assembly** for high performance.

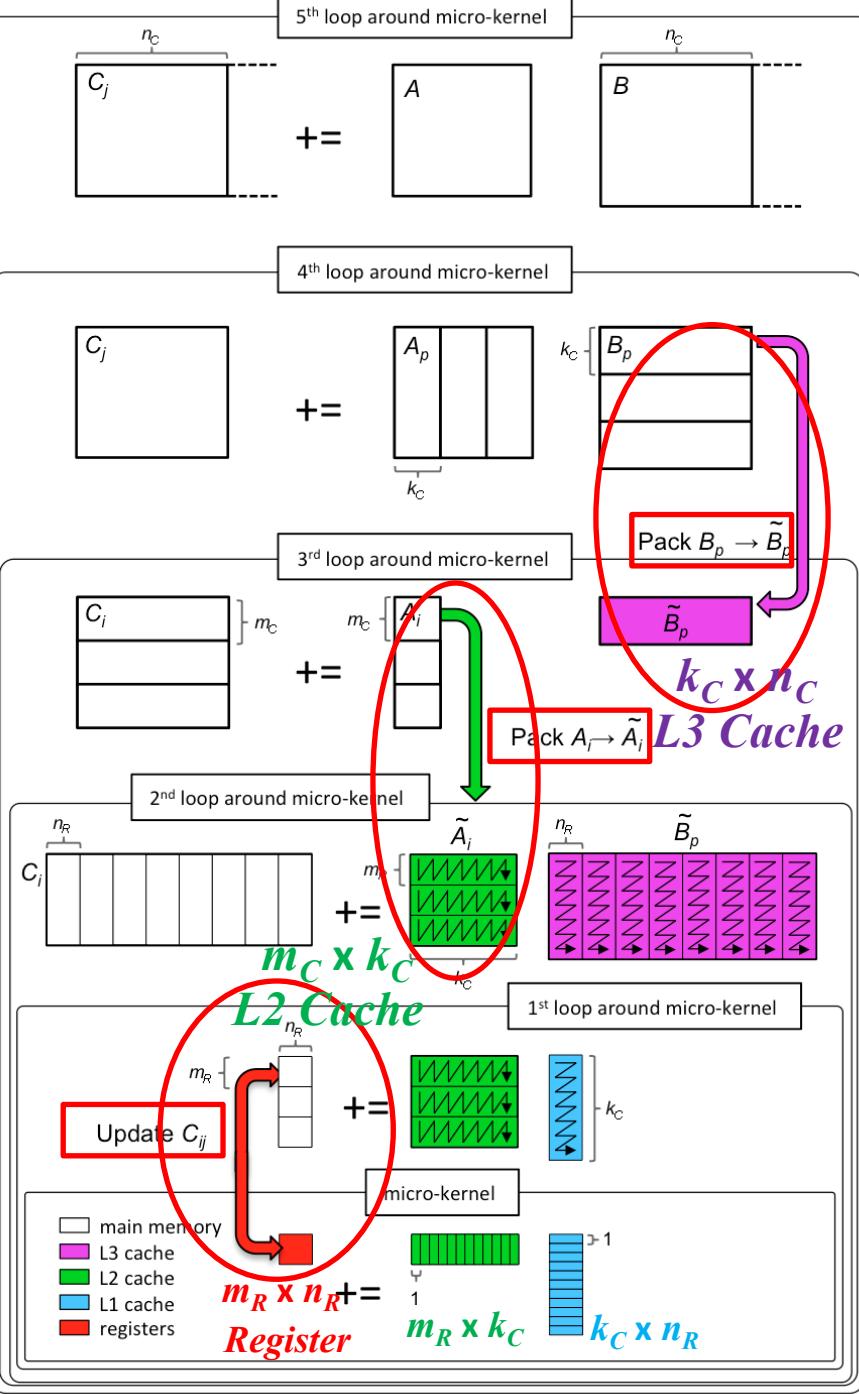
# GotoBLAS algorithm for GEMM in BLIS



```

n      k      n
|---| |---| |---|
| C | | A | | B |
|---| |---| |---|
m      + = m   k   n
|---| |---| |---|
| C | + = | A | x | B |
|---| |---| |---|
Loop 5  for  $j_c = 0 : n - 1$  steps of  $n_c$ 
         $\mathcal{J}_c = j_c : j_c + n_c - 1$ 
        for  $p_c = 0 : k - 1$  steps of  $k_c$ 
             $\mathcal{P}_c = p_c : p_c + k_c - 1$ 
             $B(\mathcal{P}_c, \mathcal{J}_c) \rightarrow \tilde{B}_p$ 
            for  $i_c = 0 : m - 1$  steps of  $m_c$ 
                 $\mathcal{I}_c = i_c : i_c + m_c - 1$ 
                 $A(\mathcal{I}_c, \mathcal{P}_c) \rightarrow \tilde{A}_i$ 
                // macro-kernel
                for  $j_r = 0 : n_c - 1$  steps of  $n_r$ 
                     $\mathcal{J}_r = j_r : j_r + n_r - 1$ 
                    for  $i_r = 0 : m_c - 1$  steps of  $m_r$ 
                         $\mathcal{I}_r = i_r : i_r + m_r - 1$ 
                        // micro-kernel
                        for  $p_r = 0 : p_c - 1$  steps of 1
                             $C_c(\mathcal{I}_r, \mathcal{J}_r) += \tilde{A}_i(\mathcal{I}_r, p_r) \tilde{B}_p(p_r, \mathcal{J}_r)$ 
                        endfor
                    endfor
                endfor
            endfor
        endfor
    endfor
endfor

```

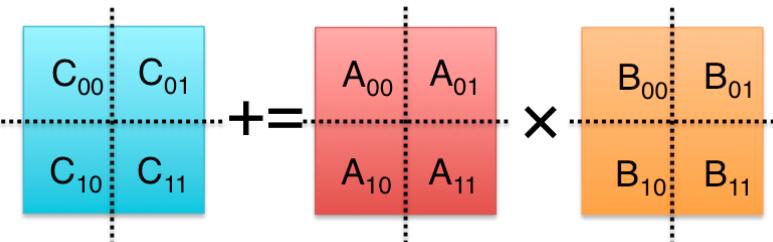


\*Field G. Van Zee, and Tyler M. Smith. “Implementing high-performance complex matrix multiplication via the 3m and 4m methods.” In ACM Transactions on Mathematical Software (TOMS), accepted.

# High-performance Strassen

# Strassen's Algorithm Reloaded

$M_0 := (A_{00} + A_{11})(B_{00} + B_{11});$   
 $M_1 := (A_{10} + A_{11})B_{00};$   
 $M_2 := A_{00}(B_{01} - B_{11});$   
 $M_3 := A_{11}(B_{10} - B_{00});$   
 $M_4 := (A_{00} + A_{01})B_{11};$   
 $M_5 := (A_{10} - A_{00})(B_{00} + B_{01});$   
 $M_6 := (A_{01} - A_{11})(B_{10} + B_{11});$   
 $C_{00} += M_0 + M_3 - M_4 + M_6$   
 $C_{01} += M_2 + M_4$   
 $C_{10} += M_1 + M_3$   
 $C_{11} += M_0 - M_1 + M_2 + M_5$



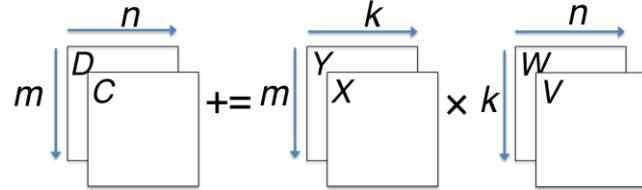
General operation for one-level Strassen:

$M_0 := (A_{00} + A_{11})(B_{00} + B_{11}); \quad C_{00} += M_0; \quad C_{11} += M_0;$   
 $M_1 := (A_{10} + A_{11})B_{00}; \quad C_{10} += M_1; \quad C_{11} -= M_1;$   
 $M_2 := A_{00}(B_{01} - B_{11}); \quad C_{01} += M_2; \quad C_{11} += M_2;$   
 $M_3 := A_{11}(B_{10} - B_{00}); \quad C_{00} += M_3; \quad C_{10} += M_3;$   
 $M_4 := (A_{00} + A_{01})B_{11}; \quad C_{01} += M_4; \quad C_{00} -= M_4;$   
 $M_5 := (A_{10} - A_{00})(B_{00} + B_{01}); \quad C_{11} += M_5;$   
 $M_6 := (A_{01} - A_{11})(B_{10} + B_{11}); \quad C_{00} += M_6;$

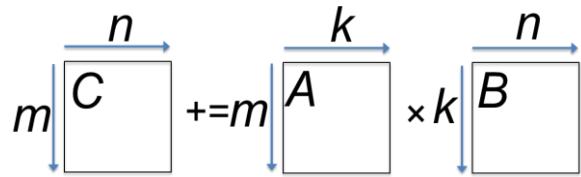
$M := (X + Y)(V + W); \quad C += M; \quad D += M;$

$M := (X + \delta Y)(V + \varepsilon W); \quad C += \gamma_0 M; \quad D += \gamma_1 M;$   
 $\gamma_0, \gamma_1, \delta, \varepsilon \in \{-1, 0, 1\}.$

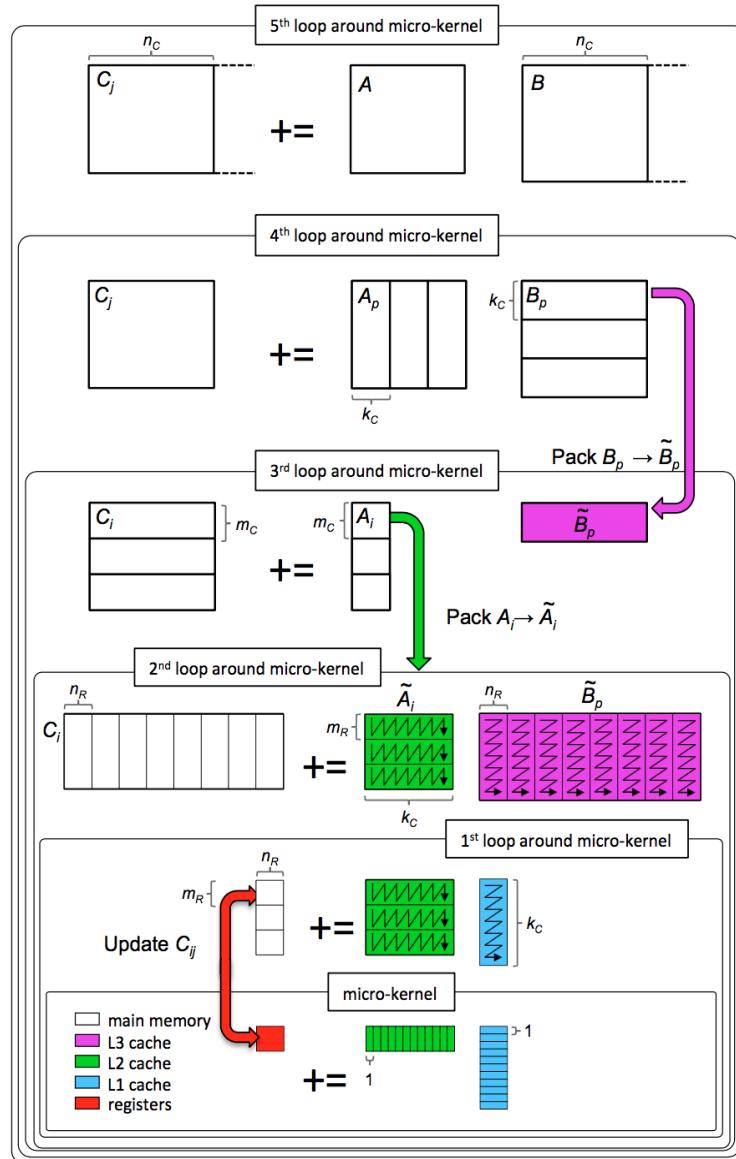
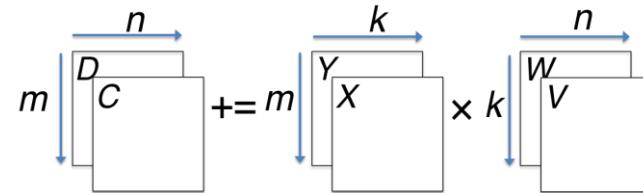
$$M := (X+Y)(V+W); \quad C += M; \quad D += M;$$



