

# A set of building blocks for tensor operations: transposition, summation, and contraction

Paul Springer, Paolo Bientinesi  
*Markus Höhnerbach*

SIAM PP 2018



1979: BLAS Level 1

Vector-Vector operations (e.g., AXPY:  $\mathbf{y} \leftarrow \alpha \mathbf{x} + \mathbf{y}$ )

1979: BLAS Level 1

Vector-Vector operations (e.g., AXPY:  $\mathbf{y} \leftarrow \alpha \mathbf{x} + \mathbf{y}$ )

1988: BLAS Level 2

Matrix-Vector operations (e.g., GEMV:  $\mathbf{y} \leftarrow \alpha A \mathbf{x} + \beta \mathbf{y}$ )

1979: BLAS Level 1

Vector-Vector operations (e.g., AXPY:  $\mathbf{y} \leftarrow \alpha \mathbf{x} + \mathbf{y}$ )

1988: BLAS Level 2

Matrix-Vector operations (e.g., GEMV:  $\mathbf{y} \leftarrow \alpha \mathbf{A}\mathbf{x} + \beta \mathbf{y}$ )

1990: BLAS Level 3

Matrix-Matrix operations (e.g., GEMM:  $\mathbf{C} \leftarrow \alpha \mathbf{A}\mathbf{B} + \beta \mathbf{C}$ )

1979: BLAS Level 1

Vector-Vector operations (e.g., AXPY:  $\mathbf{y} \leftarrow \alpha \mathbf{x} + \mathbf{y}$ )

1988: BLAS Level 2

Matrix-Vector operations (e.g., GEMV:  $\mathbf{y} \leftarrow \alpha \mathbf{A}\mathbf{x} + \beta \mathbf{y}$ )

1990: BLAS Level 3

Matrix-Matrix operations (e.g., GEMM:  $\mathbf{C} \leftarrow \alpha \mathbf{A}\mathbf{B} + \beta \mathbf{C}$ )

2018: BLAS Level 4 ?

Tensor operations

## Goals

- Readily usable open-source software (library and/or compiler)
- Substantial speedups over state-of-the-art algorithms

## Goals

- Readily usable open-source software (library and/or compiler)
- Substantial speedups over state-of-the-art algorithms

## Challenges for High-Performance

## Goals

- Readily usable open-source software (library and/or compiler)
- Substantial speedups over state-of-the-art algorithms

## Challenges for High-Performance

- High order of tensors
  - Strided memory accesses
  - Difficult to exploit *spatial-* and *temporal-locality*



## Goals

- Readily usable open-source software (library and/or compiler)
- Substantial speedups over state-of-the-art algorithms

## Challenges for High-Performance

- High order of tensors
  - Strided memory accesses
  - Difficult to exploit *spatial-* and *temporal-locality*
- Vectorization

## Goals

- Readily usable open-source software (library and/or compiler)
- Substantial speedups over state-of-the-art algorithms

## Challenges for High-Performance

- High order of tensors
  - Strided memory accesses
  - Difficult to exploit *spatial-* and *temporal-locality*
- Vectorization
- Parallelization

1 Tensor Transpositions

2 Spin Summations

3 Tensor Contractions

- Transpositions of the general form:

$$\mathcal{B}_{i_1, i_2, \dots, i_N} \leftarrow \alpha \mathcal{A}_{\pi(i_1, i_2, \dots, i_N)} + \beta \mathcal{B}_{i_1, i_2, \dots, i_N}$$

---

<sup>1</sup>P. Springer, T. Su and P. Bientinesi, "HPTT: A High-Performance Tensor Transposition C++ Library", ARRAY 2017

- Transpositions of the general form:

$$\mathcal{B}_{i_1, i_2, \dots, i_N} \leftarrow \alpha \mathcal{A}_{\pi(i_1, i_2, \dots, i_N)} + \beta \mathcal{B}_{i_1, i_2, \dots, i_N}$$

- HPTT: C++ 11 library for tensor transpositions<sup>1</sup>
  - Explicitly vectorized
  - Multi-threaded
  - Autotuning

---

<sup>1</sup>P. Springer, T. Su and P. Bientinesi, "HPTT: A High-Performance Tensor Transposition C++ Library", ARRAY 2017

- Transpositions of the general form:

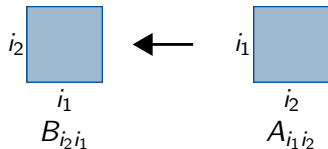
$$\mathcal{B}_{i_1, i_2, \dots, i_N} \leftarrow \alpha \mathcal{A}_{\pi(i_1, i_2, \dots, i_N)} + \beta \mathcal{B}_{i_1, i_2, \dots, i_N}$$

- HPTT: C++ 11 library for tensor transpositions<sup>1</sup>
  - Explicitly vectorized
  - Multi-threaded
  - Autotuning
- Dynamic data structure called `plan` (similar to FFTW)
  - Encodes the execution of a tensor transposition
  - Loop Order
  - Parallelism

---

<sup>1</sup>P. Springer, T. Su and P. Bientinesi, "HPTT: A High-Performance Tensor Transposition C++ Library", ARRAY 2017

```
for  $i_1 = 0 : N$   
  for  $i_2 = 0 : N$   
     $B[i_2, i_1] \leftarrow A[i_1, i_2]$ 
```



```

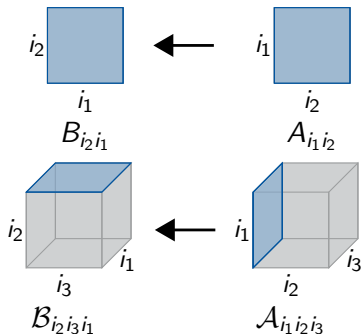
for  $i_1 = 0 : N$ 
  for  $i_2 = 0 : N$ 
     $B[i_2, i_1] \leftarrow A[i_1, i_2]$ 

```

```

for  $i_2 = 0 : N$ 
  for  $i_3 = 0 : N$ 
    for  $i_1 = 0 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

```





```

for  $i_1 = 0 : N$ 
  for  $i_2 = 0 : N$ 
     $B[i_2, i_1] \leftarrow A[i_1, i_2]$ 

```

```

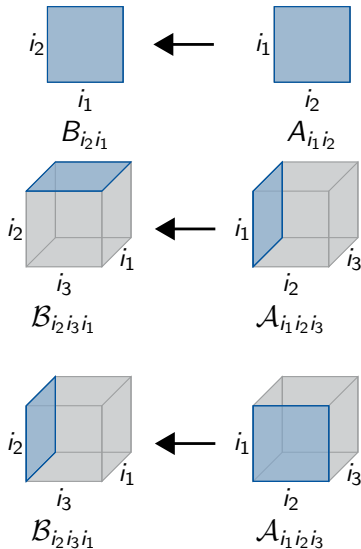
for  $i_2 = 0 : N$ 
  for  $i_3 = 0 : N$ 
    for  $i_1 = 0 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

```

```

for  $i_3 = 0 : N$ 
  for  $i_2 = 0 : N$ 
    for  $i_1 = 0 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

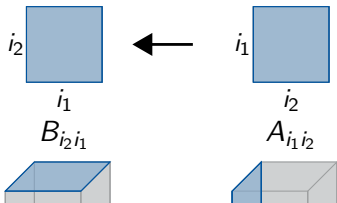
```



```

for  $i_1 = 0 : N$ 
  for  $i_2 = 0 : N$ 
     $B[i_2, i_1] \leftarrow A[i_1, i_2]$ 

```



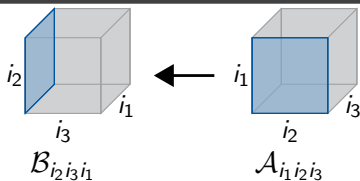
## Key Insight

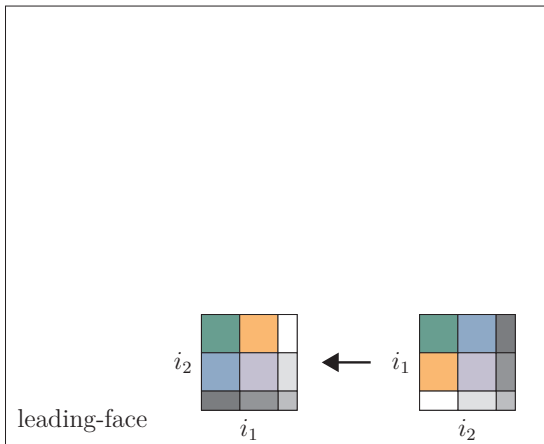
- *Leading-face*: Spanned by the two stride-1 modes
- Reduction to 2D is always possible

