

The Landscape of High-Performance Tensor Contractions

Paul Springer and Paolo Bientinesi

Aachen Institute for Advanced Study in
Computational Engineering Science

Atlanta, Feb. 24th 2017



- A tensors is a multidimensional array:
 - 0-order tensor: scalar α

- A tensors is a multidimensional array:
 - 0-order tensor: scalar α
 - 1-order tensor: vector \mathcal{A}_{i_1}

- A tensors is a multidimensional array:
 - 0-order tensor: scalar α
 - 1-order tensor: vector \mathcal{A}_{i_1}
 - 2-order tensor: matrix \mathcal{A}_{i_1, i_2}

- A tensors is a multidimensional array:
 - 0-order tensor: scalar α
 - 1-order tensor: vector \mathcal{A}_{i_1}
 - 2-order tensor: matrix \mathcal{A}_{i_1, i_2}
 - n -order tensor: $\mathcal{A}_{i_1, i_2, \dots, i_n}$

¹Paul Springer et al. "Design of a high-performance GEMM-like Tensor-Tensor Multiplication"

- A tensors is a multidimensional array:
 - 0-order tensor: scalar α
 - 1-order tensor: vector \mathcal{A}_{i_1}
 - 2-order tensor: matrix \mathcal{A}_{i_1, i_2}
 - n -order tensor: $\mathcal{A}_{i_1, i_2, \dots, i_n}$
- Tensor contractions can be thought of as generalized GEMMs

¹Paul Springer et al. "Design of a high-performance GEMM-like Tensor-Tensor Multiplication"

- A tensors is a multidimensional array:
 - 0-order tensor: scalar α
 - 1-order tensor: vector \mathcal{A}_{i_1}
 - 2-order tensor: matrix \mathcal{A}_{i_1, i_2}
 - n -order tensor: $\mathcal{A}_{i_1, i_2, \dots, i_n}$
- Tensor contractions can be thought of as generalized GEMMs
- Three approaches to tensor contractions:
 - Nested loops
 - Loops over GEMM (LoG)
 - Transpose-Transpose-GEMM-Transpose (TTGT)

¹Paul Springer et al. "Design of a high-performance GEMM-like Tensor-Tensor Multiplication"

- A tensors is a multidimensional array:
 - 0-order tensor: scalar α
 - 1-order tensor: vector \mathcal{A}_{i_1}
 - 2-order tensor: matrix \mathcal{A}_{i_1, i_2}
 - n -order tensor: $\mathcal{A}_{i_1, i_2, \dots, i_n}$
- Tensor contractions can be thought of as generalized GEMMs
- Three approaches to tensor contractions:
 - Nested loops
 - Loops over GEMM (LoG)
 - Transpose-Transpose-GEMM-Transpose (TTGT)
- We propose a novel approach: GETT¹
 - Akin to a high-performance GEMM implementation

¹Paul Springer et al. "Design of a high-performance GEMM-like Tensor-Tensor Multiplication"

- Approaches to Tensor Contractions:
 - Loops over GEMM (LoG)
 - Transpose-Transpose-GEMM-Transpose (TTGT)
 - GEMM-like Tensor-Tensor Multiply (GETT)
- Tensor Contraction Code Generator²
- Performance Evaluation

²Source code available at: <https://github.com/HPAC/tccg>

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

Conceptual Idea

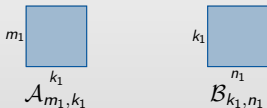
Identify 2D subtensors and contract them via GEMM

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

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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



```
gemm(M1, N1, K1, A[:, :], B[:, :], C[:, :])
```

Conceptual Idea

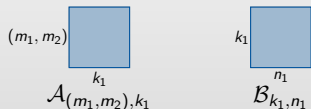
Identify 2D subtensors and contract them via GEMM

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

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

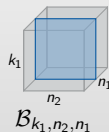
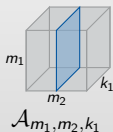


```
gemm(M1 x M2, N1, K1, A[:, :], B[:, :], C[:, :])
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

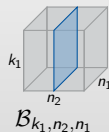
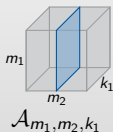


