

Design of a High-Performance GEMM-like Tensor-Tensor Multiplication

Paul Springer and Paolo Bientinesi

Aachen Institute for Advanced Study in
Computational Engineering Science

Austin, Sep. 20th 2016



- 1 Introduction
- 2 GEMM-like Tensor-Tensor Multiplication
- 3 Tensor Contraction Code Generator
- 4 Performance
- 5 Conclusion and Future Work

- Tensors can be thought of as higher dimensional matrices
- Tensor contraction can be thought of as higher dimensional GEMMs

¹Paul Springer and Paolo Bientinesi. “Design of a high-performance GEMM-like Tensor-Tensor Multiplication”. In: *TOMS, in review* ().

- Tensors can be thought of as higher dimensional matrices
- Tensor contraction can be thought of as higher dimensional GEMMs
- Essentially three approaches:
 - Nested loops
 - Transpose-Transpose-GEMM-Transpose (TTGT)
 - Loops over GEMM (LoG)

¹Paul Springer and Paolo Bientinesi. “Design of a high-performance GEMM-like Tensor-Tensor Multiplication”. In: *TOMS, in review* ().

- Tensors can be thought of as higher dimensional matrices
- Tensor contraction can be thought of as higher dimensional GEMMs
- Essentially three approaches:
 - Nested loops
 - Transpose-Transpose-GEMM-Transpose (TTGT)
 - Loops over GEMM (LoG)
- We propose a novel approach: GETT¹
 - Akin to a high-performance GEMM implementation
 - Adopts the BLIS methodology: **Breaking through the BLAS layer**

¹Paul Springer and Paolo Bientinesi. “Design of a high-performance GEMM-like Tensor-Tensor Multiplication”. In: *TOMS, in review* ().

- Tensors can be thought of as higher dimensional matrices
- Tensor contraction can be thought of as higher dimensional GEMMs
- Essentially three approaches:
 - Nested loops
 - Transpose-Transpose-GEMM-Transpose (TTGT)
 - Loops over GEMM (LoG)
- We propose a novel approach: GETT¹
 - Akin to a high-performance GEMM implementation
 - Adopts the BLIS methodology: **Breaking through the BLAS layer**
- Tensor Contraction Code Generator (TCCG)
 - combine GETT, TTGT and LoG into a unified tool

¹Paul Springer and Paolo Bientinesi. “Design of a high-performance GEMM-like Tensor-Tensor Multiplication”. In: *TOMS, in review* ().

Matrix-Matrix Multiplication

$A \in \mathbb{R}^{M \times K}$, $B \in \mathbb{R}^{K \times N}$ and $C \in \mathbb{R}^{M \times N}$ be 2D tensors:

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

Matrix-Matrix Multiplication (Einstein notation)

$A \in \mathbb{R}^{M \times K}$, $B \in \mathbb{R}^{K \times N}$ and $C \in \mathbb{R}^{M \times N}$ be 2D tensors:

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

Matrix-Matrix Multiplication (Einstein notation)

$A \in \mathbb{R}^{M \times K}$, $B \in \mathbb{R}^{K \times N}$ and $C \in \mathbb{R}^{M \times N}$ be 2D tensors:

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

```
// N-Loop
for j = 0 : N - 1
  // M-Loop
  for i = 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)

$A \in \mathbb{R}^{M \times K}$, $B \in \mathbb{R}^{K \times N}$ and $C \in \mathbb{R}^{M \times N}$ be 2D tensors:

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