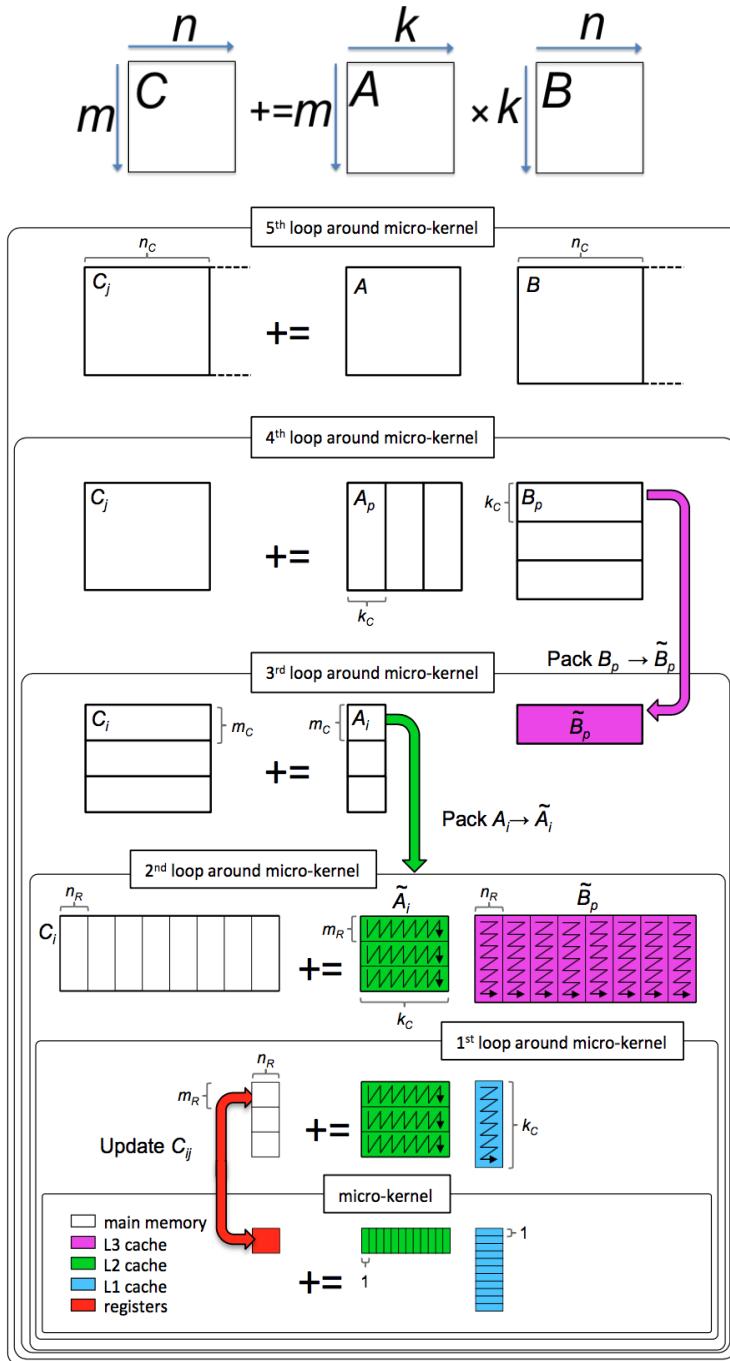
$$C += AB;$$



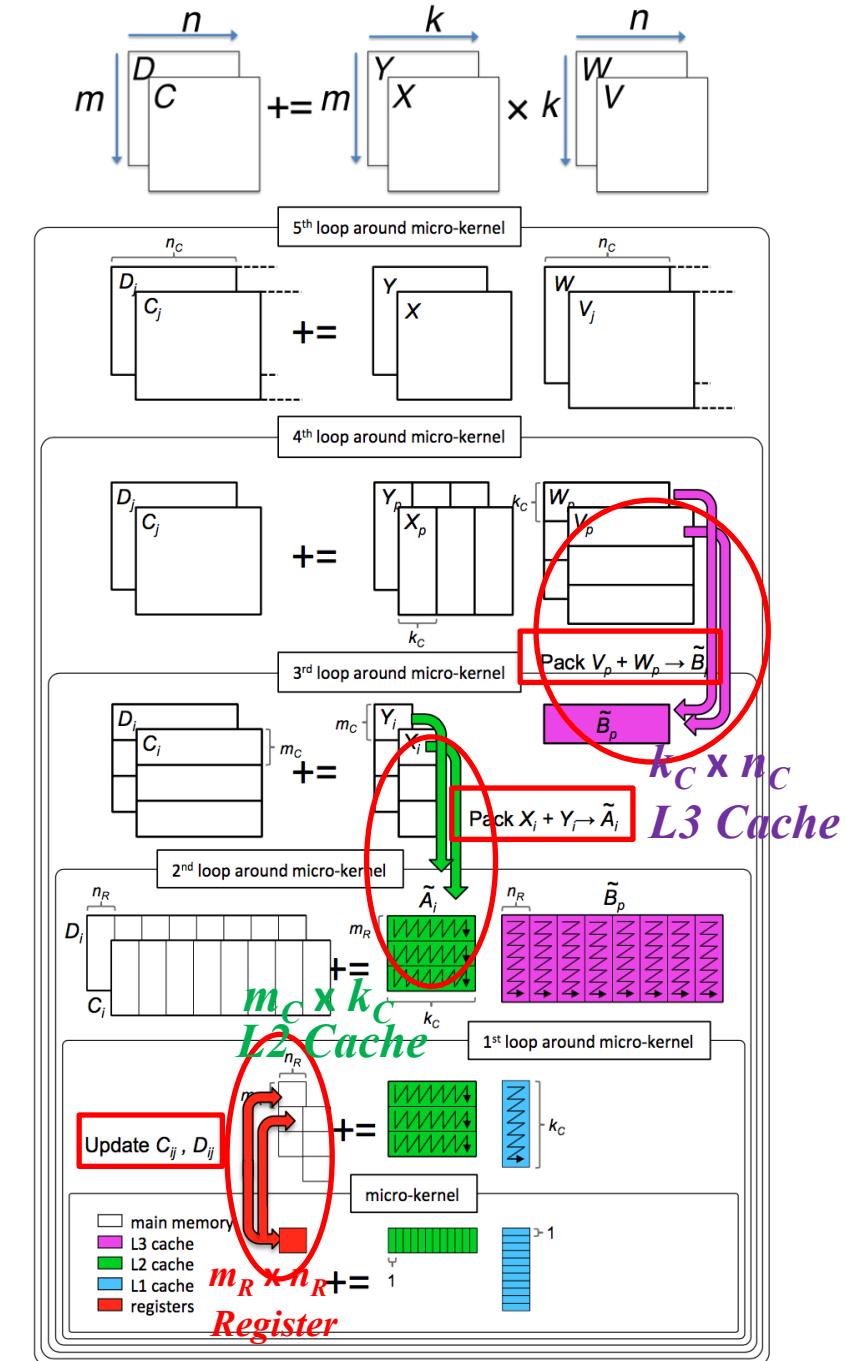
$$M := (X+Y)(V+W); \quad C += M; \quad D += M;$$



$$C += AB;$$



$$M := (X+Y)(V+W); \quad C += M; \quad D += M;$$



# High-performance Tensor Contraction

# Matrix vs. Tensor

## Matrix Multiplication

$$\begin{matrix} \text{C} \\ \text{A} \\ \text{B} \end{matrix} \quad += \quad \times$$

$$C += AB$$

$$C_{i,j} += \sum_k A_{i,k} B_{k,j}$$

BLAS/BLIS!

## Tensor Contraction

$$\begin{matrix} \text{C} \\ \mathcal{A} \\ \mathcal{B} \end{matrix} \quad += \quad \times$$