```
for( m2 = 0; m2 < M2; m2++ )
  for( n1 = 0; n1 < N1; n1++ )
    gemm( M1, N2, K1, A[:, m2, :], B[:, :, n1], C[:, n1, :, m2] )
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

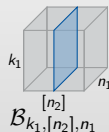
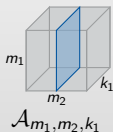


```
for( m2 = 0; m2 < M2; m2++ )
  for( n2 = 0; n2 < N2; n2++ )
    gemm( M1, N1, K1, A[:, m2, :], B[:, n2, :], C[:, :, n2, m2] )
```


Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

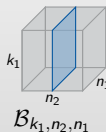
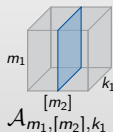


```
for( m2 = 0; m2 < M2; m2++ )
    gemm_batch( M1, N1, K1, A[:, m2, :], B[:, n2, :], C[:, :, n2, m2], N2 )
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

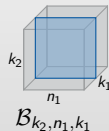
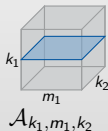


```
for( n2 = 0; n2 < N2; n2++ )
    gemm_batch( M1, N1, K1, A[:, m2, :], B[:, n2, :], C[:, :, n2, m2], M2 )
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

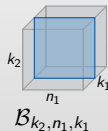
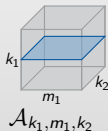


```
for( k1 = 0; k1 < K1; k1++ )
    gemm(M1, N1, K2, A[k1, :, :], B[:, :, k1], C[:, :])
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

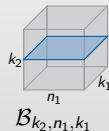
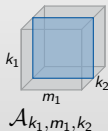


```
for ( k1 = 0; k1 < K1; k1++ )
    gemm( M1, N1, K2, A[k1, :, :], B[:, :, k1], C[:, :] )
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

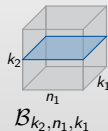
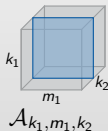


```
for ( k2 = 0; k2 < K2; k2++ )
    gemm( M1, N1, K1, A[:, :, k2]^T, B[k2, :, :]^T, C[:, :] )
```

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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

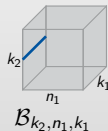
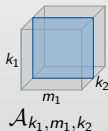


~~for ($k_2 = 0$; $k_2 < K_2$; k_2++)
~~gemm($M_1, N_1, K_1, A[:, :, k_2]^T, B[k_2, :, :]^T, C[:, :]$)~~~~

Conceptual Idea

Identify 2D subtensors and contract them via GEMM

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



```
for( n = 0; n1 < N1; n1++ )
  for( k2 = 0; k2 < K2; k2++ )
    gemv( M1, K1, A[:, :, k2]^T, B[k2, n1, :], C[:, n] )
```

- Search space:

- Search space:
 - GEMM indices: m, n, k

- Search space:
 - GEMM indices: m, n, k
 - Loop order

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:
 - Easy to implement

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - No additional memory required

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - No additional memory required
- Disadvantages:

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - No additional memory required
- Disadvantages:
 - Some contractions cannot be implemented via straight LoG

- Search space:
 - GEMM indices: m, n, k
 - Loop order
 - Batched index
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - No additional memory required
- Disadvantages:
 - Some contractions cannot be implemented via straight LoG
 - GEMM's arithmetic intensity can be suboptimal

- Free indices of \mathcal{A}

- $I_m := \{m_1, m_2, \dots, m_\gamma\} = I_{\mathcal{A}} \cap I_c$

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$
- Contracted indices
 - $l_k := \{k_1, k_2, \dots, k_\xi\} = l_{\mathcal{A}} \cap l_{\mathcal{B}}$

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$
- Contracted indices
 - $l_k := \{k_1, k_2, \dots, k_\xi\} = l_{\mathcal{A}} \cap l_{\mathcal{B}}$
- Tensor contractions can be mapped to BLAS routines^{3,4}:
 - GEMM: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k \neq \emptyset$.