```
// N-Loop
for j = 0 : N - 1
  // M-Loop
  for i = 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
     $\hat{B}$  = identify_submatrix(B, n, k)
    // pack  $\hat{B}$  into  $\tilde{B}$ 
     $\tilde{B}$  = packB( $\hat{B}$ ) //  $\tilde{B} \in \mathbb{R}^{kc \times nc}$ 
    // M-Loop
    for m = 0 : mc : M - 1
       $\hat{A}$  = identify_submatrix(A, m, k)
      // pack  $\hat{A}$  into  $\tilde{A}$ 
       $\tilde{A}$  = packA( $\hat{A}$ ) //  $\tilde{A} \in \mathbb{R}^{mc \times kc}$ 
       $\hat{C}$  = identify_submatrix(C, m, n)
      // matrix-matrix product:  $\tilde{A}\tilde{B}$ 
      macroKernel( $\tilde{A}$ ,  $\tilde{B}$ ,  $\hat{C}$ , α, β)
```

High-performance GEMM.

- Tensor contraction examples:
 - $C_{m,n} \leftarrow A_{m,k} B_{k,n}$

- Tensor contraction examples:

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

- $C_{m_1,m_2,n} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- Tensor contraction examples:

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

- $C_{m_1,m_2,n} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- $C_{m_1,n,m_2} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- Tensor contraction examples:

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

- $C_{m_1,m_2,n} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- $C_{m_1,n,m_2} \leftarrow A_{m_1,m_2,k} B_{k,n}$

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

- Tensor contraction examples:

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

- $C_{m_1,m_2,n} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- $C_{m_1,n,m_2} \leftarrow A_{m_1,m_2,k} B_{k,n}$

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

- $C_{m_1,n_1,n_2,m_2} \leftarrow A_{m_1,k_1,m_2,k_2} B_{k_2,n_2,k_1,n_1}$

- Tensor contraction examples:

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

- $C_{m_1,m_2,n} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- $C_{m_1,n,m_2} \leftarrow A_{m_1,m_2,k} B_{k,n}$

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

- $C_{m_1,n_1,n_2,m_2} \leftarrow A_{m_1,k_1,m_2,k_2} B_{k_2,n_2,k_1,n_1}$

- $C_{m_1,n_1,n_2,m_2,n_3} \leftarrow A_{m_1,k_1,m_2,k_2} B_{n_3,k_2,n_2,k_1,n_1}$

- ...

- Tensor contraction examples:

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

- $C_{m_1,m_2,n} \leftarrow A_{m_1,m_2,k} B_{k,n}$

- $C_{m_1,n,m_2} \leftarrow A_{m_1,m_2,k} B_{k,n}$

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

- $C_{m_1,n_1,n_2,m_2} \leftarrow A_{m_1,k_1,m_2,k_2} B_{k_2,n_2,k_1,n_1}$

- $C_{m_1,n_1,n_2,m_2,n_3} \leftarrow A_{m_1,k_1,m_2,k_2} B_{n_3,k_2,n_2,k_1,n_1}$

- ...

⇒ Quite similar to GEMM.

Tensor-Tensor Multiplication (Einstein notation)

Let the input tensors $\mathcal{A} \in \mathbb{R}^{s_1^A \times s_2^A \times \dots \times s_{r_A}^A}$, and $\mathcal{B} \in \mathbb{R}^{s_1^B \times s_2^B \times \dots \times s_{r_B}^B}$ update the output tensor $\mathcal{C} \in \mathbb{R}^{s_1^C \times s_2^C \times \dots \times s_{r_C}^C}$:

$$\mathcal{C}_{\Pi^C(I_m \cup I_n)} \leftarrow \alpha \mathcal{A}_{\Pi^A(I_m \cup I_k)} \mathcal{B}_{\Pi^B(I_n \cup I_k)} + \beta \mathcal{C}_{\Pi^C(I_m \cup I_n)}.$$

Tensor-Tensor Multiplication (Einstein notation)

Let the input tensors $\mathcal{A} \in \mathbb{R}^{S_1^A \times S_2^A \times \dots \times S_{r_A}^A}$, and $\mathcal{B} \in \mathbb{R}^{S_1^B \times S_2^B \times \dots \times S_{r_B}^B}$ update the output tensor $\mathcal{C} \in \mathbb{R}^{S_1^C \times S_2^C \times \dots \times S_{r_C}^C}$:

$$\mathcal{C}_{\Pi^C(I_m \cup I_n)} \leftarrow \alpha \mathcal{A}_{\Pi^A(I_m \cup I_k)} \mathcal{B}_{\Pi^B(I_n \cup I_k)} + \beta \mathcal{C}_{\Pi^C(I_m \cup I_n)}.$$

- These index sets I_m , I_n and I_k are critical
 - $I_m := \{m_1, m_2, \dots, m_\gamma\}$: free indices of \mathcal{A}
 - $I_n := \{n_1, n_2, \dots, n_\zeta\}$: free indices of \mathcal{B}
 - $I_k := \{k_1, k_2, \dots, k_\xi\}$: contracted indices

```

1 // N-Loops
2 for  $n_1 = 1 : S_{n_1}$ 
3 // ... remaining N-loops omitted ...
4 for  $n_\zeta = 1 : S_{n_\zeta}$ 
5 // M-Loops
6 for  $m_1 = 1 : S_{m_1}$ 
7 // ... remaining M-loops omitted ...
8 for  $m_\gamma = 1 : S_{m_\gamma}$ 
9     tmp = 0
10 // K-Loops (contracted)
11 for  $k_1 = 1 : S_{k_1}$ 
12 // ... remaining K-loops omitted ...
13 for  $k_\xi = 1 : S_{k_\xi}$ 
14     tmp +=  $\mathcal{A}_{\Pi^A(m_1, \dots, m_\gamma, k_1, \dots, k_\xi)} \mathcal{B}_{\Pi^B(k_1, \dots, k_\xi, n_1, \dots, n_\zeta)}$ 
15 // update  $\mathcal{C}$ 
16  $\mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)} = \alpha \text{ tmp} + \beta \mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)}$ 

```

Naive GETT.