$$C += \mathcal{A}\mathcal{B}$$

$$C_{a,b,c} += \sum_d \mathcal{A}_{d,c,a} \mathcal{B}_{d,b}$$

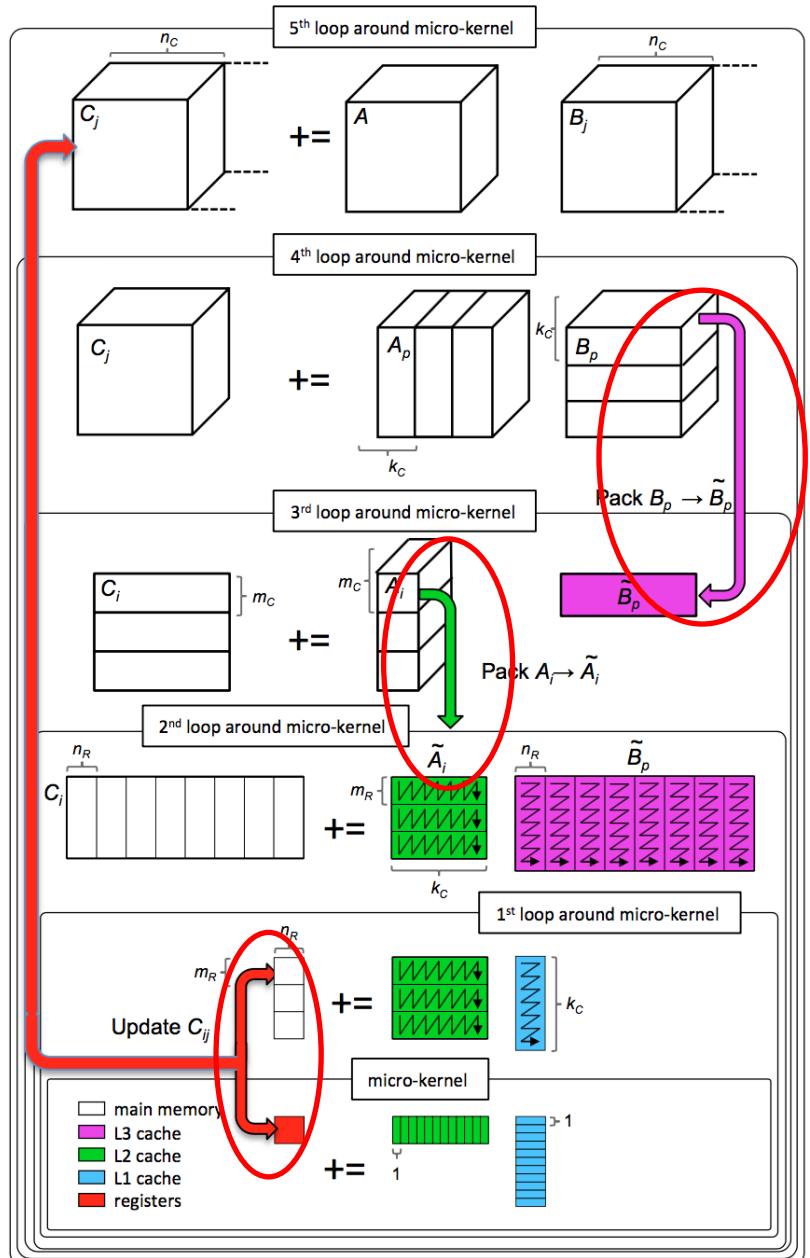
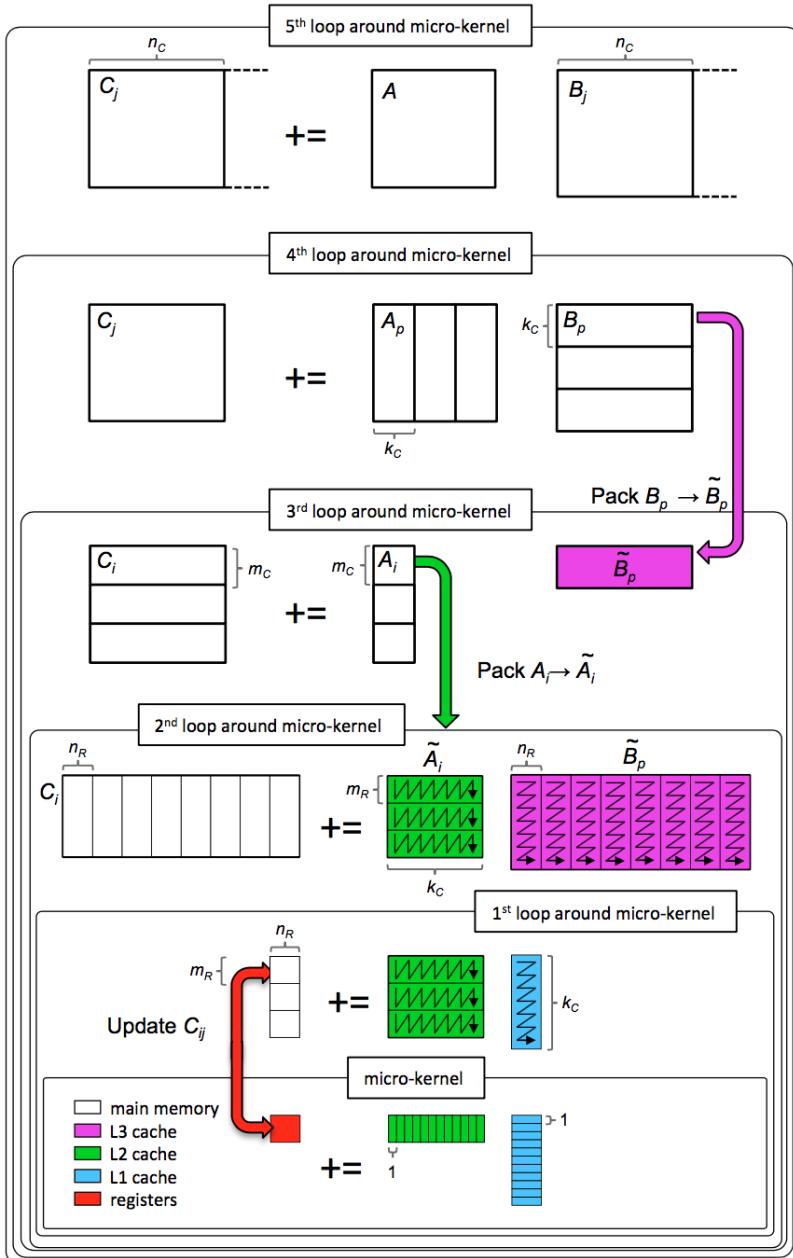
TBLIS/GETT

Paul Springer, and Paolo Bientinesi. "Design of a high-performance GEMM-like tensor-tensor multiplication." arXiv preprint arXiv:1607.00145 (2016).

Devin A. Matthews. "High-Performance Tensor Contraction without Transposition." Accepted in *SISC*.

$$C := AB + C$$

$$\mathcal{C}_{\Pi_{\mathbf{C}}(I_m J_n)} := \sum_{P_k \in N_{p_0} \times \dots \times N_{p_{k-1}}} \mathcal{A}_{\Pi_{\mathbf{A}}(I_m P_k)} \cdot \mathcal{B}_{\Pi_{\mathbf{B}}(P_k J_n)} + \mathcal{C}_{\Pi_{\mathbf{C}}(I_m J_n)}$$



# Outline

- Background
  - High-performance GEMM
  - High-performance Strassen
  - High-performance Tensor Contraction
- **Strassen's Algorithm for Tensor Contraction**
- Performance Model
- Experiments
- Conclusion



# Matrix vs. Tensor

## Matrix Multiplication

$$\begin{matrix} \text{C} \\ \text{A} \\ \text{B} \end{matrix} \quad += \quad \times$$

$$C += AB$$

$$C_{i,j} += \sum_k A_{i,k} B_{k,j}$$

BLAS/BLIS!

## Tensor Contraction

$$\begin{matrix} \text{C} \\ \mathcal{A} \\ \mathcal{B} \end{matrix} \quad += \quad \times$$

$$C += \mathcal{A}\mathcal{B}$$

$$C_{a,b,c} += \sum_d \mathcal{A}_{d,c,a} \mathcal{B}_{d,b}$$

TBLIS/GETT

Paul Springer, and Paolo Bientinesi. "Design of a high-performance GEMM-like tensor-tensor multiplication." arXiv preprint arXiv:1607.00145 (2016).

Devin A. Matthews. "High-Performance Tensor Contraction without Transposition." Accepted in *SISC*.

# Matrix vs. Tensor

Matrix Multiplication

$$\begin{matrix} \text{C} \\ \text{A} \\ \text{B} \end{matrix} \quad += \quad \times$$

$$C += AB$$

$$C_{i,j} += \sum_k A_{i,k} B_{k,j}$$

BLAS/BLIS!

$$\begin{array}{|c|c|} \hline C_0 & C_1 \\ \hline C_2 & C_3 \\ \hline \end{array} \quad += \quad \begin{array}{|c|c|} \hline A_0 & A_1 \\ \hline A_2 & A_3 \\ \hline \end{array} \quad \times \quad \begin{array}{|c|c|} \hline B_0 & B_1 \\ \hline B_2 & B_3 \\ \hline \end{array}$$

Tensor Contraction

$$\begin{matrix} \mathcal{C} \\ \mathcal{A} \\ \mathcal{B} \end{matrix} \quad += \quad \times$$

$$\mathcal{C} += \mathcal{A}\mathcal{B}$$

$$C_{a,b,c} += \sum_d \mathcal{A}_{d,c,a} \mathcal{B}_{d,b}$$

TBLIS/GETT

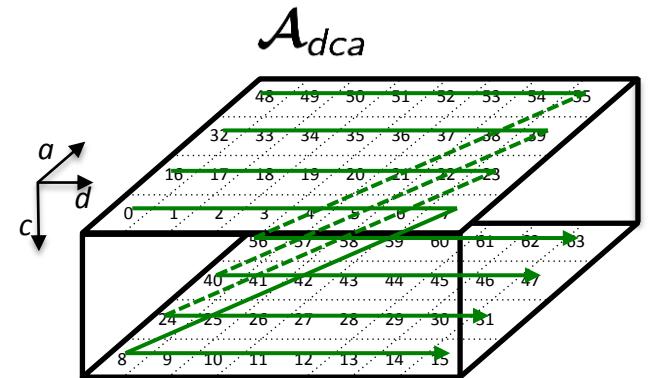
$$\begin{array}{|c|c|} \hline C_0 & C_1 \\ \hline C_2 & C_3 \\ \hline \end{array} \quad += \quad \begin{array}{|c|c|} \hline \mathcal{A}_0 & \mathcal{A}_1 \\ \hline \mathcal{A}_2 & \mathcal{A}_3 \\ \hline \end{array} \quad \times \quad \begin{array}{|c|c|} \hline \mathcal{B}_0 & \mathcal{B}_1 \\ \hline \mathcal{B}_2 & \mathcal{B}_3 \\ \hline \end{array}$$

# Tensors As Matrices: Block-Scatter-Matrix View

➤ **Tensor:**  $\mathcal{A}_{dca}$ , 8x2x4

with  $N_a = 4, N_c = 2, N_d = 8$

□ “d” dimension is stride-1, other dimensions have increasing strides (8, 16).



# Tensors As Matrices: Block-Scatter-Matrix View

- **Tensor:**  $\mathcal{A}_{dca}$ , 8x2x4

with  $N_a = 4, N_c = 2, N_d = 8$

- “d” dimension is stride-1, other dimensions have increasing strides (8, 16).

- **Matrix:**  $A_{(ac)(d)}$ , 8x8

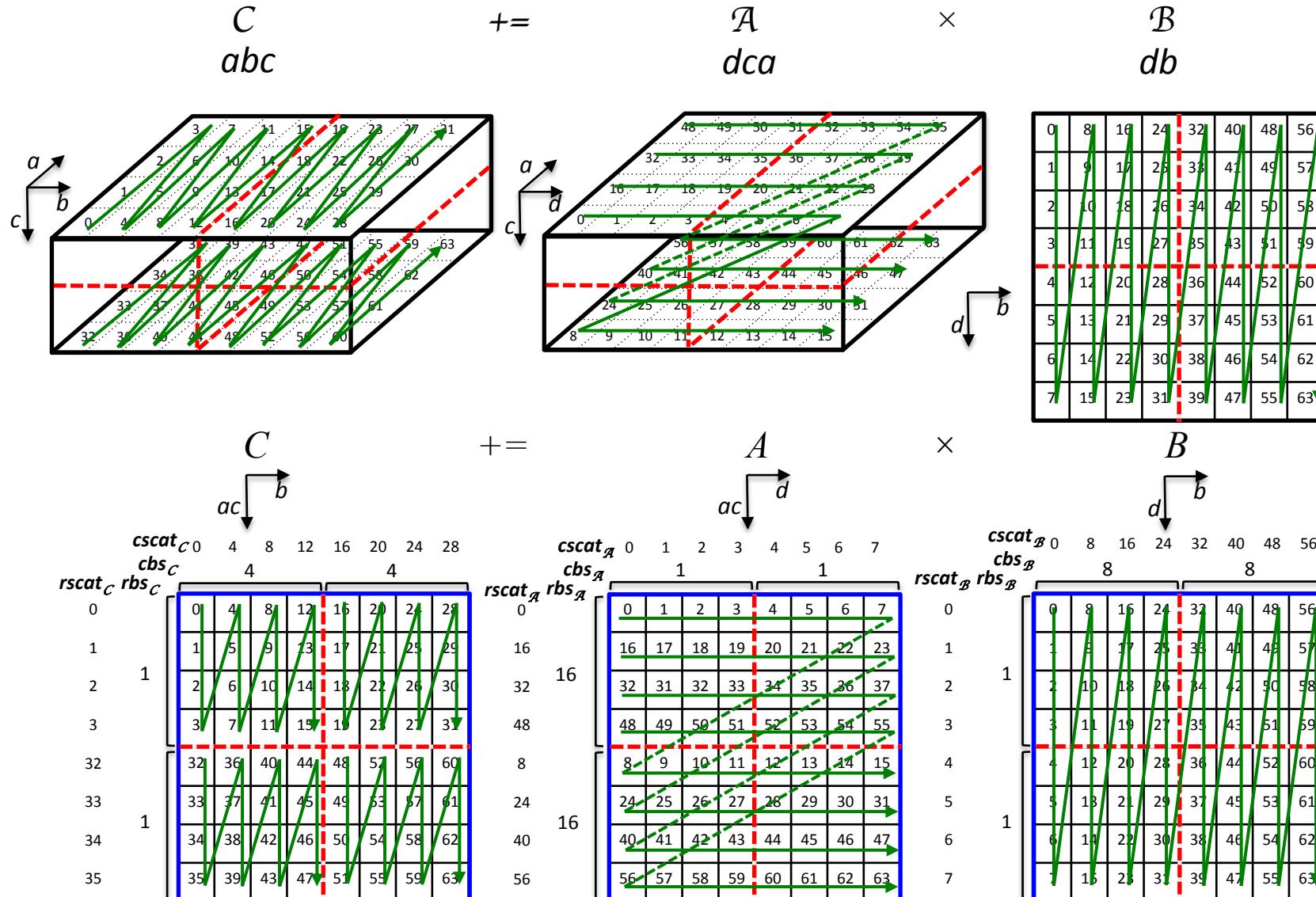
with  $N_{I_m} = N_a \cdot N_c = 8, N_{P_k} = N_d = 8$

- Column “ac” dimension has stride of “c” (8x2=16).  
Row “d” dimension has stride-1. (i.e. A is row-major.)
- $rscat_{\mathcal{A}}, cscat_{\mathcal{A}}$  store offset for each position in rows or columns: **Scatter-Matrix Vector**  
 $\text{OFFSET}_{\mathcal{A}_{d,c,a}} = rscat_{A;(ac)} + cscat_{A;(d)}$
- $rbs_{\mathcal{A}}, cbs_{\mathcal{A}}$  store stride for each block or zero for irregular blocks: **Block-Scatter-Matrix Vector**
  - vector load/store instructions for stride-1 index
  - vector gather/scatter instructions for stride-n index.