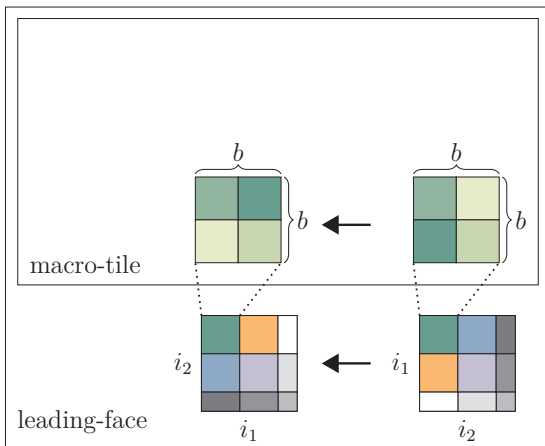
```

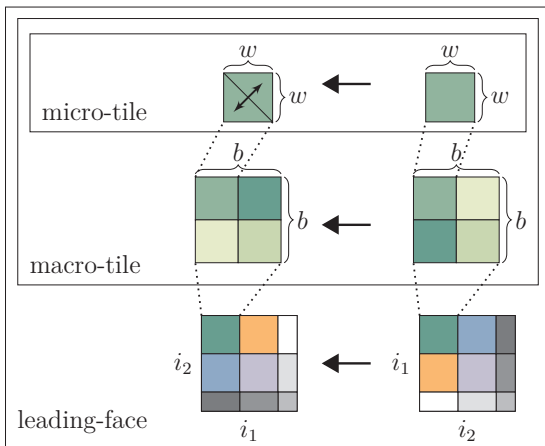
for  $i_3 = 0 : N$ 
  for  $i_2 = 0 : N$ 
    for  $i_1 = 0 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

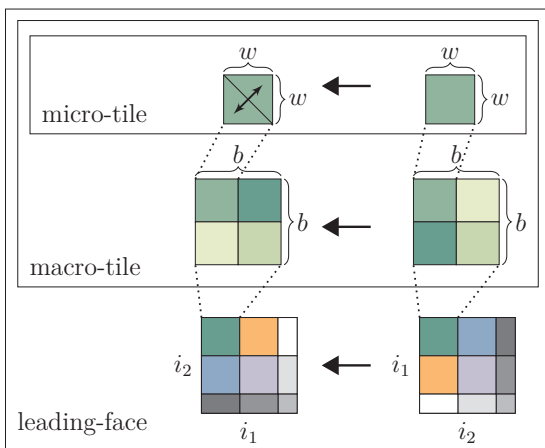
```



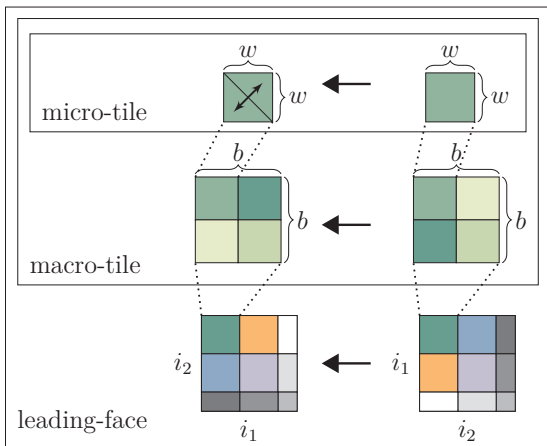




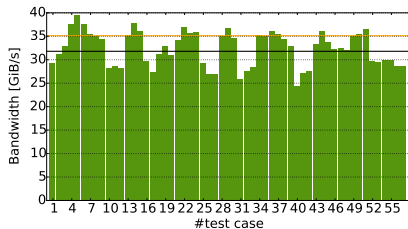




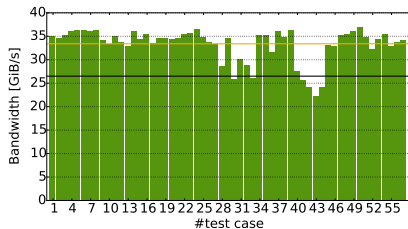
- Decompose a macro-tile into micro-tiles
  - Vectorized micro-tiles



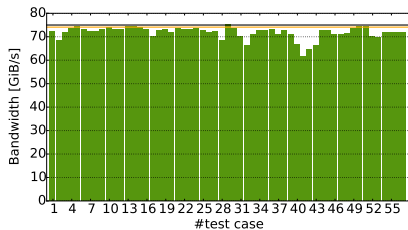
- Decompose a macro-tile into micro-tiles
  - Vectorized micro-tiles
- Decompose a transposition into macro-tiles
  - Parallel over macro-tiles



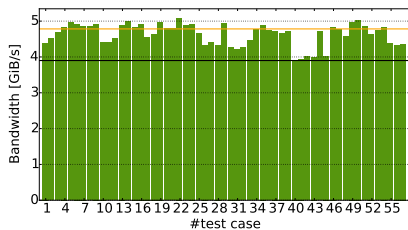
(a) Intel Ivy Bridge E5-2670 v2.



(b) IBM Power7.



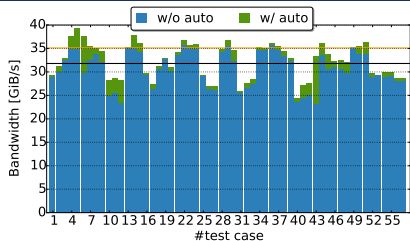
(c) Intel KNL Xeon Phi 7210.



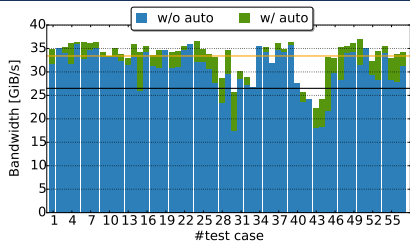
(d) ARMv7-A.

HPTT's bandwidth. Black and orange lines respectively denote STREAM and AXPY bandwidth.

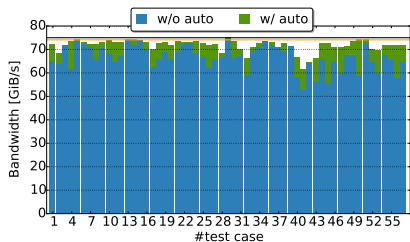




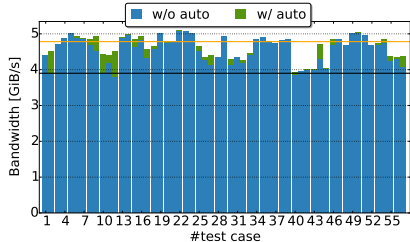
(a) Intel Ivy Bridge E5-2670 v2.



(b) IBM Power7.



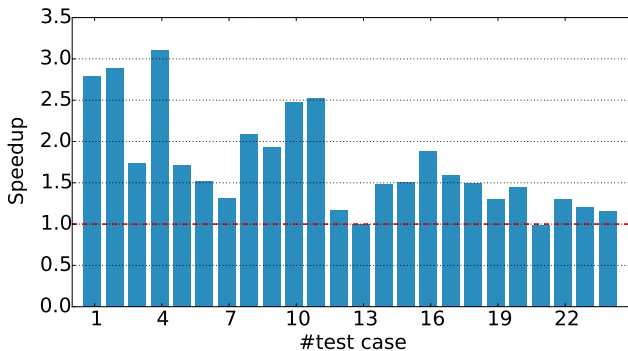
(c) Intel KNL Xeon Phi 7210.



(d) ARMv7-A.

HPTT's bandwidth. Black and orange lines respectively denote STREAM and AXPY bandwidth.

- Cyclops Tensor Framework  
Tensor Contractions  $\rightarrow$  Transpositions



HPTT's impact on CTF's performance.