```

1 // N-Loops
2 for  $n_1 = 1 : S_{n_1}$ 
3   // ... remaining N-loops omitted ...
4   for  $n_\zeta = 1 : S_{n_\zeta}$ 
5     // M-Loops
6     for  $m_1 = 1 : S_{m_1}$ 
7       // ... remaining M-loops omitted ...
8       for  $m_\gamma = 1 : S_{m_\gamma}$ 
9         tmp = 0
10        // K-Loops (contracted)
11        for  $k_1 = 1 : S_{k_1}$ 
12          // ... remaining K-loops omitted ...
13          for  $k_\xi = 1 : S_{k_\xi}$ 
14            tmp +=  $\mathcal{A}_{\Pi^A(m_1, \dots, m_\gamma, k_1, \dots, k_\xi)} \mathcal{B}_{\Pi^B(k_1, \dots, k_\xi, n_1, \dots, n_\zeta)}$ 
15            // update  $\mathcal{C}$ 
16             $\mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)} = \alpha \text{ tmp} + \beta \mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)}$ 

```

Naive GETT.

```

1 // N-Loops
2 for  $n_1 = 1 : S_{n_1}$ 
3 // ... remaining N-loops omitted ...
4 for  $n_\zeta = 1 : S_{n_\zeta}$ 
5 // M-Loops
6 for  $m_1 = 1 : S_{m_1}$ 
7 // ... remaining M-loops omitted ...
8 for  $m_\gamma = 1 : S_{m_\gamma}$ 
9     tmp = 0
10    // K-Loops (contracted)
11    for  $k_1 = 1 : S_{k_1}$ 
12        // ... remaining K-loops omitted ...
13        for  $k_\xi = 1 : S_{k_\xi}$ 
14            tmp +=  $\mathcal{A}_{\Pi^A(m_1, \dots, m_\gamma, k_1, \dots, k_\xi)} \mathcal{B}_{\Pi^B(k_1, \dots, k_\xi, n_1, \dots, n_\zeta)}$ 
15            // update  $\mathcal{C}$ 
16             $\mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)} = \alpha \text{ tmp} + \beta \mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)}$ 

```

Naive GETT.

```

1  // N-Loops
2  for  $n_1 = 1 : S_{n_1}$ 
3      // ... remaining N-loops omitted ...
4      for  $n_\zeta = 1 : S_{n_\zeta}$ 
5          // M-Loops
6          for  $m_1 = 1 : S_{m_1}$ 
7              // ... remaining M-loops omitted ...
8              for  $m_\gamma = 1 : S_{m_\gamma}$ 
9                  tmp = 0
10                 // K-Loops (contracted)
11                 for  $k_1 = 1 : S_{k_1}$ 
12                     // ... remaining K-loops omitted ...
13                     for  $k_\xi = 1 : S_{k_\xi}$ 
14                         tmp +=  $\mathcal{A}_{\Pi^A(m_1, \dots, m_\gamma, k_1, \dots, k_\xi)} \mathcal{B}_{\Pi^B(k_1, \dots, k_\xi, n_1, \dots, n_\zeta)}$ 
15                         // update  $\mathcal{C}$ 
16                          $\mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)} = \alpha \text{ tmp} + \beta \mathcal{C}_{\Pi^C(m_1, \dots, m_\gamma, n_1, \dots, n_\zeta)}$ 

```

Naive GETT.