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$
- Contracted indices
 - $l_k := \{k_1, k_2, \dots, k_\xi\} = l_{\mathcal{A}} \cap l_{\mathcal{B}}$
- Tensor contractions can be mapped to BLAS routines^{3,4}:
 - GEMM: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k \neq \emptyset$.
 - GEMV: ($l_m = \emptyset$ **or** $l_n = \emptyset$) **and** $l_k \neq \emptyset$

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$
- Contracted indices
 - $l_k := \{k_1, k_2, \dots, k_\xi\} = l_{\mathcal{A}} \cap l_{\mathcal{B}}$
- Tensor contractions can be mapped to BLAS routines^{3,4}:
 - GEMM: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k \neq \emptyset$.
 - GEMV: $(l_m = \emptyset$ **or** $l_n = \emptyset)$ **and** $l_k \neq \emptyset$
 - GER: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k = \emptyset$

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$
- Contracted indices
 - $l_k := \{k_1, k_2, \dots, k_\xi\} = l_{\mathcal{A}} \cap l_{\mathcal{B}}$
- Tensor contractions can be mapped to BLAS routines^{3,4}:
 - GEMM: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k \neq \emptyset$.
 - GEMV: $(l_m = \emptyset$ **or** $l_n = \emptyset)$ **and** $l_k \neq \emptyset$
 - GER: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k = \emptyset$
 - AXPY: $(l_m = \emptyset$ **or** $l_n = \emptyset)$ **and** $l_k = \emptyset$

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

- Free indices of \mathcal{A}
 - $l_m := \{m_1, m_2, \dots, m_\gamma\} = l_{\mathcal{A}} \cap l_{\mathcal{C}}$
- Free indices of \mathcal{B}
 - $l_n := \{n_1, n_2, \dots, n_\zeta\} = l_{\mathcal{B}} \cap l_{\mathcal{C}}$
- Contracted indices
 - $l_k := \{k_1, k_2, \dots, k_\xi\} = l_{\mathcal{A}} \cap l_{\mathcal{B}}$
- Tensor contractions can be mapped to BLAS routines^{3,4}:
 - GEMM: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k \neq \emptyset$.
 - GEMV: $(l_m = \emptyset$ **or** $l_n = \emptyset)$ **and** $l_k \neq \emptyset$
 - GER: $l_m \neq \emptyset$ **and** $l_n \neq \emptyset$ **and** $l_k = \emptyset$
 - AXPY: $(l_m = \emptyset$ **or** $l_n = \emptyset)$ **and** $l_k = \emptyset$
 - DOT: else.

³Di Napoli et al. "Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions"

⁴Yang Shi et al. "Tensor Contractions with Extended BLAS Kernels on CPU and GPU"

Conceptual Idea

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

Conceptual Idea

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

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

Conceptual Idea

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

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

$$\begin{aligned} \tilde{A}_{m_1, (k_1, k_2)} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}, \tilde{B}, C) \end{aligned}$$

Conceptual Idea

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

$$\bullet C_{m_1, n_1} \leftarrow A_{k_1, m_1, k_2} B_{k_2, n_1, k_1}$$

$$\begin{aligned} \tilde{A}_{m_1, (k_1, k_2)} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}^T, \tilde{B}, C) \end{aligned}$$

Conceptual Idea

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

$$\bullet C_{m_1, n_1} \leftarrow A_{k_1, m_1, k_2} B_{k_2, n_1, k_1}$$

$$\begin{aligned} \tilde{A}_{m_1, (k_1, k_2)} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}^T, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{B}^T, \tilde{A}, \tilde{C}) \\ C_{m_1, n_1} &\leftarrow \tilde{C}_{m_1, m_1} \end{aligned}$$

Conceptual Idea

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

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

$$\begin{aligned} \tilde{A}_{m_1, (k_1, k_2)} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}^T, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{B}^T, \tilde{A}, \tilde{C}) \\ C_{m_1, n_1} &\leftarrow \tilde{C}_{m_1, m_1} \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_2, k_1), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_2, k_1), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{B}^T, \tilde{A}, \tilde{C}) \\ C_{m_1, n_1} &\leftarrow \tilde{C}_{n_1, m_1} \end{aligned}$$