1 Tensor Transpositions

2 Spin Summations

3 Tensor Contractions

## Linear summation over tensor transpositions

- $\mathcal{B}_{i_0 i_1 i_2} = 2\mathcal{A}_{i_0 i_1 i_2} - \mathcal{A}_{i_2 i_1 i_0} - \mathcal{A}_{i_0 i_2 i_1}$
- $\mathcal{B}_{i_0 i_1 i_2} = 4\mathcal{A}_{i_0 i_1 i_2} - 2\mathcal{A}_{i_1 i_0 i_2} - 2\mathcal{A}_{i_2 i_1 i_0} + \mathcal{A}_{i_1 i_2 i_0} - 2\mathcal{A}_{i_0 i_2 i_1} + \mathcal{A}_{i_2 i_0 i_1}$
- $\mathcal{B}_{i_0 i_1 i_2 i_3} = 2\mathcal{A}_{i_0 i_1 i_2 i_3} - \mathcal{A}_{i_2 i_1 i_0 i_3} - \mathcal{A}_{i_0 i_2 i_1 i_3} - \mathcal{A}_{i_0 i_1 i_3 i_2}$
- $\mathcal{B}_{i_0 i_1 i_2 i_3} = 4\mathcal{A}_{i_0 i_1 i_2 i_3} - 2\mathcal{A}_{i_0 i_3 i_2 i_1} - 2\mathcal{A}_{i_0 i_1 i_3 i_2} - 2\mathcal{A}_{i_1 i_0 i_2 i_3} + \mathcal{A}_{i_1 i_3 i_2 i_0} + \mathcal{A}_{i_1 i_0 i_3 i_2} - 2\mathcal{A}_{i_2 i_1 i_0 i_3} + \mathcal{A}_{i_2 i_3 i_0 i_1} + \mathcal{A}_{i_2 i_1 i_3 i_0} - 2\mathcal{A}_{i_3 i_1 i_2 i_0} + \mathcal{A}_{i_3 i_0 i_2 i_1} + \mathcal{A}_{i_3 i_1 i_0 i_2}$

## Linear summation over tensor transpositions

- $\mathcal{B}_{i_0 i_1 i_2} = 2\mathcal{A}_{i_0 i_1 i_2} - \mathcal{A}_{i_2 i_1 i_0} - \mathcal{A}_{i_0 i_2 i_1}$
- $\mathcal{B}_{i_0 i_1 i_2} = 4\mathcal{A}_{i_0 i_1 i_2} - 2\mathcal{A}_{i_1 i_0 i_2} - 2\mathcal{A}_{i_2 i_1 i_0} + \mathcal{A}_{i_1 i_2 i_0} - 2\mathcal{A}_{i_0 i_2 i_1} + \mathcal{A}_{i_2 i_0 i_1}$
- $\mathcal{B}_{i_0 i_1 i_2 i_3} = 2\mathcal{A}_{i_0 i_1 i_2 i_3} - \mathcal{A}_{i_2 i_1 i_0 i_3} - \mathcal{A}_{i_0 i_2 i_1 i_3} - \mathcal{A}_{i_0 i_1 i_3 i_2}$
- $\mathcal{B}_{i_0 i_1 i_2 i_3} = 4\mathcal{A}_{i_0 i_1 i_2 i_3} - 2\mathcal{A}_{i_0 i_3 i_2 i_1} - 2\mathcal{A}_{i_0 i_1 i_3 i_2} - 2\mathcal{A}_{i_1 i_0 i_2 i_3} + \mathcal{A}_{i_1 i_3 i_2 i_0} + \mathcal{A}_{i_1 i_0 i_3 i_2} - 2\mathcal{A}_{i_2 i_1 i_0 i_3} + \mathcal{A}_{i_2 i_3 i_0 i_1} + \mathcal{A}_{i_2 i_1 i_3 i_0} - 2\mathcal{A}_{i_3 i_1 i_2 i_0} + \mathcal{A}_{i_3 i_0 i_2 i_1} + \mathcal{A}_{i_3 i_1 i_0 i_2}$

### Central Challenge

Spatial and temporal locality in both  $\mathcal{A}$  and  $\mathcal{B}$

## Linear summation over tensor transpositions

- $\mathcal{B}_{i_0 i_1 i_2} = 2\mathcal{A}_{i_0 i_1 i_2} - \mathcal{A}_{i_2 i_1 i_0} - \mathcal{A}_{i_0 i_2 i_1}$
- $\mathcal{B}_{i_0 i_1 i_2} = 4\mathcal{A}_{i_0 i_1 i_2} - 2\mathcal{A}_{i_1 i_0 i_2} - 2\mathcal{A}_{i_2 i_1 i_0} + \mathcal{A}_{i_1 i_2 i_0} - 2\mathcal{A}_{i_0 i_2 i_1} + \mathcal{A}_{i_2 i_0 i_1}$
- $\mathcal{B}_{i_0 i_1 i_2 i_3} = 2\mathcal{A}_{i_0 i_1 i_2 i_3} - \mathcal{A}_{i_2 i_1 i_0 i_3} - \mathcal{A}_{i_0 i_2 i_1 i_3} - \mathcal{A}_{i_0 i_1 i_3 i_2}$
- $\mathcal{B}_{i_0 i_1 i_2 i_3} = 4\mathcal{A}_{i_0 i_1 i_2 i_3} - 2\mathcal{A}_{i_0 i_3 i_2 i_1} - 2\mathcal{A}_{i_0 i_1 i_3 i_2} - 2\mathcal{A}_{i_1 i_0 i_2 i_3} + \mathcal{A}_{i_1 i_3 i_2 i_0} + \mathcal{A}_{i_1 i_0 i_3 i_2} - 2\mathcal{A}_{i_2 i_1 i_0 i_3} + \mathcal{A}_{i_2 i_3 i_0 i_1} + \mathcal{A}_{i_2 i_1 i_3 i_0} - 2\mathcal{A}_{i_3 i_1 i_2 i_0} + \mathcal{A}_{i_3 i_0 i_2 i_1} + \mathcal{A}_{i_3 i_1 i_0 i_2}$

### Central Challenge

Spatial and temporal locality in both  $\mathcal{A}$  and  $\mathcal{B}$

P. Springer, D. Matthews and P. Bientinesi,

“Spin Summations: A High-Performance Perspective”, ACM TOMS’18

1 Tensor Transpositions

2 Spin Summations

3 Tensor Contractions

- Prior Work
  - Nested loops
  - Loops over GEMM (LoG)
  - Transpose-Transpose-GEMM-Transpose (TTGT)

---

<sup>2</sup>P. Springer and P. Bientinesi, "Design of a high-performance GEMM-like Tensor-Tensor Multiplication", ACM TOMS'18

<sup>3</sup>Field Van Zee et al. "BLIS: A Framework for Rapidly Instantiating BLAS Functionality", ACM TOMS'15



- Prior Work
  - Nested loops
  - Loops over GEMM (LoG)
  - Transpose-Transpose-GEMM-Transpose (TTGT)
- Disadvantages
  - Suboptimal I/O cost
  - Poor cache utilization
  - Additional memory requirements

---

<sup>2</sup>P. Springer and P. Bientinesi, "Design of a high-performance GEMM-like Tensor-Tensor Multiplication", ACM TOMS'18

<sup>3</sup>Field Van Zee et al. "BLIS: A Framework for Rapidly Instantiating BLAS Functionality", ACM TOMS'15

- Prior Work
  - Nested loops
  - Loops over GEMM (LoG)
  - Transpose-Transpose-GEMM-Transpose (TTGT)
- Disadvantages
  - Suboptimal I/O cost
  - Poor cache utilization
  - Additional memory requirements
- We propose a fourth approach: GETT<sup>2</sup>
  - Akin to a high-performance GEMM implementation
  - Adopts the BLIS<sup>3</sup> methodology

---

<sup>2</sup>P. Springer and P. Bientinesi, "Design of a high-performance GEMM-like Tensor-Tensor Multiplication", ACM TOMS'18

<sup>3</sup>Field Van Zee et al. "BLIS: A Framework for Rapidly Instantiating BLAS Functionality", ACM TOMS'15

## Matrix-Matrix Multiplication

$$C_{m,n} \leftarrow \sum_k A_{m,k} B_{k,n}$$

Matrix-Matrix Multiplication (Einstein notation)

$$C_{m,n} \leftarrow A_{m,k} B_{k,n}$$

## Matrix-Matrix Multiplication (Einstein notation)

$$C_{m,n} \leftarrow A_{m,k} B_{k,n}$$

```
// N-Loop
for n = 0 : N - 1
  // M-Loop
  for m = 0 : M - 1
    tmp = 0
    // K-Loop (contracted)
    for k = 0 : K - 1
      tmp += Ai,k Bk,j
    // update C
    Ci,j = α tmp + β Ci,j
```