```

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  $\tilde{B}$ 
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  $\tilde{A}$ 
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.


```

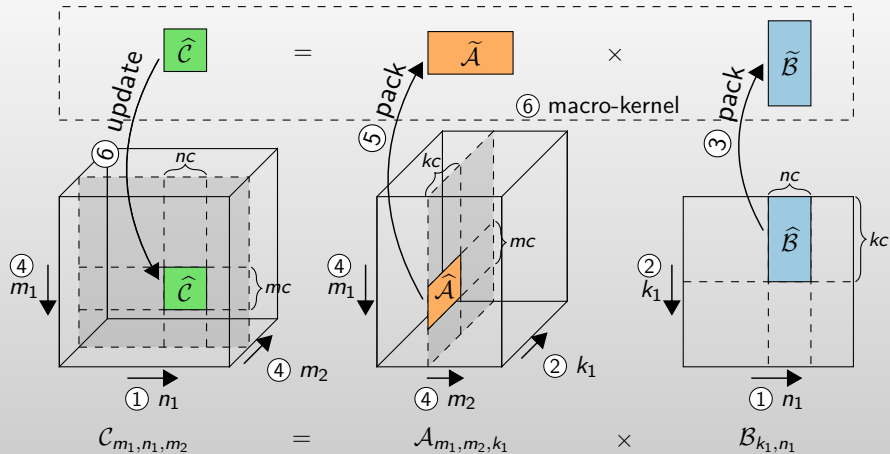
1 // N-Loop
2 for n = 1 : nc : Sln
3   // K-Loop (contracted)
4   for k = 1 : kc : Slk
5      $\hat{B}$  = identify_subtensor(B, n, k)
6     // pack  $\hat{B}$  into  $\tilde{B}$ 
7      $\tilde{B}$  = packB( $\hat{B}$ )
8     // M-Loop
9     for m = 1 : mc : Slm
10       $\hat{A}$  = identify_subtensor(A, m, k)
11      // pack  $\hat{A}$  into  $\tilde{A}$ 
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( $\tilde{A}, \tilde{B}, \hat{C}, \alpha, \beta$ )

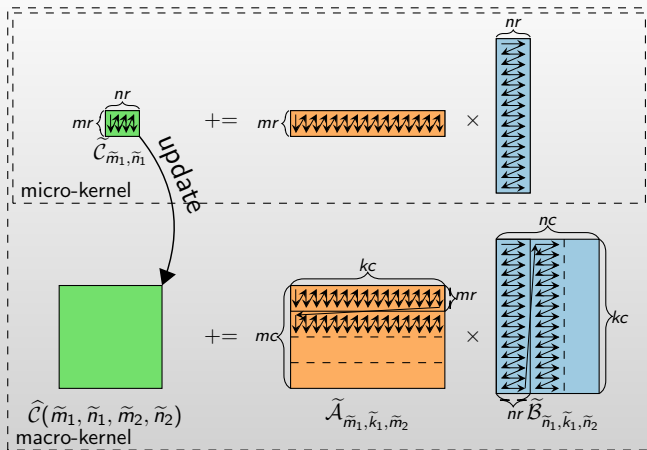
```

High-performance GETT.

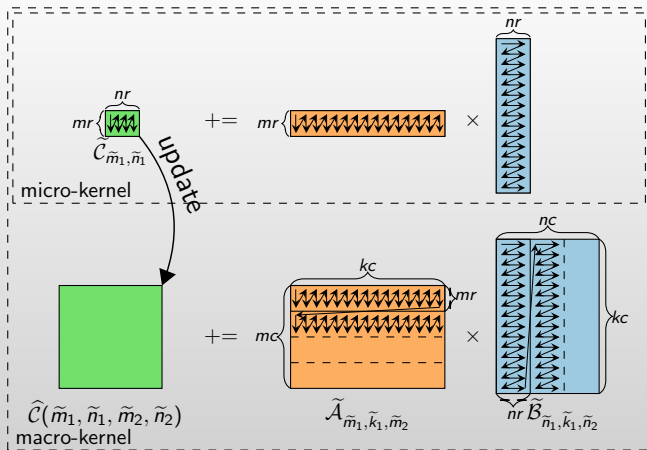
Key Idea

Pack-and-transpose while moving data into the caches

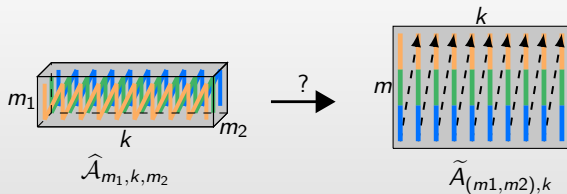


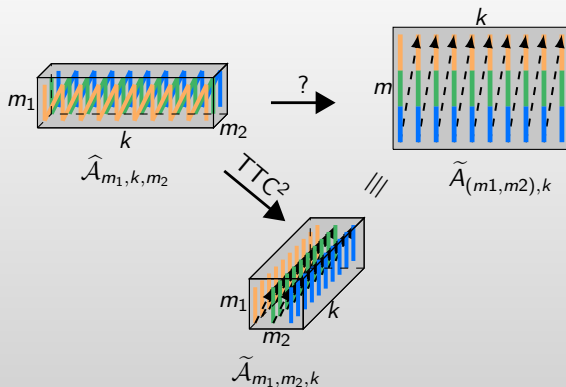


- Blocking for L3, L2, L1 cache as well as registers



- Blocking for L3, L2, L1 cache as well as registers
- Written in AVX2 intrinsics





- Preserve stride-1 index
 ⇒ Efficient packing routines

²Paul Springer, Jeff R. Hammond, and Paolo Bientinesi. "TTC: A high-performance Compiler for Tensor Transpositions". In: *TOMS, in review* ().