Conceptual Idea

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

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

$$\begin{aligned} \tilde{A}_{m_1, (k_1, k_2)} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{A}^T, \tilde{B}, C) \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_1, k_2), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_1, k_2), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{B}^T, \tilde{A}, \tilde{C}) \\ C_{m_1, n_1} &\leftarrow \tilde{C}_{m_1, m_1} \end{aligned}$$

$$\begin{aligned} \tilde{A}_{(k_2, k_1), m_1} &\leftarrow A_{k_1, m_1, k_2} \\ \tilde{B}_{(k_2, k_1), n_1} &\leftarrow B_{k_2, n_1, k_1} \\ \text{gemm}(M_1, N_1, K_1 \times K_2, \tilde{B}^T, \tilde{A}, \tilde{C}) \\ C_{m_1, n_1} &\leftarrow \tilde{C}_{m_1, m_1} \end{aligned}$$

... and more.

- Search space:

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of l_m, l_n, l_k

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of I_m, I_n, I_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of I_m, I_n, I_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - All TCs can be implemented via TTGT

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of I_m, I_n, I_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - All TCs can be implemented via TTGT
 - Large GEMM \Rightarrow good performance?

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of I_m, I_n, I_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - All TCs can be implemented via TTGT
 - Large GEMM \Rightarrow good performance?
- Disadvantages:

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - All TCs can be implemented via TTGT
 - Large GEMM \Rightarrow good performance?
- Disadvantages:
 - Transpositions⁵ account for pure overhead

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- Search space:
 - Any permutation of l_m, l_n, l_k
 - Transposed \mathcal{A}
 - Transposed \mathcal{B}
 - Interchange \mathcal{A} and \mathcal{B} within GEMM
- Advantages:
 - Easy to implement
 - Exploits existing BLAS libraries
 - All TCs can be implemented via TTGT
 - Large GEMM \Rightarrow good performance?
- Disadvantages:
 - Transpositions⁵ account for pure overhead
 - Additional memory required

⁵Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

Key Idea

- Eliminate explicit transpositions
- Pack-and-transpose while moving data into the caches⁵
 - ⇒ Complexity offloaded into packing routines

⁵Sharan Chetlur et al. "CUDA NN: Efficient Primitives for Deep Learning"

Key Idea

- Eliminate explicit transpositions
- Pack-and-transpose while moving data into the caches⁵
 - ⇒ Complexity offloaded into packing routines

```

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( $\tilde{A}$ ,  $\tilde{B}$ ,  $\hat{C}$ ,  $\alpha$ ,  $\beta$ )

```

⁵Sharan Chetlur et al. "CUDA: Efficient Primitives for Deep Learning"

Key Idea

- Eliminate explicit transpositions
- Pack-and-transpose while moving data into the caches⁵
 - ⇒ Complexity offloaded into packing routines

```

1 // N-Loop
2 for n = 1 : nc : Sn

```

Monday, February 27 10:15 - 10:35

TCCG: Tensor Contraction Code Generator

```

3 // ...
4 // ...
5 // ...
6 // ...
7 // ...
8 // ...
9 // ...
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( $\tilde{A}, \tilde{B}, \hat{C}, \alpha, \beta$ )

```

⁵Sharan Chetlur et al. "CUDA: Efficient Primitives for Deep Learning"

- Search space:

- Search space:
 - Blocking parameters: mc, nc, kc

- Search space:
 - Blocking parameters: mc, nc, kc
 - Subtensors $\hat{\mathcal{A}}, \hat{\mathcal{B}}, \hat{\mathcal{C}}$

- Search space:
 - Blocking parameters: mc, nc, kc
 - Subtensors $\hat{\mathcal{A}}, \hat{\mathcal{B}}, \hat{\mathcal{C}}$
- Advantages:

- Search space:
 - Blocking parameters: mc, nc, kc
 - Subtensors $\hat{\mathcal{A}}, \hat{\mathcal{B}}, \hat{\mathcal{C}}$
- Advantages:
 - Same arithmetic intensity as GEMM