Naive GEMM

## Matrix-Matrix Multiplication (Einstein notation)

$$C_{m,n} \leftarrow A_{m,k} B_{k,n}$$

```
// N-Loop
for n = 0 : N - 1
  // M-Loop
  for m = 0 : M - 1
    tmp = 0
    // K-Loop (contracted)
    for k = 0 : K - 1
      tmp += Ai,k Bk,j
    // update C
    Ci,j = α tmp + β Ci,j
```

Naive GEMM

```
// N-Loop
for n = 0 : nc : N - 1
  // K-Loop (contracted)
  for k = 0 : kc : K - 1
    B̂ = identify_submatrix(B, n, k)
    // pack B̂ into B̃
    B̃ = packB(B̂) // B̃ ∈ ℝkc × nc
    // M-Loop
    for m = 0 : mc : M - 1
      Â = identify_submatrix(A, m, k)
      // pack Â into Ã
      Ã = packA(Â) // Ã ∈ ℝmc × kc
      Ĉ = identify_submatrix(C, m, n)
      // matrix-matrix product: ÃB̃
      macroKernel(Ã, B̃, Ĉ, α, β)
```

High-performance GEMM

## Key Idea

Pack-and-transpose while moving data into the caches

## Key Idea

Pack-and-transpose while moving data into the caches

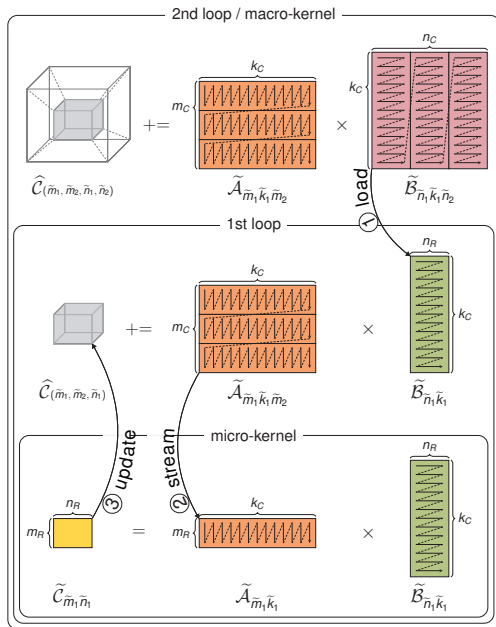
```

1 // N-Loop
2 for n = 1 : nc : Sln
3   // K-Loop (contracted)
4   for k = 1 : kc : Slk
5      $\hat{B}$  = identify_subtensor(B, n, k)
6     // pack  $\hat{B}$  into  $\tilde{B}$  (L3 cache)
7      $\tilde{B}$  = packB( $\hat{B}$ )
8     // M-Loop
9     for m = 1 : mc : Slm
10       $\hat{A}$  = identify_subtensor(A, m, k)
11      // pack  $\hat{A}$  into  $\tilde{A}$  (L2 cache)
12       $\tilde{A}$  = packA( $\hat{A}$ )
13       $\hat{C}$  = identify_subtensor(C, m, n)
14      // compute matrix-matrix product of  $\tilde{A}\tilde{B}$ 
15      macroKernel( $\tilde{A}$ ,  $\tilde{B}$ ,  $\hat{C}$ ,  $\alpha$ ,  $\beta$ )

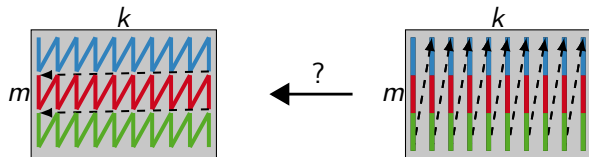
```

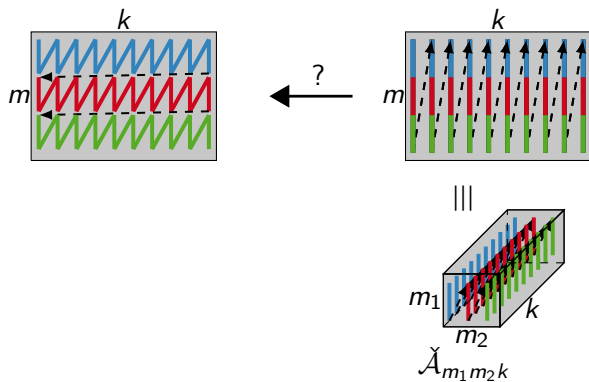
High-performance GETT

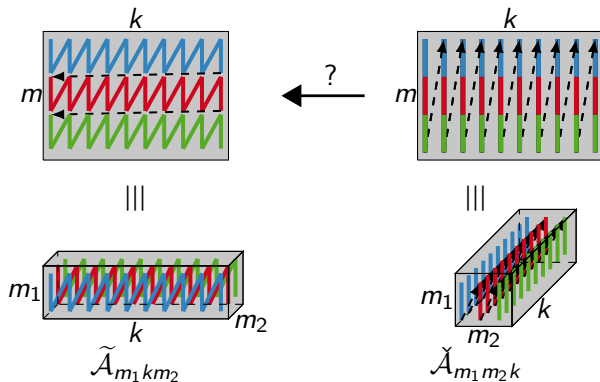


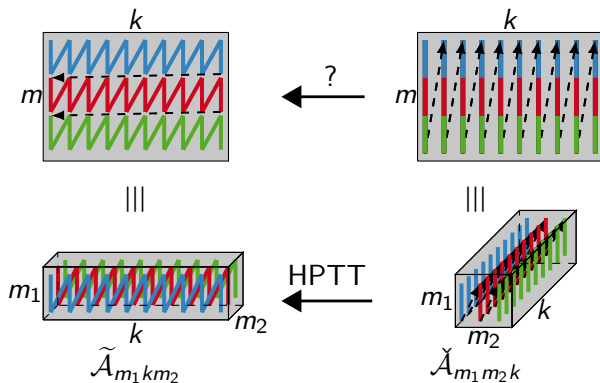


- Similar to BLIS
- Blocking
  - L3 cache
  - L2 cache
  - L1 cache
  - Registers
- Vectorized *micro-kernel*









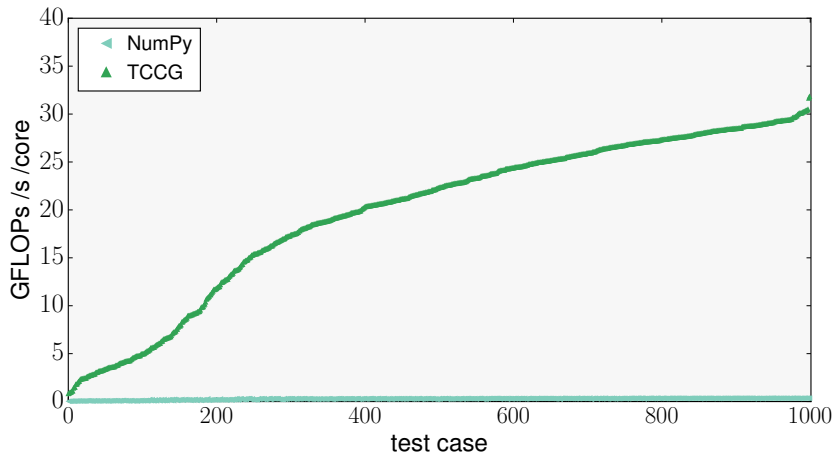
- Multi-dimensional packing routines
  - Preserve stride-1 index
  - Exploit spatial locality  $\Rightarrow$  high efficiency

- Tensor Contraction Code Generator (TCCG)<sup>4</sup>
  - Research Tool
  - Supports: Naive, LoG, TTGT, and **GETT**
  - Autotuning
  
- Tensor Contraction Library (TCL)<sup>5</sup>
  - C++ library
  - Based on TTGT + HPTT
  - Offers C and Python interfaces
  - Supports a large variety of TCs