Devin A. Matthews. “High-Performance Tensor Contraction without Transposition.” Accepted in SISC.

19

# Strassen's Algorithm for Tensor Contraction



Jianyu Huang, Devin A. Matthews, and Robert A. van de Geijn. "Strassen's Algorithm for Tensor Contraction." arXiv:1704.03092 (2017).

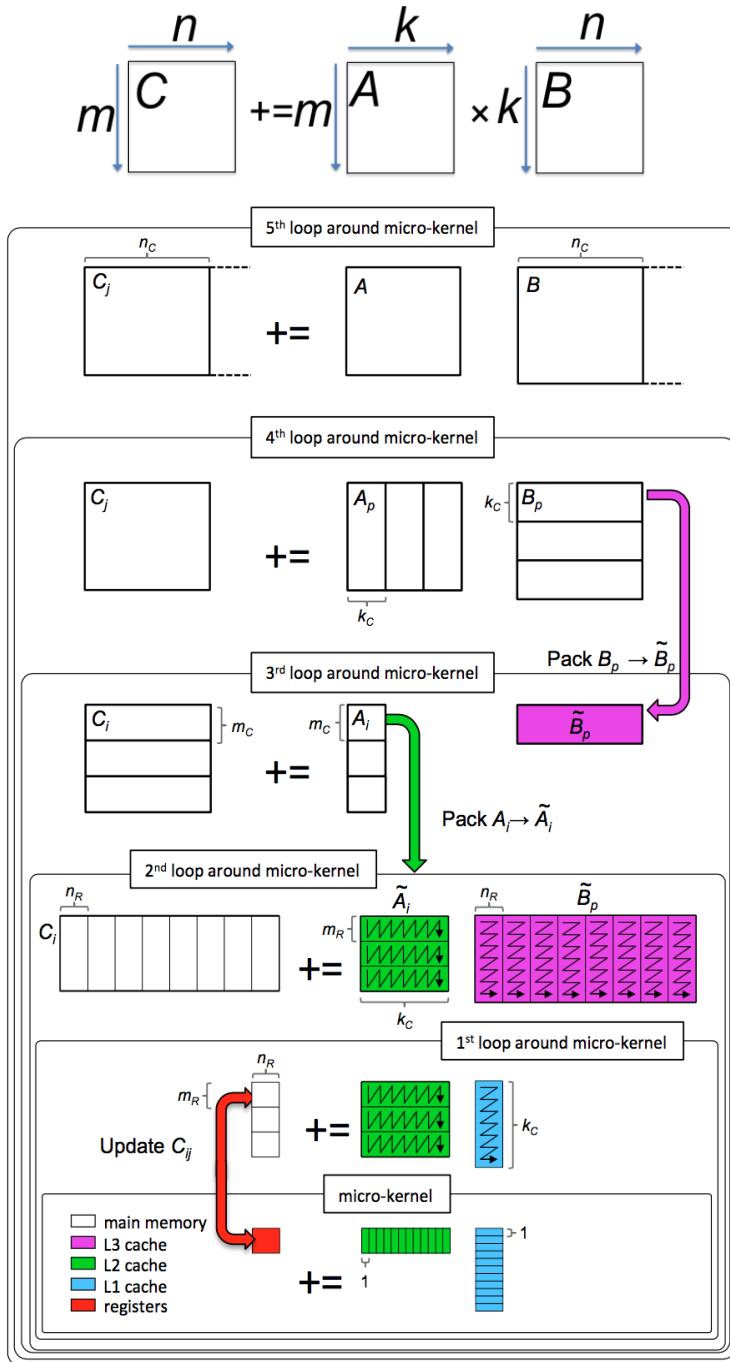
# Modifications to GEMM

$$M_0 := (A_{00} + A_{11})(B_{00} + B_{11}); C_{00} += M_0; C_{11} += M_0;$$

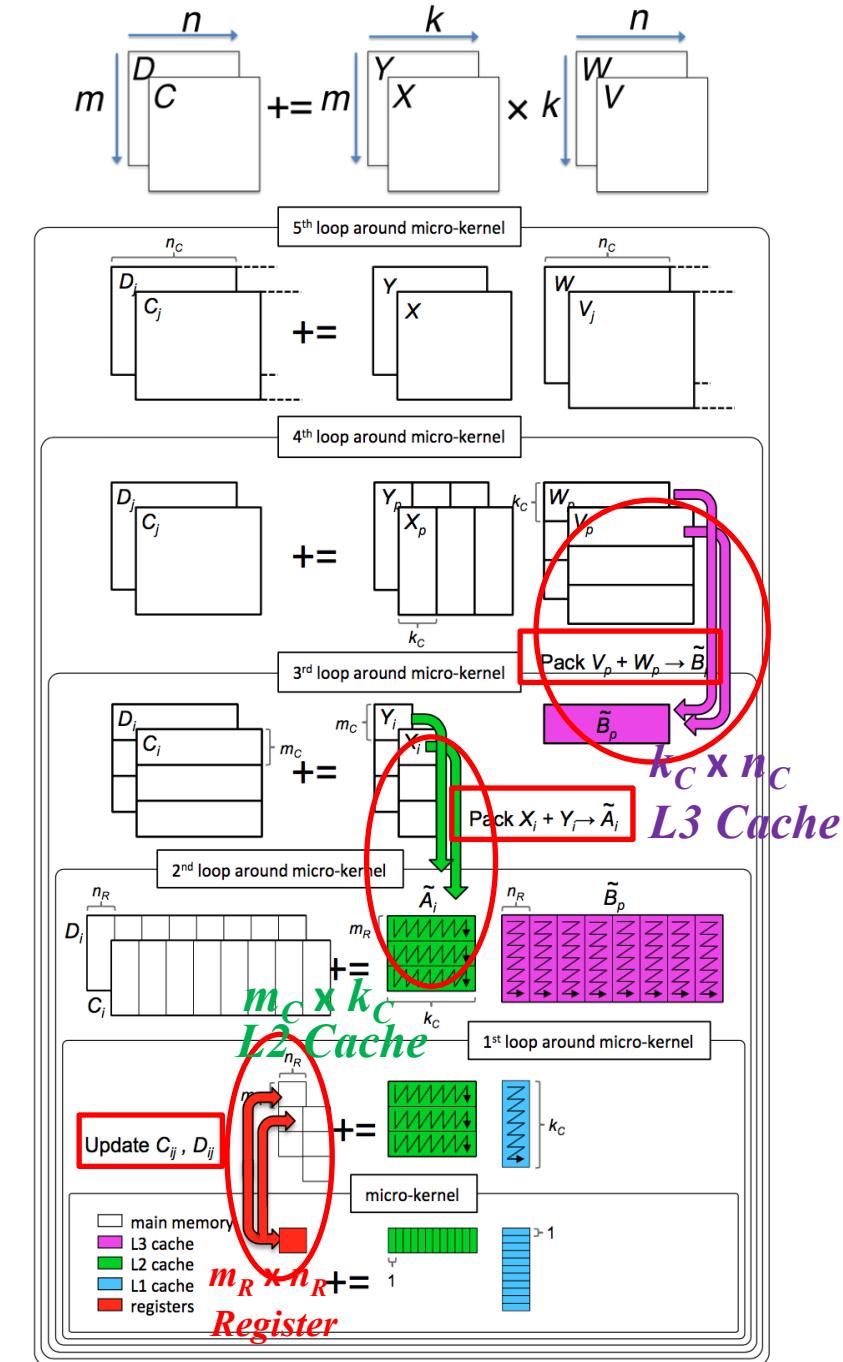
- Packing routines:
  - Implicit tensor-to-matrix permutations
  - Addition of submatrices of A and B.
- Micro-kernel:
  - Implicit matrix-to-tensor transformations
  - Scatter update of submatrices of C.