- Blocking for caches
- Blocking for registers
- Explicitly vectorized
- Use TTC to generate high-performance packing routines
 - Exploits full cache line (avoids non-stride-one memory accesses)
- Explore large search-space:
 - Different GEMM-variants (e.g., panel-matrix, matrix-panel)
 - Different permutations
 - Different values for mc , nc and kc
- Prune the search space via a performance model

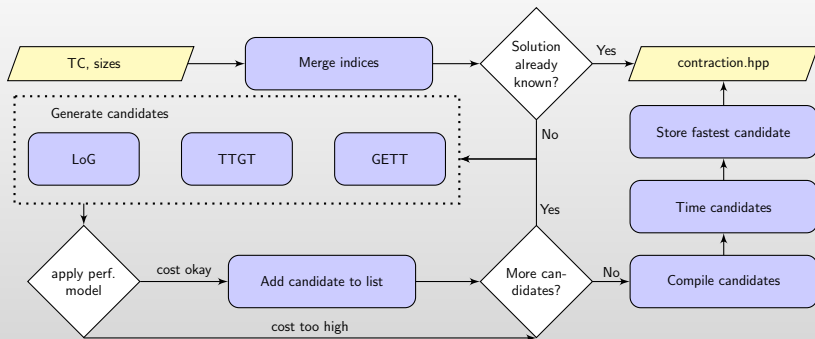
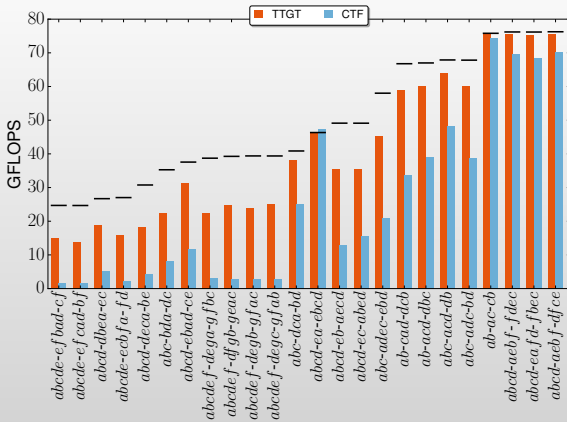
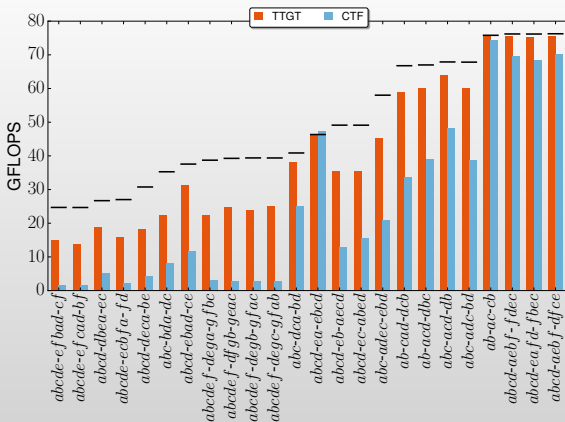


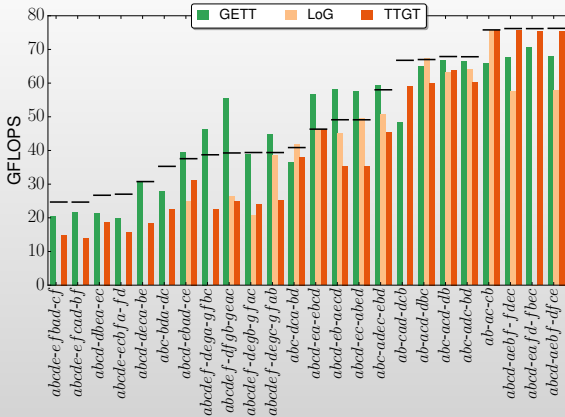
Figure: Schematic overview of TCCG.