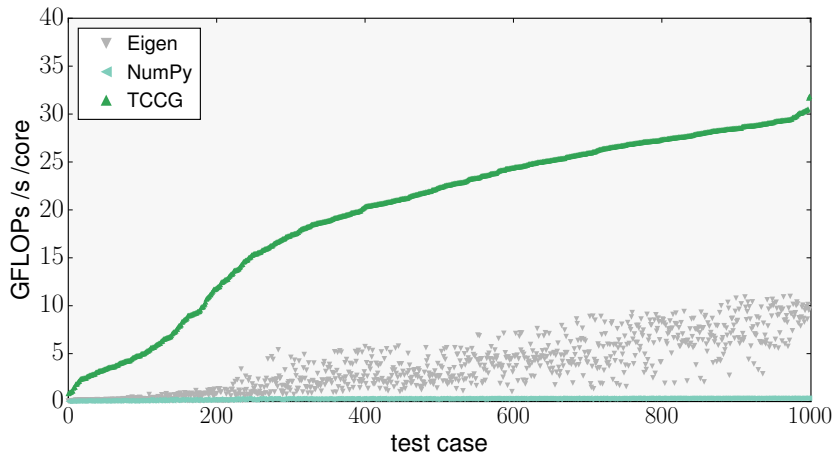
---

<sup>4</sup>Available at: <https://github.com/hpac/tccg>

<sup>5</sup>Available at: <https://github.com/springer13/tcl>

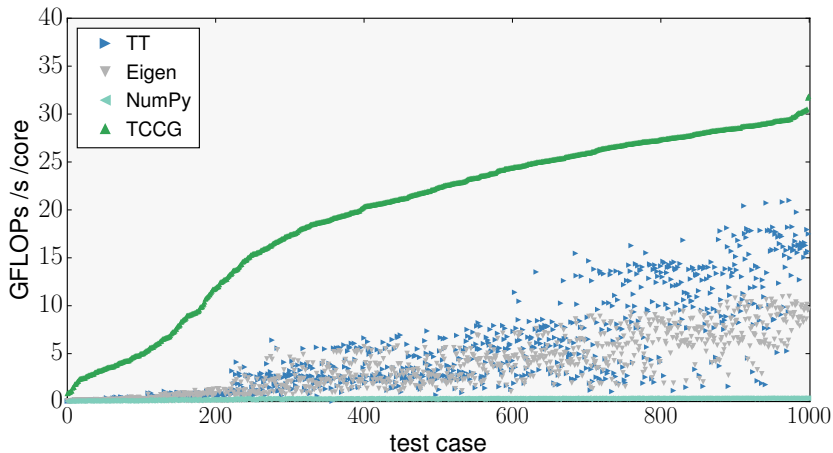


DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.

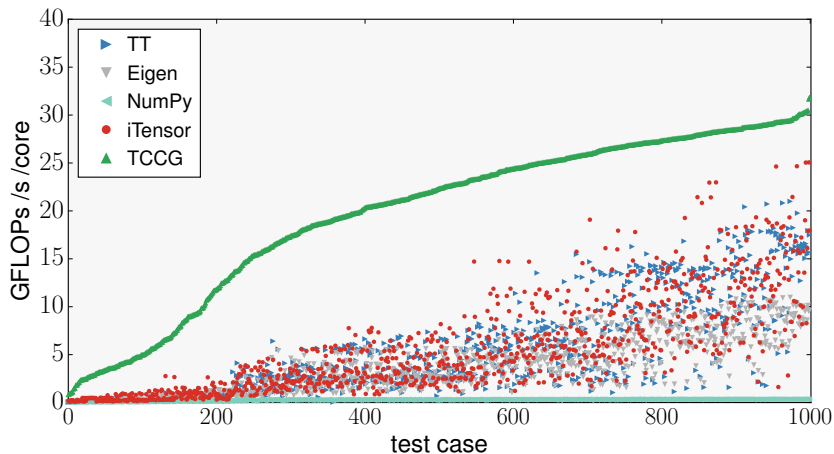


DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.

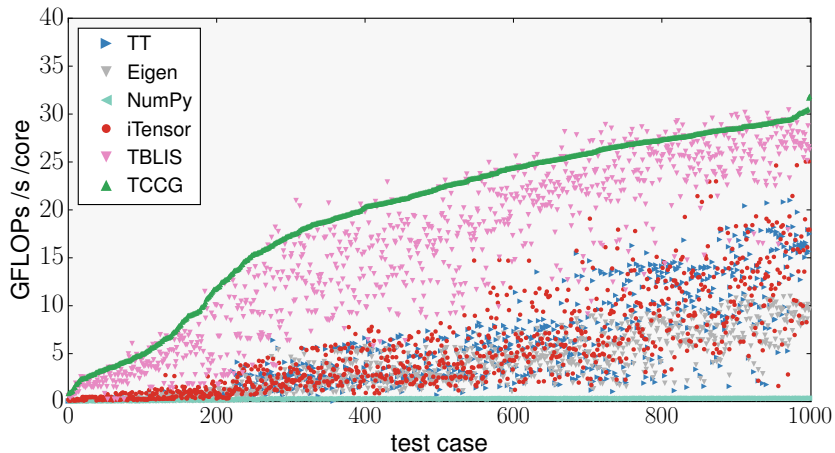




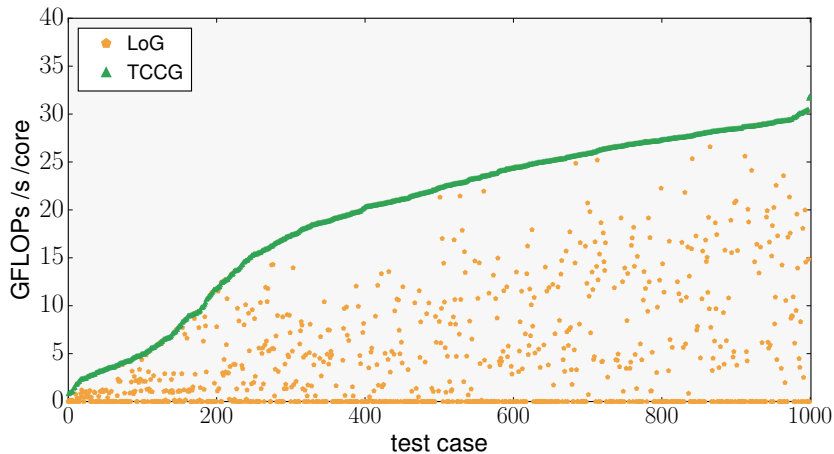
DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.



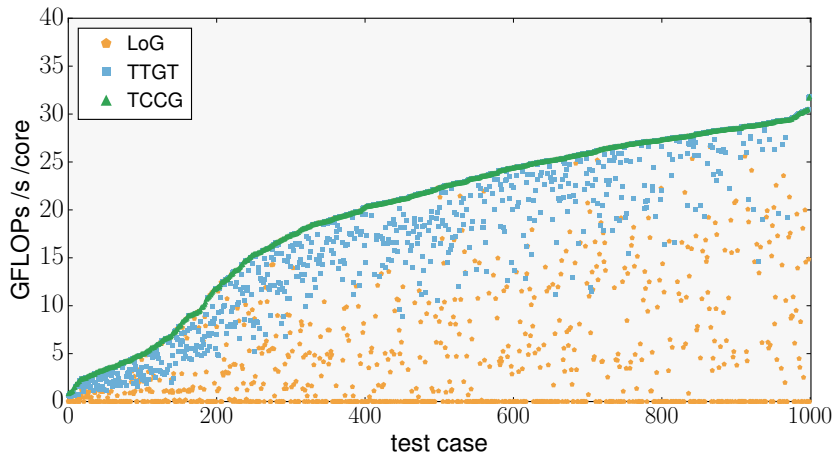
DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.



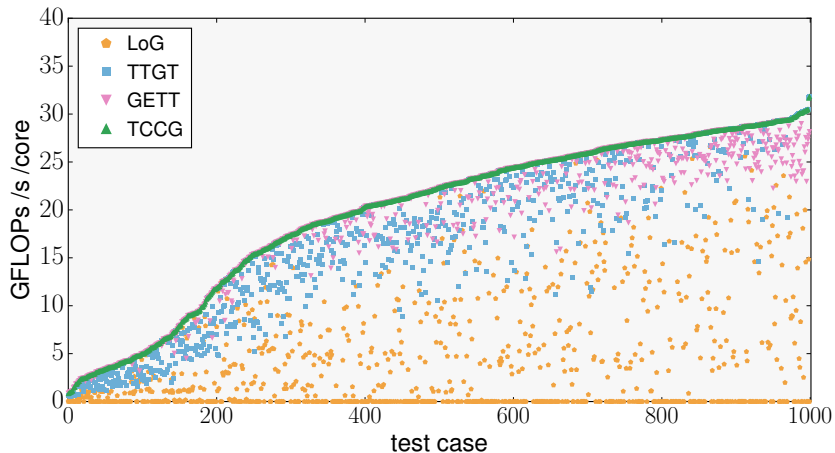
DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.



DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.



DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.



DP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.

- Noticeable Speedups over state-of-the-art solutions ✓
  - Proper memory accesses are critical for tensor operations
- Released open-source code ✓
  - High-Performance Tensor Transpose C++ Library (HPTT)<sup>6</sup>
  - Code Generator for Spin Summations<sup>7</sup>
  - GEMM-like Tensor-Tensor contraction (GETT)<sup>8</sup>
  - Tensor Contraction C++ Library (TCL)<sup>9</sup>

---

<sup>6</sup>Available at <https://github.com/springer13/hptt>, BSD-3

<sup>7</sup>Available at <https://github.com/springer13/spin>, BSD-3

<sup>8</sup>Available at <https://github.com/hpac/tccg>, LGPLv3

<sup>9</sup>Available at <https://github.com/springer13/tcl>, LGPLv3

- Noticeable Speedups over state-of-the-art solutions ✓
  - Proper memory accesses are critical for tensor operations
- Released open-source code ✓
  - High-Performance Tensor Transpose C++ Library (HPTT)<sup>6</sup>
  - Code Generator for Spin Summations<sup>7</sup>
  - GEMM-like Tensor-Tensor contraction (GETT)<sup>8</sup>
  - Tensor Contraction C++ Library (TCL)<sup>9</sup>

Thank you for your attention

---

<sup>6</sup> Available at [hptts://github.com/springer13/hptt](https://github.com/springer13/hptt), BSD-3

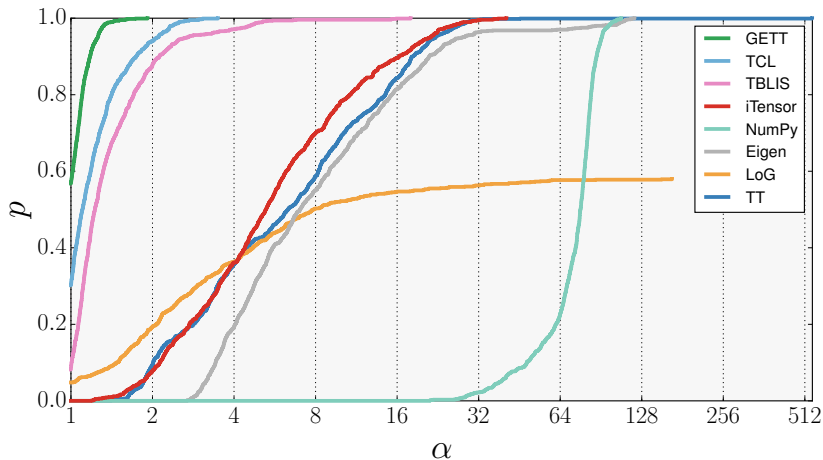
<sup>7</sup> Available at [hptts://github.com/springer13/spin](https://github.com/springer13/spin), BSD-3

<sup>8</sup> Available at [hptts://github.com/hpac/tccg](https://github.com/hpac/tccg), LGPLv3

<sup>9</sup> Available at [hptts://github.com/springer13/tcl](https://github.com/springer13/tcl), LGPLv3

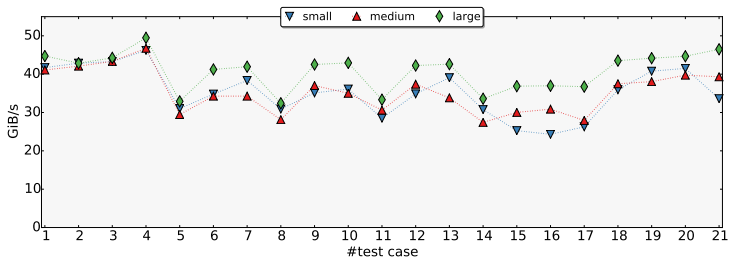


- [1] Paul Springer and Paolo Bientinesi. Design of a high-performance gemm-like tensor-tensor multiplication. *ACM Trans. Math. Softw.*, 44(3):28:1–28:29, January 2018.
- [2] Paul Springer, Jeff R. Hammond, and Paolo Bientinesi. TTC: A high-performance compiler for tensor transpositions. *ACM Trans. Math. Softw.*, 44(2):15:1–15:21, August 2017.
- [3] Paul Springer, Devin Matthews, and Paolo Bientinesi. Spin Summations: A High-Performance Perspective. apr 2017. Under review for ACM Transactions on Mathematical Software.
- [4] Paul Springer, Aravind Sankaran, and Paolo Bientinesi. Ttc: A tensor transposition compiler for multiple architectures. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Libraries, Languages, and Compilers for Array Programming*, ARRAY 2016, pages 41–46, New York, NY, USA, 2016. ACM.
- [5] Paul Springer, Tong Su, and Paolo Bientinesi. HPTT: A High-Performance Tensor Transposition C++ Library. In *Proceedings of the 4th ACM SIGPLAN International Workshop on Libraries, Languages, and Compilers for Array Programming*, ARRAY 2017, pages 56–62, New York, NY, USA, 2017. ACM.

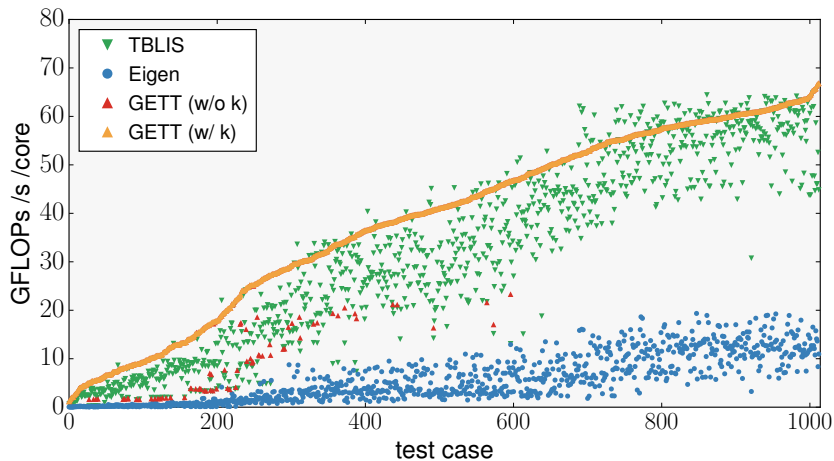


Performance Profile.





Bandwidth; small, medium and large respectively correspond tensors of size 70 MiB, 320 MiB, and 1200 MiB.



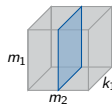
SP tensor contractions. Host: 2× Haswell-EP E5-2680 v3 @ 24 threads.

## Key Idea

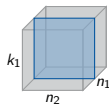
Identify 2D subtensors and contract them via GEMM

$$\bullet C_{m_1, n_1, n_2, m_2} \leftarrow A_{m_1, m_2, k_1} B_{k_1, n_2, n_1}$$

```
for( m2 = 0; m2 < M2; m2++ )  
  for( n1 = 0; n1 < N1; n1++ )  
    gemm( M1, N2, K1, A[:, m2, :], B[:, :, n1], C[:, n1, :, m2] )
```



$A_{m_1, m_2, k_1}$



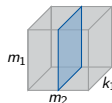
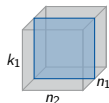
$B_{k_1, n_2, n_1}$

## Key Idea

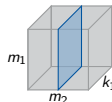
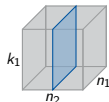
Identify 2D subtensors and contract them via GEMM

$$\bullet C_{m_1, n_1, n_2, m_2} \leftarrow A_{m_1, m_2, k_1} B_{k_1, n_2, n_1}$$

```
for( m2 = 0; m2 < M2; m2++ )
  for( n1 = 0; n1 < N1; n1++ )
    gemm( M1, N2, K1, A[:, m2, :], B[:, :, n1], C[:, n1, :, m2] )
```


 $A_{m_1, m_2, k_1}$ 

 $B_{k_1, n_2, n_1}$ 

```
for( m2 = 0; m2 < M2; m2++ )
  for( n2 = 0; n2 < N2; n2++ )
    gemm( M1, N1, K1, A[:, m2, :], B[:, n2, :], C[:, :, n2, m2] )
```