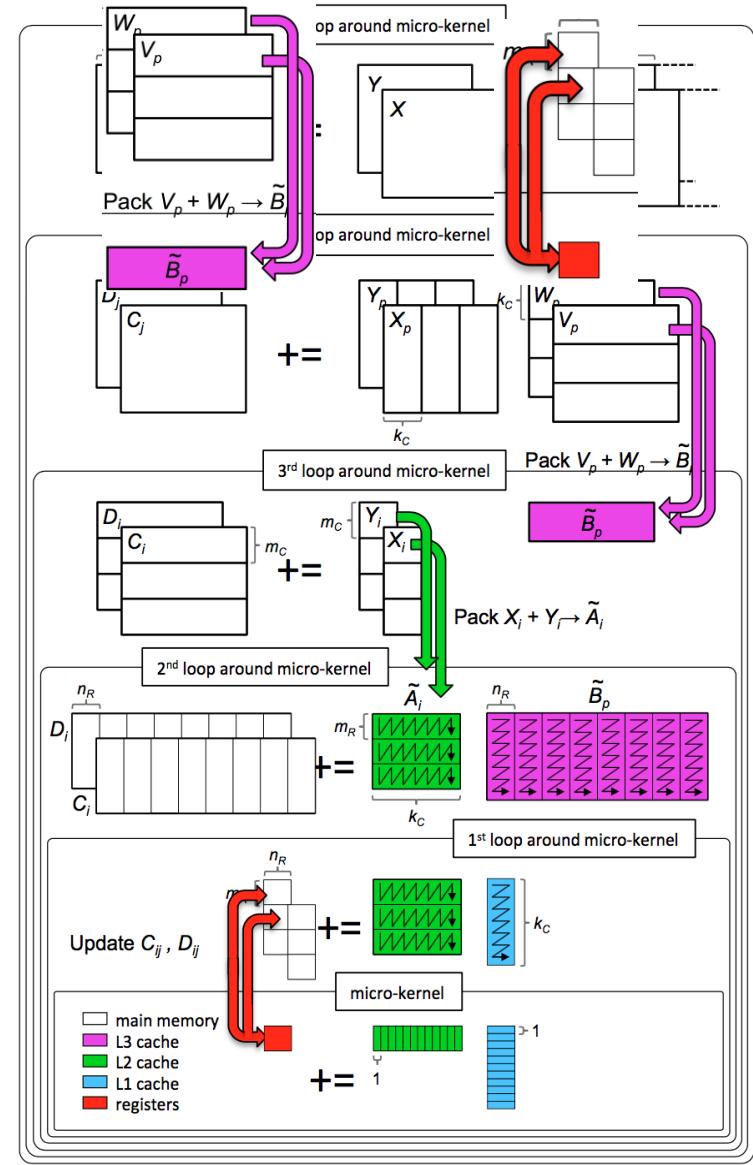
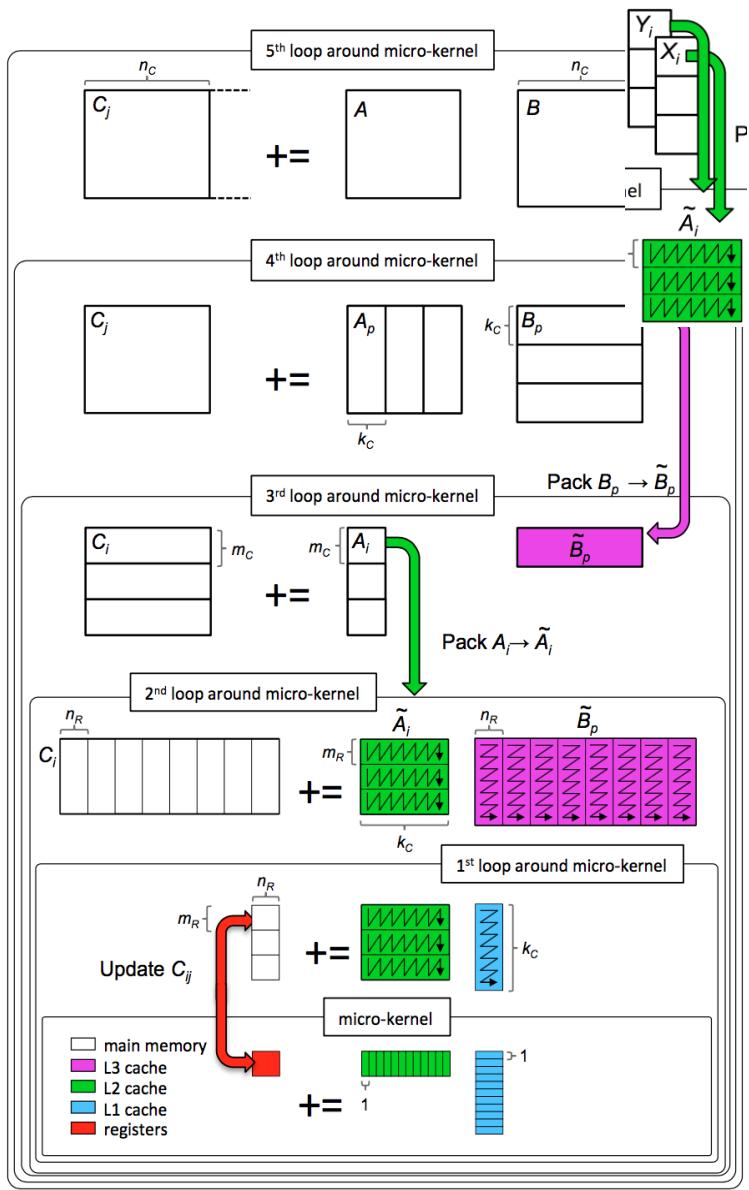
- 
- Additional workspace for Transposition (Tensor Contraction)
  - Additional Workspace for Summation (Strassen)

$$C += AB;$$

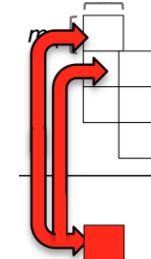
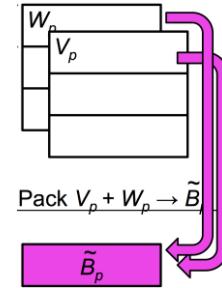
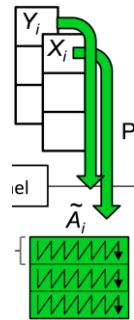
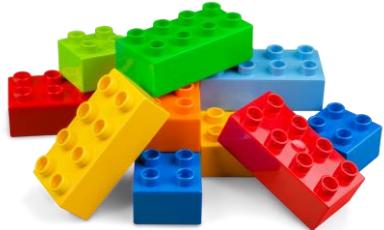


$$M := (X+Y)(V+W); \quad C += M; \quad D += M;$$





# Variations on a theme



- Naïve Strassen

✗

✗

✗

- AB Strassen

✓

✓

✗

- ABC Strassen

✓

✓

✓

# Outline

- Background
  - High-performance GEMM
  - High-performance Strassen
  - High-performance Tensor Contraction
- Strassen's Algorithm for Tensor Contraction
- Performance Model
- Experiments
- Conclusion



# Performance Model

- Performance Metric

$$\mathcal{C}_{\Pi_{\mathcal{C}}(I_m J_n)} := \sum_{P_k \in N_{p_0} \times \dots \times N_{p_{k-1}}} \mathcal{A}_{\Pi_{\mathcal{A}}(I_m P_k)} \cdot \mathcal{B}_{\Pi_{\mathcal{B}}(P_k J_n)} + \mathcal{C}_{\Pi_{\mathcal{C}}(I_m J_n)}$$

$$N_{I_m} = \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}},$$

$$N_{J_n} = \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}},$$

$$N_{P_k} = \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}.$$

$$\text{Effective GFLOPS} = \frac{2 \cdot N_{I_m} \cdot N_{J_n} \cdot N_{P_k}}{\text{T (in seconds)}} \cdot 10^{-9}$$

- Total Time Breakdown

$$T = T_a + T_m$$

↓      ↓

Arithmetic Operations      Memory Operations

$$T_a = W_a^\times \cdot T_a^\times + W_a^{\mathcal{A}+} \cdot T_a^{\mathcal{A}+} + W_a^{\mathcal{B}+} \cdot T_a^{\mathcal{B}+} + W_a^{\mathcal{C}+} \cdot T_a^{\mathcal{C}+}$$

	type	$\tau$	TBLIS	$L$ -level
$T_a^\times$	-	$\tau_a$	$2N_{I_m} N_{J_n} N_{P_k}$	$2 \frac{N_{I_m}}{2^L} \frac{N_{J_n}}{2^L} \frac{N_{P_k}}{2^L}$
$T_a^{\mathcal{A}+}$	-	$\tau_a$	-	$2 \frac{N_{I_m}}{2^L} \frac{N_{P_k}}{2^L}$
$T_a^{\mathcal{B}+}$	-	$\tau_a$	-	$2 \frac{N_{P_k}}{2^L} \frac{N_{J_n}}{2^L}$
$T_a^{\mathcal{C}+}$	-	$\tau_a$	-	$2 \frac{N_{I_m}}{2^L} \frac{N_{J_n}}{2^L}$

	TBLIS	1-level			2-level		
		ABC	AB	Naive	ABC	AB	Naive
$W_a^\times$	1	7	7	7	49	49	49
$W_a^{\mathcal{A}+}$	-	5	5	5	95	95	95
$W_a^{\mathcal{B}+}$	-	5	5	5	95	95	95
$W_a^{\mathcal{C}+}$	-	12	12	12	144	144	144

$$\begin{aligned} N_{I_m} &= \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}}, \\ N_{J_n} &= \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}}, \\ N_{P_k} &= \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}. \end{aligned}$$

$$\begin{aligned}
T_a &= W_a^\times \cdot T_a^\times + W_a^{\mathcal{A}+} \cdot T_a^{\mathcal{A}+} + W_a^{\mathcal{B}+} \cdot T_a^{\mathcal{B}+} + W_a^{\mathcal{C}+} \cdot T_a^{\mathcal{C}+} \\
T_m &= W_m^{\mathcal{A}\times} \cdot T_m^{\mathcal{A}\times} + W_m^{\mathcal{B}\times} \cdot T_m^{\mathcal{B}\times} + W_m^{\mathcal{C}\times} \cdot T_m^{\mathcal{C}\times} + W_m^{\mathcal{A}+} \cdot T_m^{\mathcal{A}+} + W_m^{\mathcal{B}+} \cdot T_m^{\mathcal{B}+} + W_m^{\mathcal{C}+} \cdot T_m^{\mathcal{C}+}
\end{aligned}$$

	type	$\tau$	TBLIS	$L$ -level
$T_a^\times$	-	$\tau_a$	$2N_{I_m} N_{J_n} N_{P_k}$	$2 \frac{N_{I_m}}{2^L} \frac{N_{J_n}}{2^L} \frac{N_{P_k}}{2^L}$
$T_a^{\mathcal{A}+}$	-	$\tau_a$	-	$2 \frac{N_{I_m}}{2^L} \frac{N_{P_k}}{2^L}$
$T_a^{\mathcal{B}+}$	-	$\tau_a$	-	$2 \frac{N_{P_k}}{2^L} \frac{N_{J_n}}{2^L}$
$T_a^{\mathcal{C}+}$	-	$\tau_a$	-	$2 \frac{N_{I_m}}{2^L} \frac{N_{J_n}}{2^L}$
$T_m^{\mathcal{A}\times}$	r	$\tau_b$	$N_{I_m} N_{P_k} \lceil \frac{N_{J_n}}{n_c} \rceil$	$\frac{N_{I_m}}{2^L} \frac{N_{P_k}}{2^L} \lceil \frac{N_{J_n}/2^L}{n_c} \rceil$
$T_m^{\tilde{\mathcal{A}}\times}$	w	$\tau_b$	$N_{I_m} N_{P_k} \lceil \frac{N_{J_n}}{n_c} \rceil$	$\frac{N_{I_m}}{2^L} \frac{N_{P_k}}{2^L} \lceil \frac{N_{J_n}/2^L}{n_c} \rceil$
$T_m^{\mathcal{B}\times}$	r	$\tau_b$	$N_{J_n} N_{P_k}$	$\frac{N_{J_n}}{2^L} \frac{N_{P_k}}{2^L}$
$T_m^{\tilde{\mathcal{B}}\times}$	w	$\tau_b$	$N_{J_n} N_{P_k}$	$\frac{N_{J_n}}{2^L} \frac{N_{P_k}}{2^L}$
$T_m^{\mathcal{C}\times}$	r/w	$\tau_b$	$2\lambda N_{I_m} N_{J_n} \lceil \frac{N_{P_k}}{k_c} \rceil$	$2\lambda \frac{N_{I_m}}{2^L} \frac{N_{J_n}}{2^L} \lceil \frac{N_{P_k}/2^L}{k_c} \rceil$
$T_m^{\mathcal{A}+}$	r/w	$\tau_b$	$N_{I_m} N_{P_k}$	$\frac{N_{I_m}}{2^L} \frac{N_{P_k}}{2^L}$
$T_m^{\mathcal{B}+}$	r/w	$\tau_b$	$N_{J_n} N_{P_k}$	$\frac{N_{J_n}}{2^L} \frac{N_{P_k}}{2^L}$
$T_m^{\mathcal{C}+}$	r/w	$\tau_b$	$N_{I_m} N_{J_n}$	$\frac{N_{I_m}}{2^L} \frac{N_{J_n}}{2^L}$