- System: Intel *Xeon E5-2680 v3* CPU (Haswell)
 - Single core
 - Turbo Boost: disabled
- Compiler: *icpc 16.0.1 20151021*
- Benchmark
 - Collection of 48 TCs
 - Compiled from four publications
 - Each TC is at least 200 MiB
- Correctness checked against naive loop-based implementation

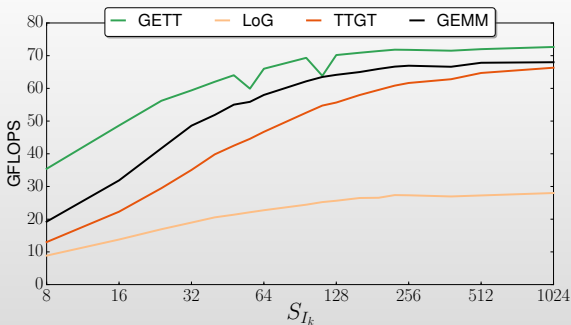




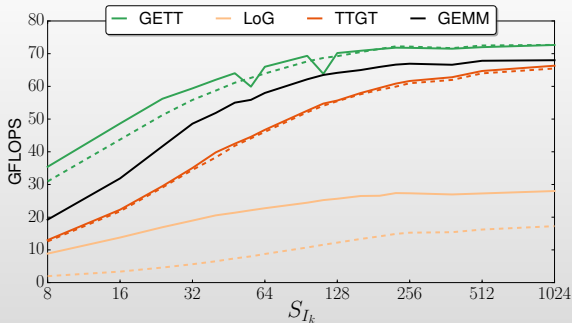
- TTGT good in compute-bound regime
- TTGT bad in bandwidth-bound regime
- TTGT faster than CTF everywhere.



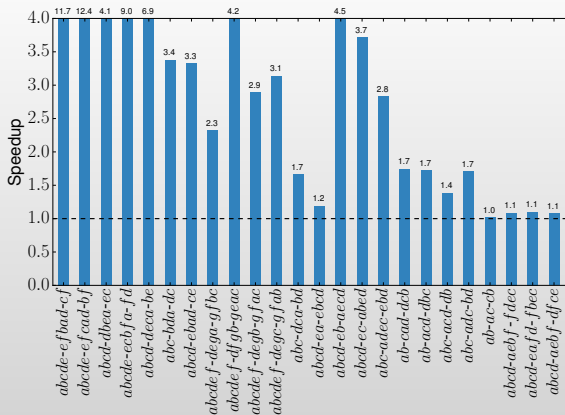
- GETT excels in bandwidth-bound regime.
- GETT slightly lags behind in compute-bound regime.



- GETT especially good in bandwidth-bound regime
 - GETT still attains up to 91.3% of peak floating-point performance
- TTGT poor in bandwidth-bound regime



- GETT especially good in bandwidth-bound regime
 - GETT still attains up to 91.3% of peak floating-point performance
- TTGT poor in bandwidth-bound regime
- LoG performance can become arbitrarily bad
- GETT and TTGT barely affected by higher dimensions



- Speedup varies between $1.0\times$ and $12.4\times$

Conclusion

- GETT: a systematic way to reduce an arbitrary TC to a GEMM-like macro-kernel
- GETT exhibits high performance across a wide range of TCs
 - It especially excels in the bandwidth-bound regime
 - It attains up to 91.3% of peak floating-point performance
- A survey of different approaches to TCs has been presented
- Give it a try: <https://github.com/HPAC/tccg>

Conclusion

- GETT: a systematic way to reduce an arbitrary TC to a GEMM-like macro-kernel
- GETT exhibits high performance across a wide range of TCs
 - It especially excels in the bandwidth-bound regime
 - It attains up to 91.3% of peak floating-point performance
- A survey of different approaches to TCs has been presented
- Give it a try: <https://github.com/HPAC/tccg>

Future Work

- Assess TCCG's performance on KNL
- Add parallelism
- Turn TCCG into a C library?