 $A_{m_1, m_2, k_1}$ 

 $B_{k_1, n_2, n_1}$

## Key Idea

- 1 “Flatten” the tensors to matrices
- 2 Use GEMM for contraction
- 3 “Unflatten” output matrix to tensor

- $C_{m_1, n_1, n_2, m_2} \leftarrow A_{m_1, m_2, k_1} B_{k_1, n_2, n_1}$



## Key Idea

- 1 “Flatten” the tensors to matrices
- 2 Use GEMM for contraction
- 3 “Unflatten” output matrix to tensor

- $\mathcal{C}_{m_1, n_1, n_2, m_2} \leftarrow \mathcal{A}_{m_1, m_2, k_1} \mathcal{B}_{k_1, n_2, n_1}$

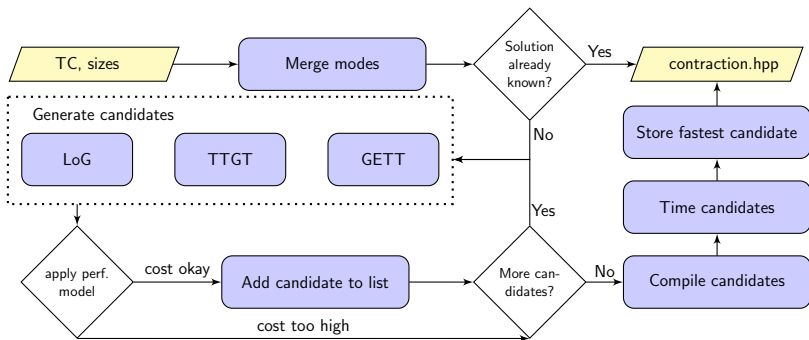
$$\begin{aligned} \tilde{\mathcal{C}}_{(m_1, m_2), (n_2, n_1)} &\leftarrow \mathcal{A}_{(m_1, m_2), k_1} \times \mathcal{B}_{k_1, (n_2, n_1)} \\ \mathcal{C}_{m_1, n_1, n_2, m_2} &\leftarrow \tilde{\mathcal{C}}_{m_1, m_2, n_2, n_1} \end{aligned}$$

- **Input:** Mathematical description of TC
  - e.g.,  $C[a,b,i,j] = A[i,k,a] * B[k,j,b];$
- **Output:** High-Performance C++ code

---

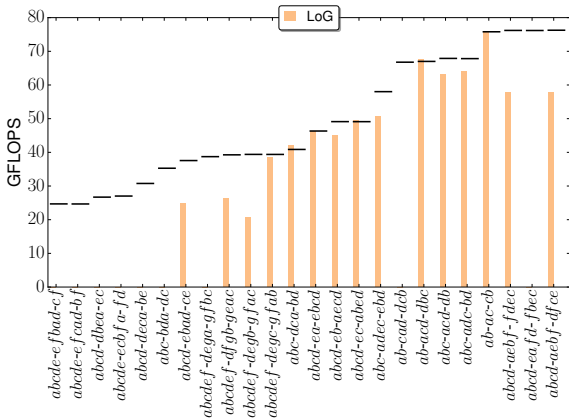
<sup>10</sup>Available at: <https://github.com/hpac/tccg>

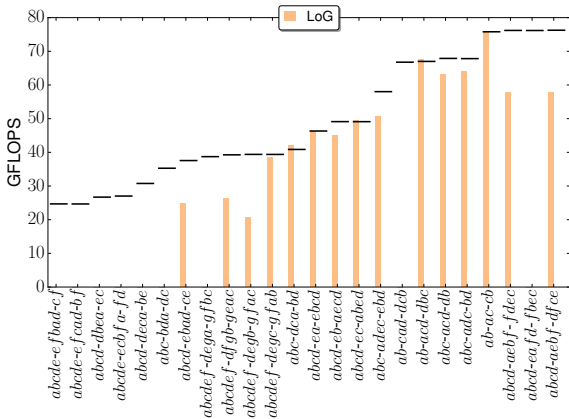
- **Input:** Mathematical description of TC
  - e.g.,  $C[a,b,i,j] = A[i,k,a] * B[k,j,b]$ ;
- **Output:** High-Performance C++ code



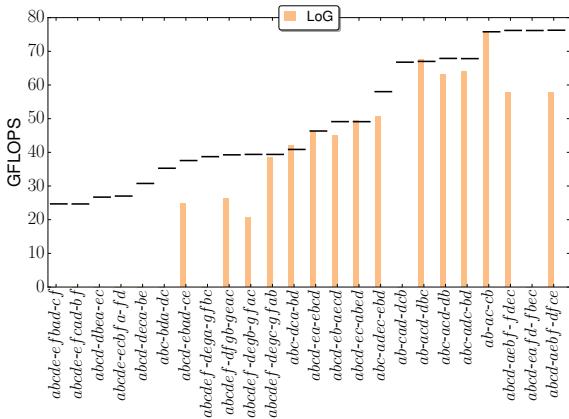
Schematic overview of TCCG<sup>10</sup>.

<sup>10</sup>Available at: <https://github.com/hpac/tccg>

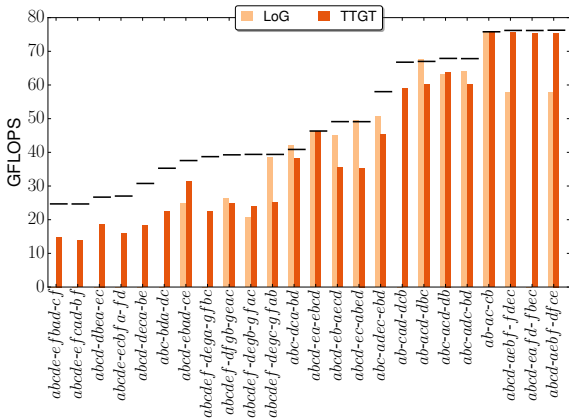


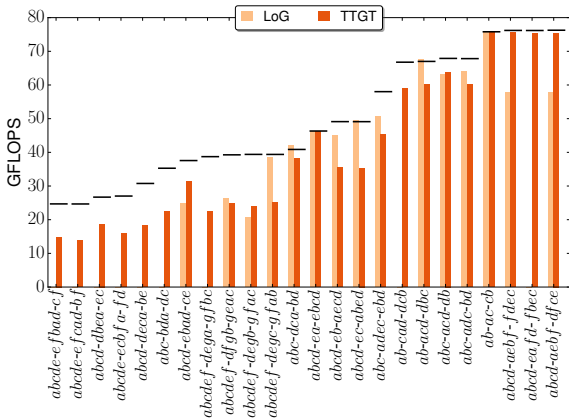


- Not all TCs can be implemented via LoG



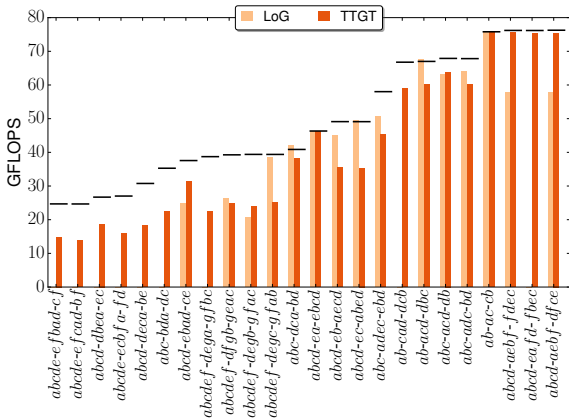
- Not all TCs can be implemented via LoG
- Mixed performance