	TBLIS	1-level			2-level		
		ABC	AB	Naive	ABC	AB	Naive
$W_a^\times$	1	7	7	7	49	49	49
$W_a^{\mathcal{A}+}$	-	5	5	5	95	95	95
$W_a^{\mathcal{B}+}$	-	5	5	5	95	95	95
$W_a^{\mathcal{C}+}$	-	12	12	12	144	144	144
$W_m^{\mathcal{A}\times}$	1	12	12	7	194	194	49
$W_m^{\tilde{\mathcal{A}}\times}$	-	-	-	-	-	-	-
$W_m^{\mathcal{B}\times}$	1	12	12	7	194	194	49
$W_m^{\tilde{\mathcal{B}}\times}$	-	-	-	-	-	-	-
$W_m^{\mathcal{C}\times}$	1	12	7	7	144	49	49
$W_m^{\mathcal{A}+}$	-	-	-	19	-	-	293
$W_m^{\mathcal{B}+}$	-	-	-	19	-	-	293
$W_m^{\mathcal{C}+}$	-	-	36	36	-	432	432

$$\begin{aligned}
N_{I_m} &= \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}}, \\
N_{J_n} &= \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}}, \\
N_{P_k} &= \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}.
\end{aligned}$$

# Outline

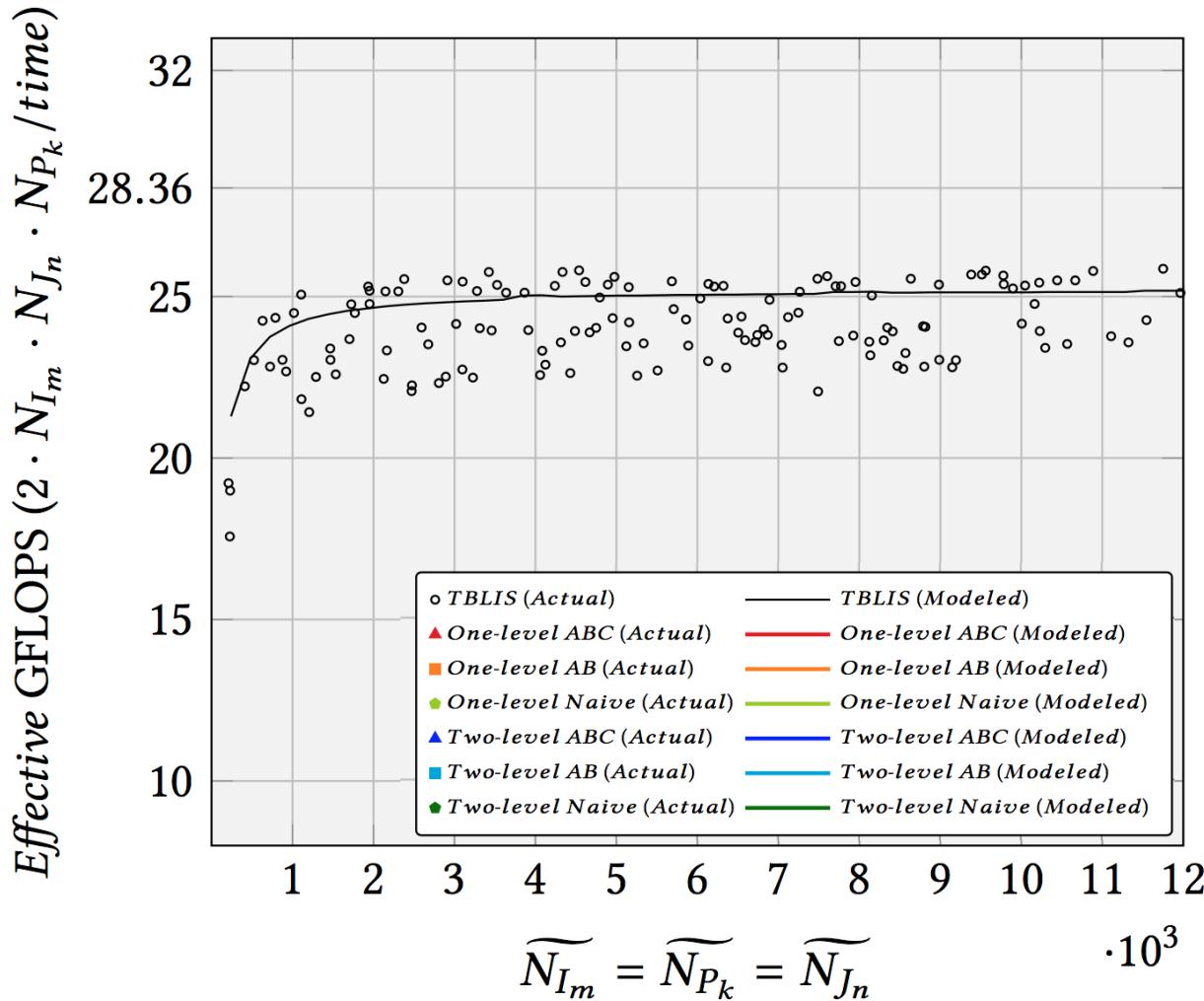
- Background
  - High-performance GEMM
  - High-performance Strassen
  - High-performance Tensor Contraction
- Strassen's Algorithm for Tensor Contraction
- Performance Model
- Experiments
- Conclusion



# Synthetic Tensor Contractions

Intel Xeon E5-2680 v2 (Ivy Bridge, 10 core/socket)

$$\widetilde{N_{I_m}} = \widetilde{N_{P_k}} = \widetilde{N_{J_n}}, \text{ 1 core, Actual vs. Modeled}$$

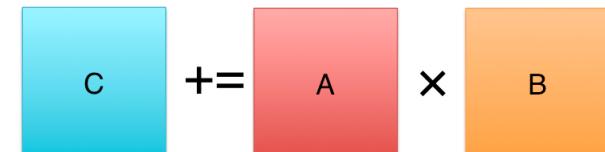


$$N_{I_m} = \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}},$$
$$N_{J_n} = \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}},$$
$$N_{P_k} = \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}.$$

*square:*

$$N_{I_m} \approx N_{J_n} \approx N_{P_k}$$

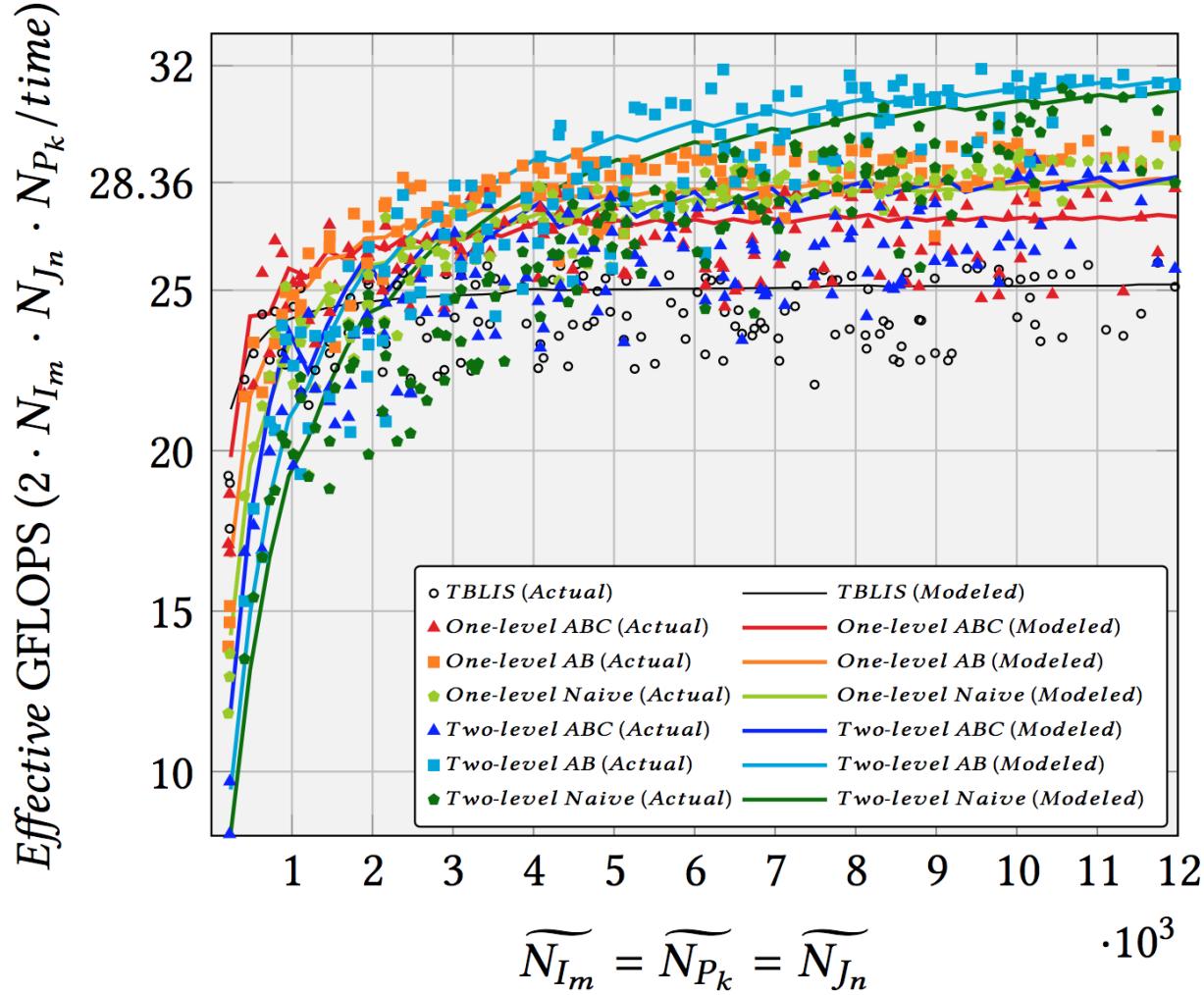
$$\begin{aligned}\widetilde{N_{I_m}} &= \widetilde{N_{J_n}} = \widetilde{N_{P_k}} \\ &= (N_{I_m} \cdot N_{J_n} \cdot N_{P_k})^{1/3}.\end{aligned}$$



# Synthetic Tensor Contractions

Intel Xeon E5-2680 v2 (Ivy Bridge, 10 core/socket)

$$\widetilde{N_{I_m}} = \widetilde{N_{P_k}} = \widetilde{N_{J_n}}, \text{ 1 core, Actual vs. Modeled}$$



$$N_{I_m} = \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}},$$

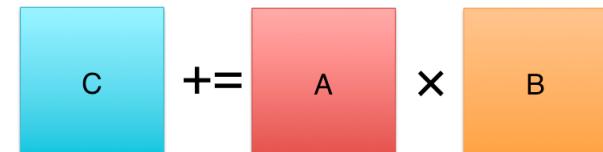
$$N_{J_n} = \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}},$$

$$N_{P_k} = \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}.$$

square:

$$N_{I_m} \approx N_{J_n} \approx N_{P_k}$$

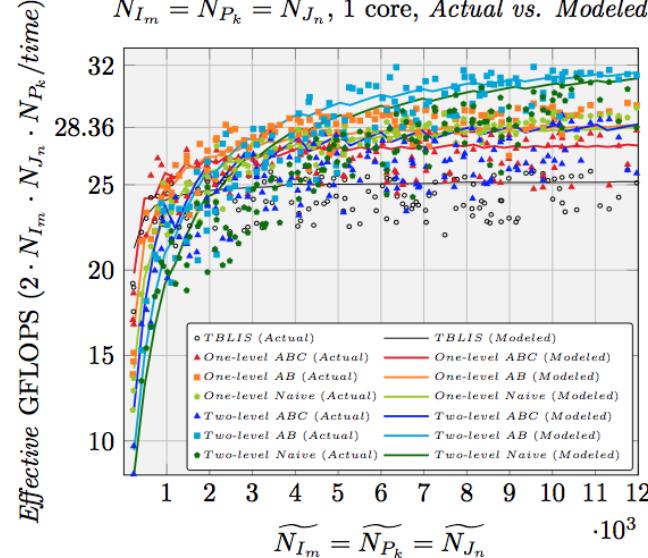
$$\begin{aligned} \widetilde{N_{I_m}} &= \widetilde{N_{J_n}} = \widetilde{N_{P_k}} \\ &= (N_{I_m} \cdot N_{J_n} \cdot N_{P_k})^{1/3}. \end{aligned}$$