Conclusion

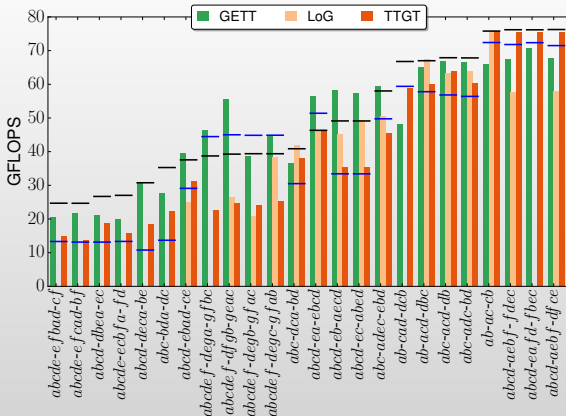
- GETT: a systematic way to reduce an arbitrary TC to a GEMM-like macro-kernel
- GETT exhibits high performance across a wide range of TCs
 - It especially excels in the bandwidth-bound regime

Thank you for your attention.

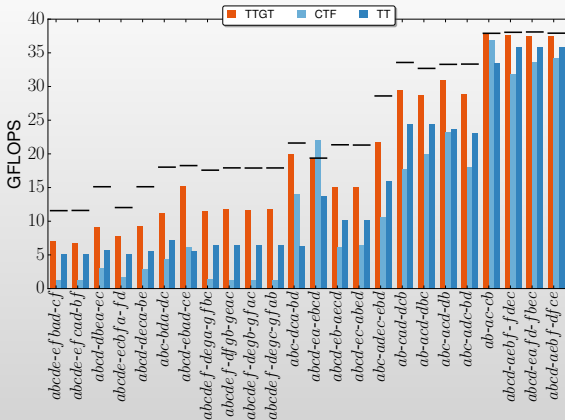
- Give it a try: <https://github.com/HPAC/tccg>

Future Work

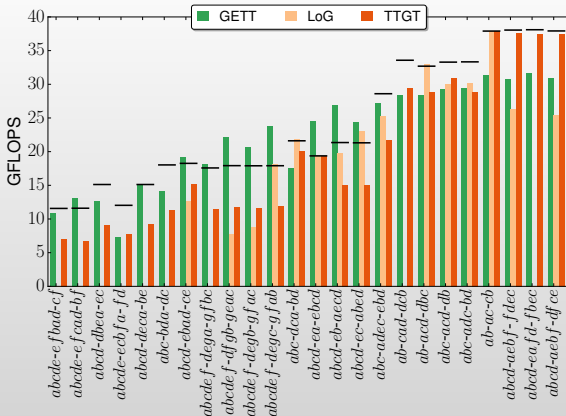
- Assess TCCG's performance on KNL
- Add parallelism
- Turn TCCG into a C library?



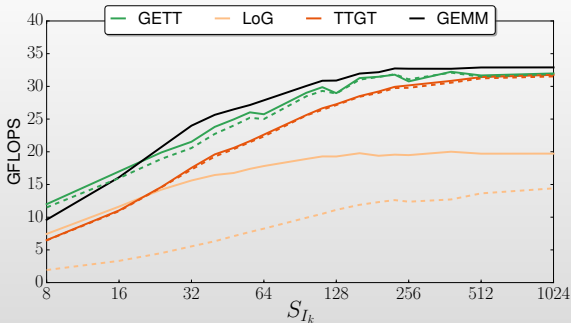
- GETT excels in bandwidth-bound regime.
- GETT slightly lags behind in compute-bound regime.
- GETT attains min/avg/max performance of GEMM:
 - SP: 72.4% / 98.1% / 141.4%
 - DP: 60.8% / 97.0% / 132.9%



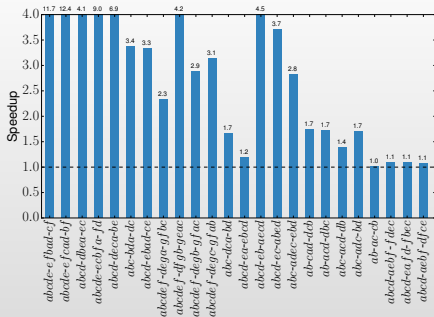
- TTGT faster than CTF everywhere.
- TTGT good in compute-bound regime
- TTGT bad in bandwidth-bound regime



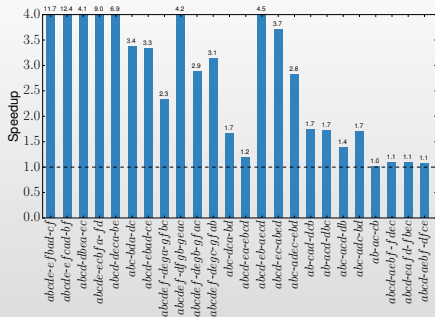
- GETT excels in bandwidth-bound regime.
- GETT slightly lags behind in compute-bound regime.
- GETT attains min/avg/max performance of GEMM:
 - SP: 72.4% / 98.1% / 141.4%
 - DP: 60.8% / 97.0% / 132.9%