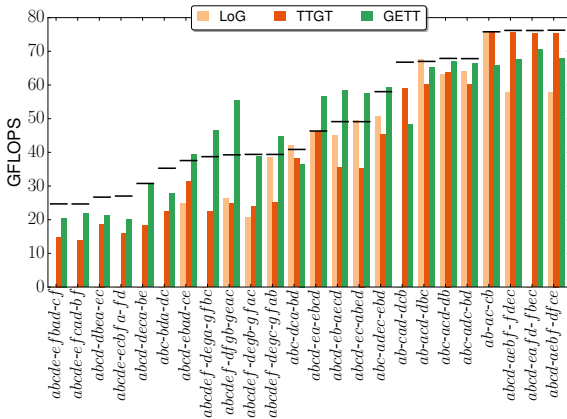


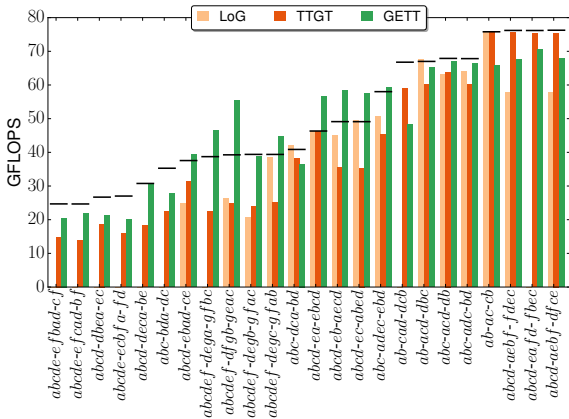
- TTGT: good for compute-bound TCs



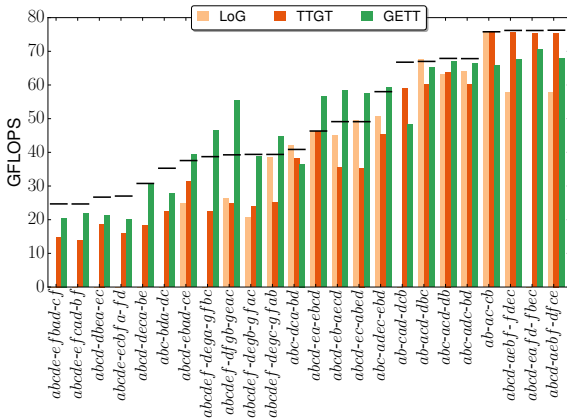


- TTGT: good for compute-bound TCs
- TTGT: bad for bandwidth-bound TCs

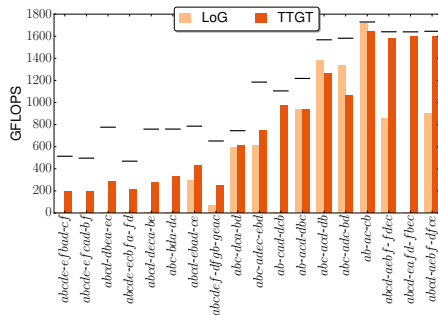




- GETT: excels for bandwidth-bound TCs

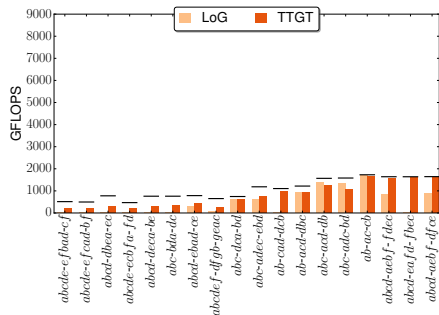


- GETT: excels for bandwidth-bound TCs
- GETT: good for compute-bound TCs

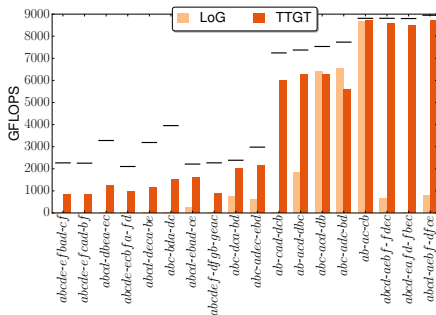


(a) 2x Intel Xeon E5-2680 v3

- Performance gap increases for bandwidth-bound TCs

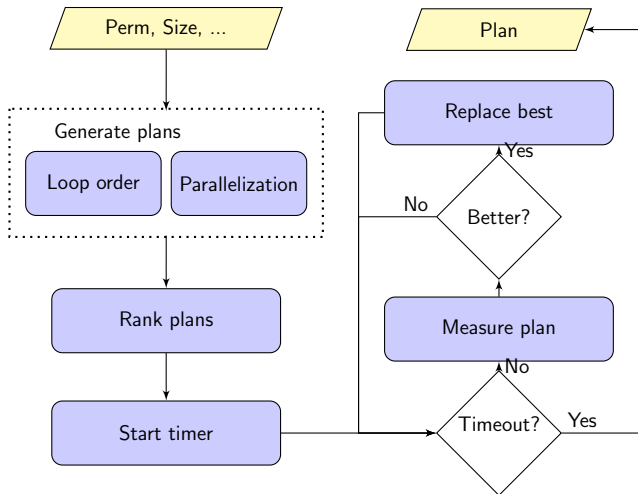


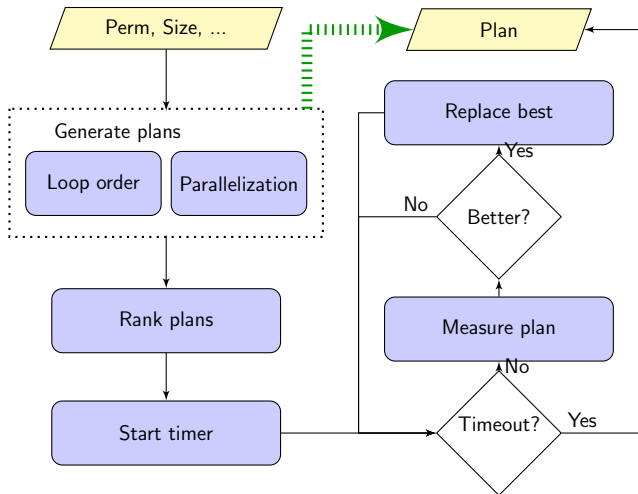
(a) 2x Intel Xeon E5-2680 v3



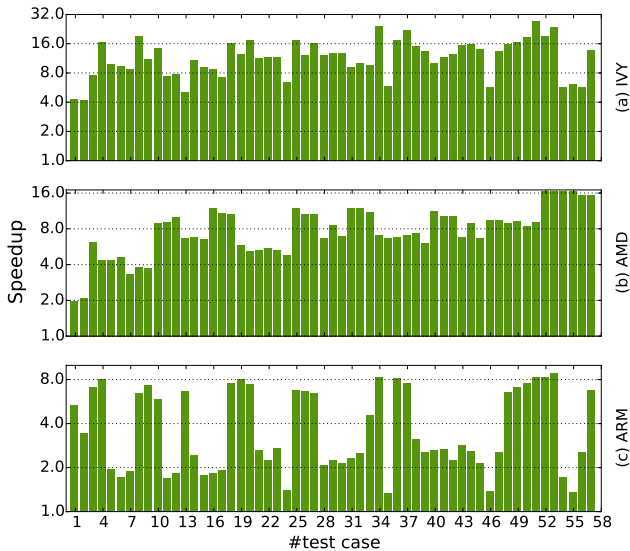
(b) NVIDIA Tesla P100

- Performance gap increases for bandwidth-bound TCs









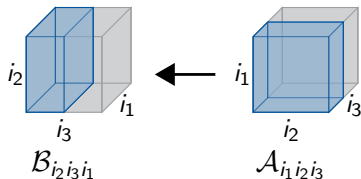
Speedup over Eigen.



```

for  $i_3 = 0 : N/2$ 
  for  $i_1 = 0 : N$ 
    for  $i_2 = 0 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

```

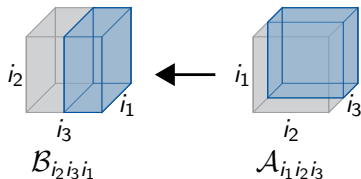


(a) Thread 1

```

for  $i_3 = N/2 : N$ 
  for  $i_1 = 0 : N$ 
    for  $i_2 = 0 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

```

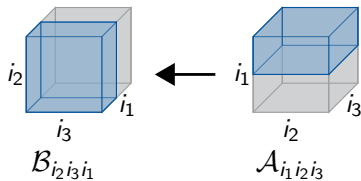


(b) Thread 2

```

for  $i_3 = 0 : N$ 
  for  $i_2 = 0 : N$ 
    for  $i_1 = 0 : N/2$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

```

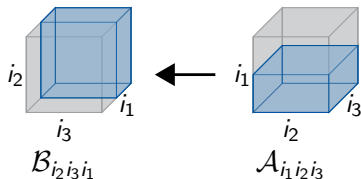


(a) Thread 1

```

for  $i_3 = N$ 
  for  $i_2 = 0 : N$ 
    for  $i_1 = N/2 : N$ 
       $B[i_2, i_3, i_1] \leftarrow A[i_1, i_2, i_3]$ 

```



(b) Thread 2