# Synthetic Tensor Contractions

Intel Xeon E5-2680 v2 (Ivy Bridge, 10 core/socket)

$$\widetilde{N_{I_m}} = \widetilde{N_{P_k}} = \widetilde{N_{J_n}}, \text{ 1 core, Actual vs. Modeled}$$



square:

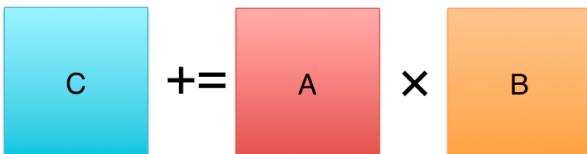
$$N_{I_m} \approx N_{J_n} \approx N_{P_k};$$

$$\widetilde{N_{I_m}} = \widetilde{N_{J_n}} = \widetilde{N_{P_k}} = (N_{I_m} \cdot N_{J_n} \cdot N_{P_k})^{1/3}.$$

$$N_{I_m} = \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}},$$

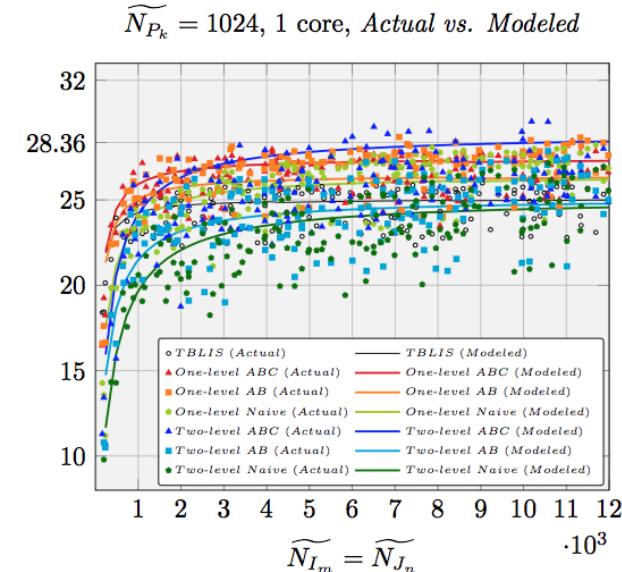
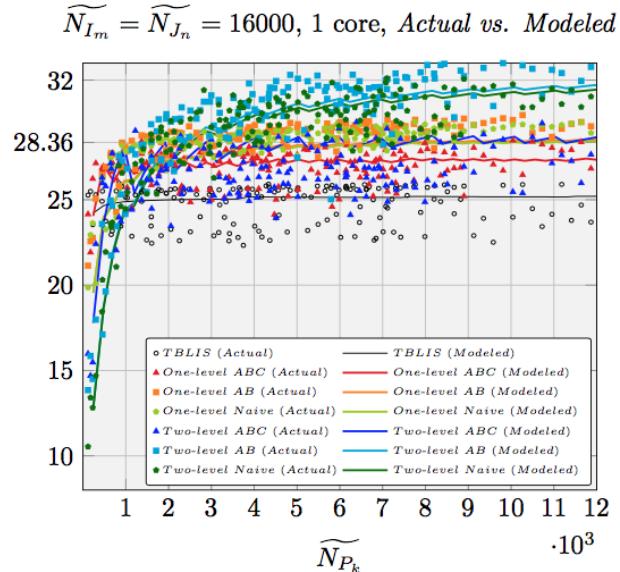
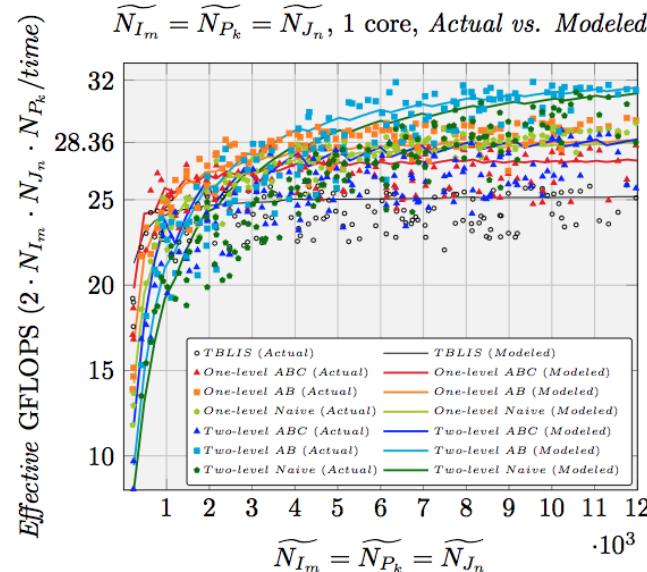
$$N_{J_n} = \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}},$$

$$N_{P_k} = \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}.$$



# Synthetic Tensor Contractions

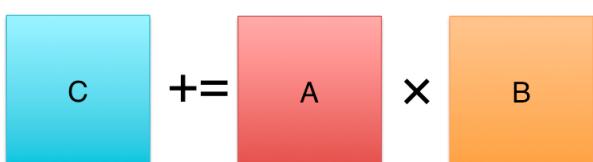
Intel Xeon E5-2680 v2 (Ivy Bridge, 10 core/socket)



square:

$$\widetilde{N_{I_m}} \approx \widetilde{N_{J_n}} \approx \widetilde{N_{P_k}}; \\ \widetilde{N_{I_m}} = \widetilde{N_{J_n}} = \widetilde{N_{P_k}} = (N_{I_m} \cdot N_{J_n} \cdot N_{P_k})^{1/3}.$$

$$N_{I_m} = \prod_{i \in I_m} N_i = N_{i_0} \cdot \dots \cdot N_{i_{m-1}}, \\ N_{J_n} = \prod_{j \in J_n} N_j = N_{j_0} \cdot \dots \cdot N_{j_{n-1}}, \\ N_{P_k} = \prod_{p \in P_k} N_p = N_{p_0} \cdot \dots \cdot N_{p_{k-1}}.$$

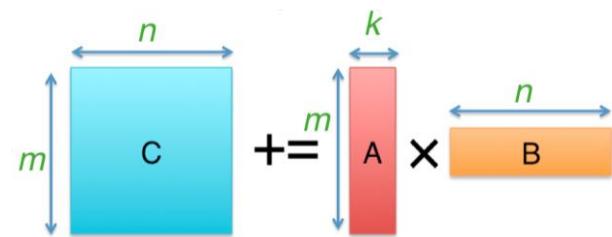


rank- $N_{P_k}$ :

$$\widetilde{N_{I_m}} \approx \widetilde{N_{J_n}} \approx 16000, N_{P_k} \text{ varies}; \\ \widetilde{N_{P_k}} = N_{I_m} \cdot N_{J_n} \cdot N_{P_k} / (16000 \cdot 16000).$$

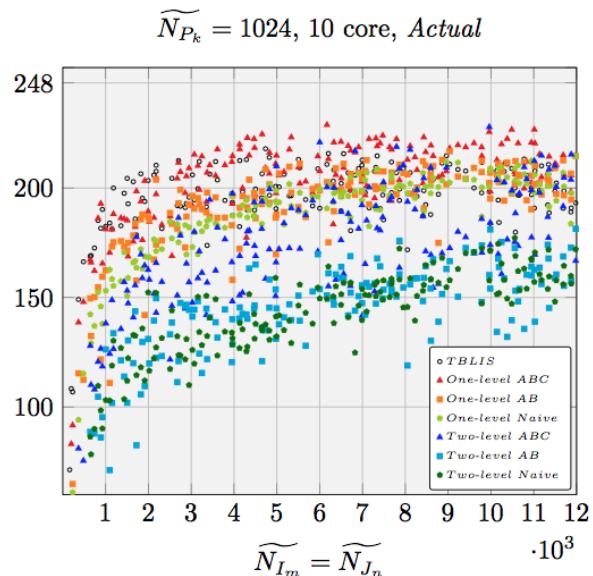
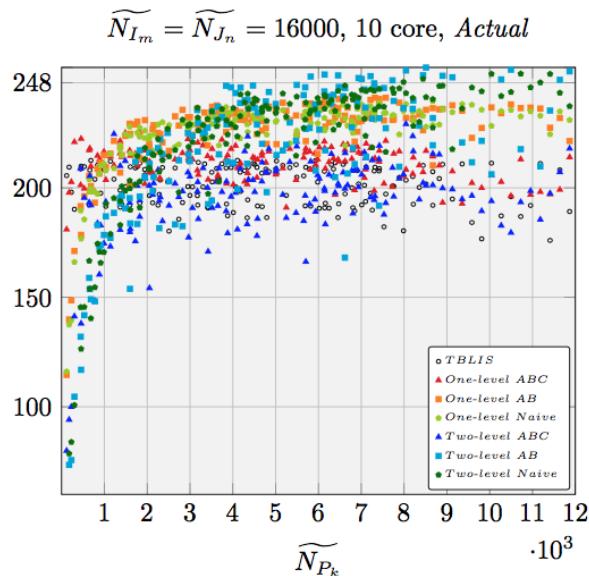
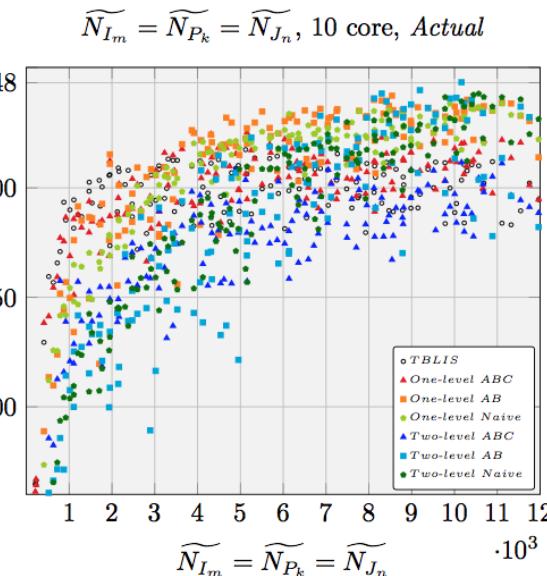
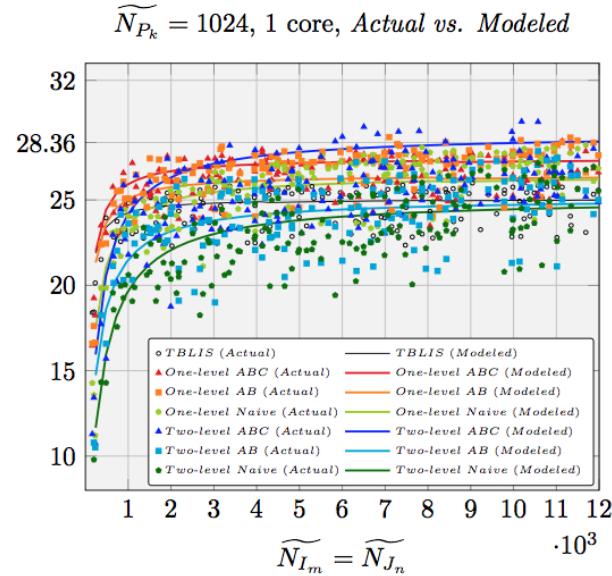
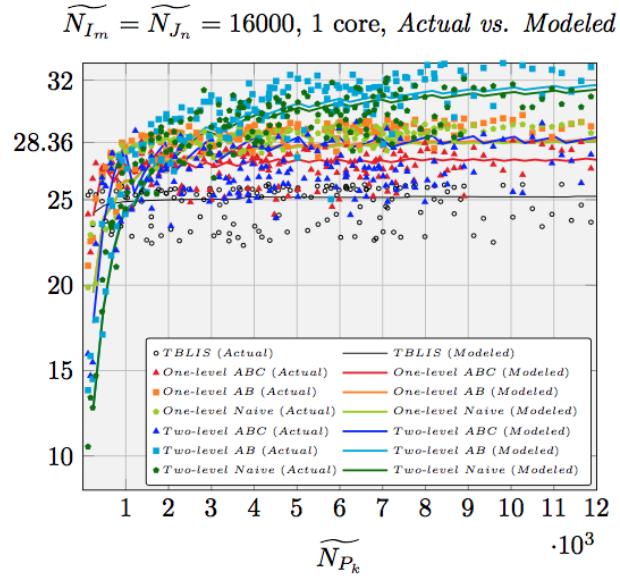
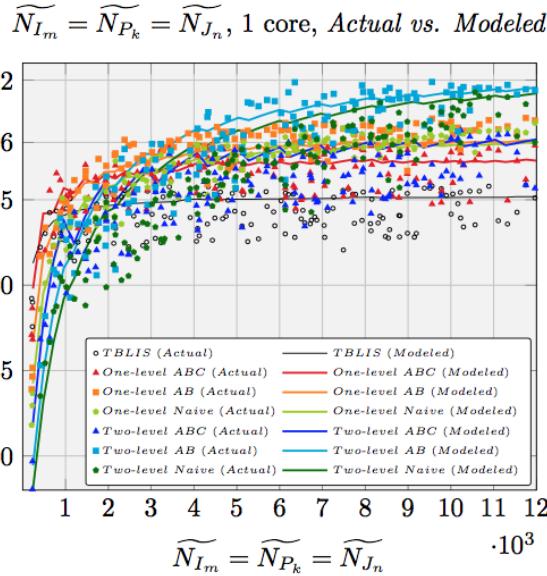
fixed- $N_{P_k}$ :