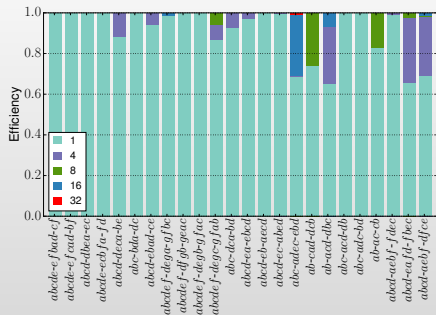
- GETT especially good in bandwidth-bound regime
 - GETT still attains up to 91.3% of peak floating-point performance
- TTGT poor in bandwidth-bound regime
- LoG performance can become arbitrarily bad
- GETT and TTGT barely affected by higher dimensions



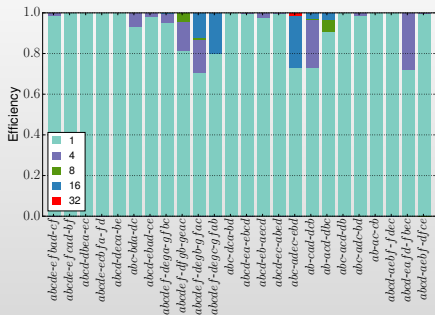
(a) Single-Precision.



(b) Double-Precision.



(a) Single-Precision.



(b) Double-Precision.

Figure: Limit the GETT candidates to 1, 4, 8, 16 or 32, respectively.

- Average performance without search: 90.7% / 92.3%
- Average performance of the four best candidates: 98.3% / 97.2%


```

C[a,b,i,j] = A[i,m,a] * B[m,j,b]
a = 24
b = 24
i = 24
j = 24
m = 24

```

Figure: Exemplary input file for TCCG.

Argument	Description
<code>--floatType=[s,d]</code>	data type
<code>--maxWorkspace=<value></code>	maximum auxiliary workspace in GB
<code>--maxImplementations=<value></code>	maximum #implementations
<code>--arch=[hsw,knl,cuda]</code>	selected architecture
<code>--numThreads=<value></code>	number of threads

Table: TCCG's command line arguments.

```

1  $A_{\Pi^m(I_m), \Pi^k(I_k)} \leftarrow \mathcal{A}_{\Pi^{\mathcal{A}}(I_m \cup I_k)}$  // unfold  $\mathcal{A}$ 
2  $B_{\Pi^k(I_k), \Pi^n(I_n)} \leftarrow \mathcal{B}_{\Pi^{\mathcal{B}}(I_n \cup I_k)}$  // unfold  $\mathcal{B}$ 
3  $X_{\Pi^m(I_m), \Pi^n(I_n)} \leftarrow \text{op}(\mathcal{A})_{\Pi^m(I_m), \Pi^k(I_k)} \times \text{op}(\mathcal{B})_{\Pi^k(I_k), \Pi^n(I_n)}$  // contract  $\mathcal{A}$  and  $\mathcal{B}$  via a GEMM
4  $C_{\Pi^c(I_m \cup I_n)} \leftarrow X_{\Pi^m(I_m), \Pi^n(I_n)}$  // fold  $X$ 
    
```

TTGT pseudo code for a general tensor contraction

$$C_{\Pi^c(I_m \cup I_n)} = \mathcal{A}_{\Pi^{\mathcal{A}}(I_m \cup I_k)} \mathcal{B}_{\Pi^{\mathcal{B}}(I_n \cup I_k)} + C_{\Pi^c(I_m \cup I_n)}.$$

- $\Pi^m(I_m)$, $\Pi^n(I_n)$ and $\Pi^k(I_k)$ represent arbitrary, but fixed, permutations
- Transpositions account for pure overhead
- Requires additional memory
- Good if GEMM dominates performance (i.e., compute-bound)
- Bad if transpositions dominate performance (i.e., bandwidth-bound)

- Loop over 2D slices of the tensors
- Contract these 2D slices via GEMM

Advantages

- Exploits GEMM's high-performance
- No additional memory

Disadvantages

- Performance can become arbitrarily poor
- Sometimes not applicable (if stride-one accesses are required)

$$C_{m_1, n_1, m_2, n_2} = A_{m_1, m_2, k_1} B_{k_1, n_1, n_2}$$

```
for m2 = 0: M2
  for n2 = 0: N2
    GEMM( &A[m2 * M1], &B[n2 * K1 * N1], &C[
      m2 * M1 * n1 + n2 * m1 * N1 * M2])
```