$$\widetilde{N_{P_k}} \approx 1024, N_{I_m} \approx N_{J_n} \text{ vary.} \\ \widetilde{N_{I_m}} = \widetilde{N_{J_n}} = (N_{I_m} \cdot N_{J_n} \cdot N_{P_k} / 1024)^{1/2}.$$



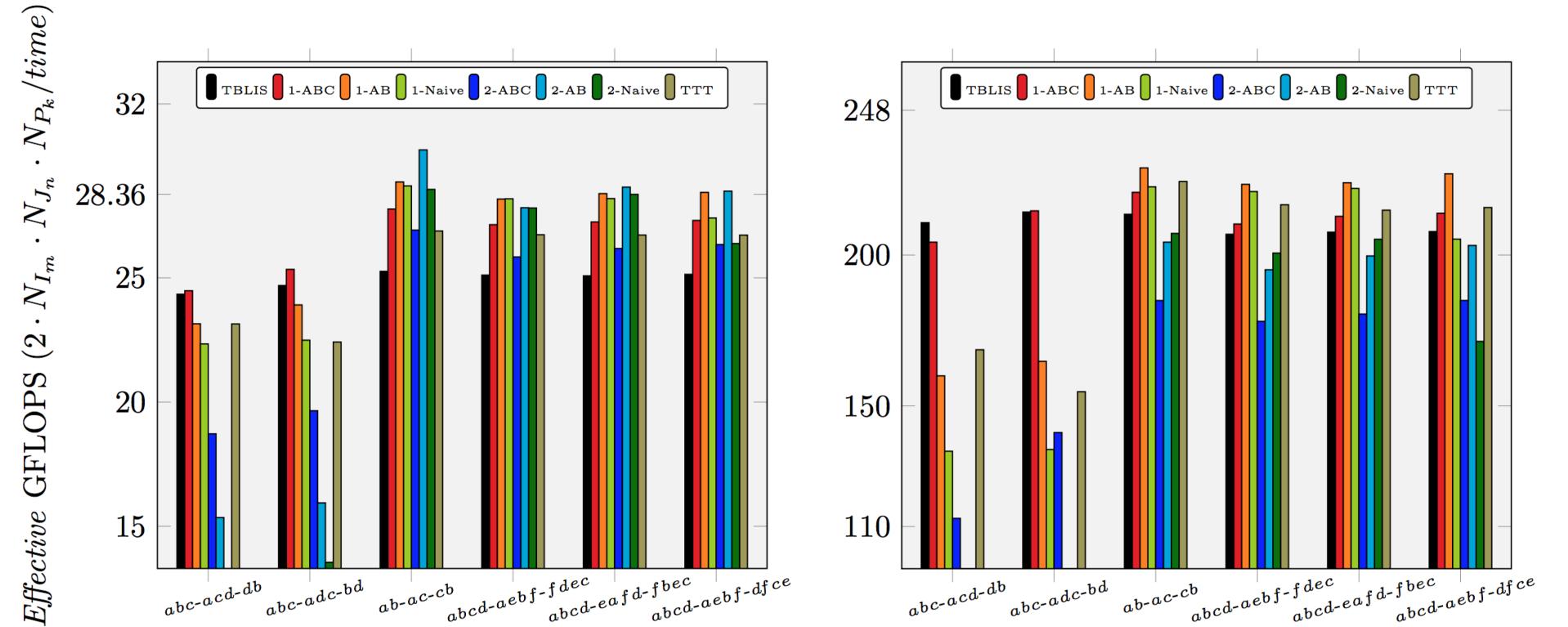
# Synthetic Tensor Contractions

Intel Xeon E5-2680 v2 (Ivy Bridge, 10 core/socket)



# Real-world Benchmark

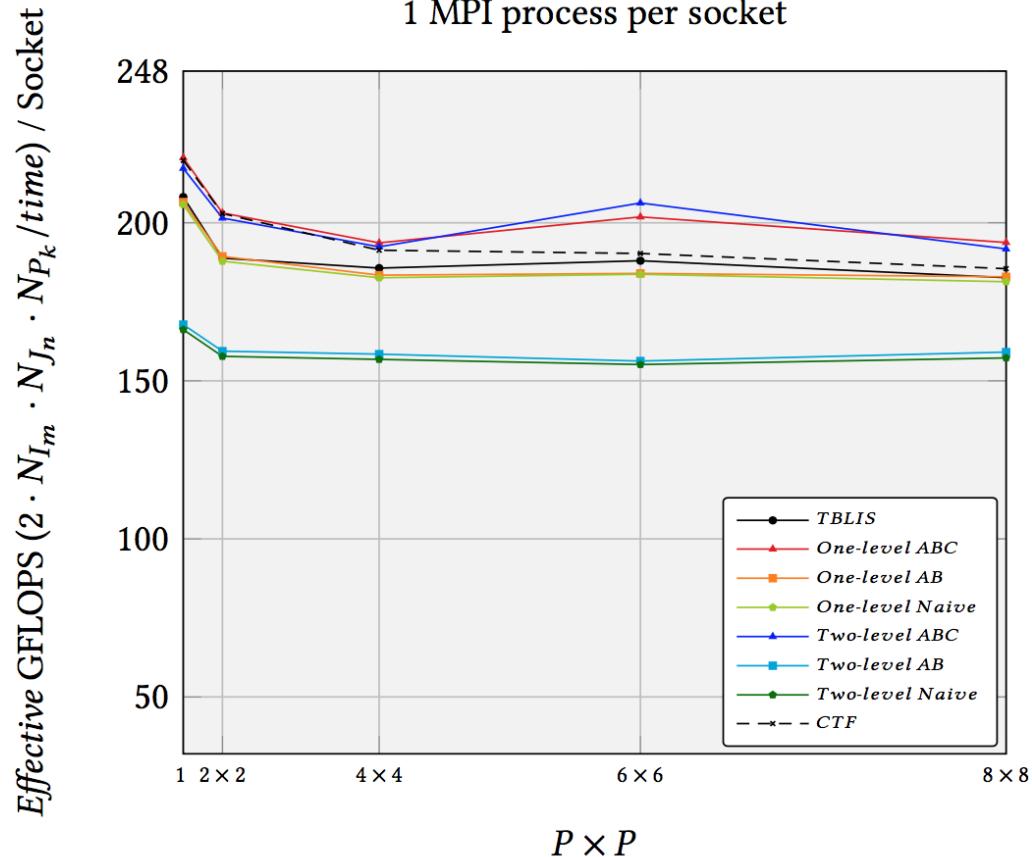
Intel Xeon E5-2680 v2 (Ivy Bridge, 10 core/socket)



$\mathcal{C}_{abcd} = \mathcal{A}_{aebf} \mathcal{B}_{dfce}$  is denoted as  $abcd-aebf-dfce$

# Distributed Memory Experiments

$N_{I_m}(N_b \cdot N_j) = N_{P_k}(N_e \cdot N_m) = N_{J_n}(N_a \cdot N_i) \approx 16000 \cdot P$   
 on  $P \times P$  MPI mesh  
 1 MPI process per socket



Weak scalability performance result of the various implementations for a 4-D tensor contraction CCSD application on distributed memory:

$$\mathcal{Z}_{abij} += \mathcal{W}_{bmej} \mathcal{T}_{aeim}$$

$$\mathcal{Z} += (\mathcal{W}_{e;0} \mid \cdots \mid \mathcal{W}_{e;K-1}) \begin{pmatrix} \frac{\mathcal{T}_{e;0}}{\vdots} \\ \vdots \\ \frac{\mathcal{T}_{e;K-1}}{\vdots} \end{pmatrix}$$

CTF shows the performance of the Cyclops Tensor Framework (linked with Intel MKL).

# Outline

- Background
  - High-performance GEMM
  - High-performance Strassen
  - High-performance Tensor Contraction
- Strassen's Algorithm for Tensor Contraction
- Performance Model
- Experiments
- Conclusion

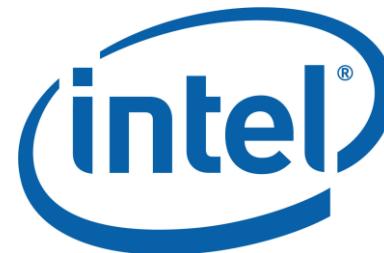


# Conclusion

- First work to leverage Strassen's Algorithm for Tensor Contraction.
- Fusing matrix summation and permutations with packing and micro-kernel operations inside GEMM.
- Avoiding explicit transpositions and extra workspace, and reducing the overhead of memory movement.
- Achieving ~1.3x speedup on single core, multicore, and distributed memory parallel architectures.



# Acknowledgement



- NSF grants ACI-1550493, CCF-1714091.
- A gift from Qualcomm Foundation.
- Intel Corporation through an Intel Parallel Computing Center (IPCC).
- Access to the Maverick and Stampede supercomputers administered by TACC.

We thank Martin Schatz for his help with distributed memory implementations, and the rest of the SHPC team (<http://shpc.ices.utexas.edu>) for their supports.

*Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.*

# Thank you!

The source code can be downloaded from:

<https://github.com/flame/tblis-strassen>