- Search space:
 - Blocking parameters: mc, nc, kc
 - Subtensors $\hat{A}, \hat{B}, \hat{C}$
- Advantages:
 - Same arithmetic intensity as GEMM
 - No memory overhead

- Search space:
 - Blocking parameters: mc, nc, kc
 - Subtensors $\hat{A}, \hat{B}, \hat{C}$
- Advantages:
 - Same arithmetic intensity as GEMM
 - No memory overhead
- Disadvantages:

- Search space:
 - Blocking parameters: mc, nc, kc
 - Subtensors $\hat{A}, \hat{B}, \hat{C}$
- Advantages:
 - Same arithmetic intensity as GEMM
 - No memory overhead
- Disadvantages:
 - Complex to implement

- **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

- **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

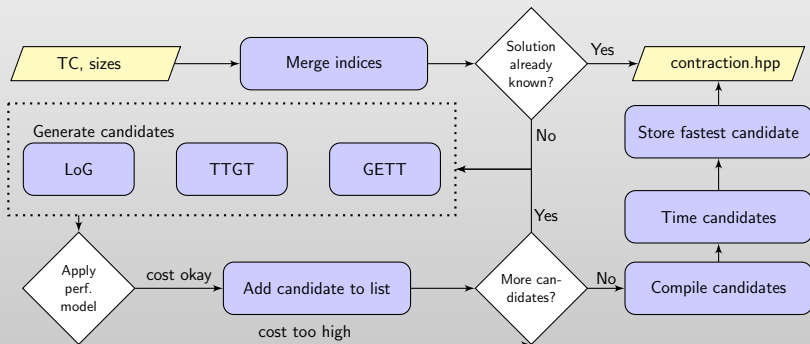
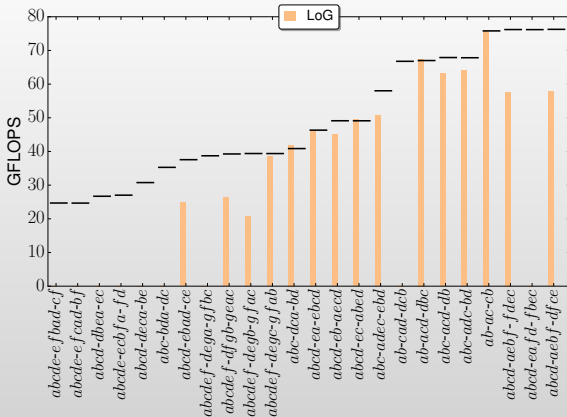
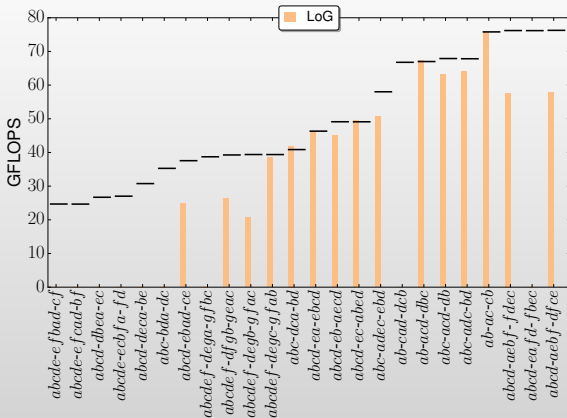
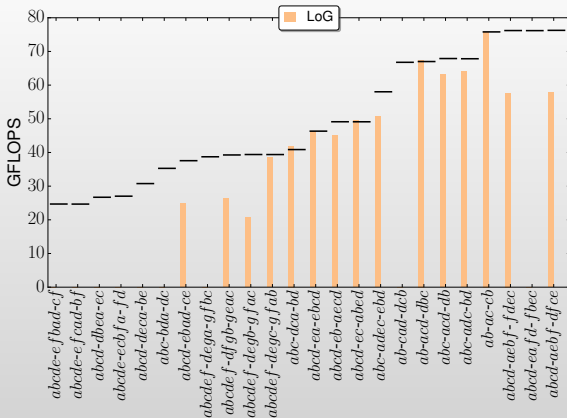


Figure: Schematic overview of TCCG.



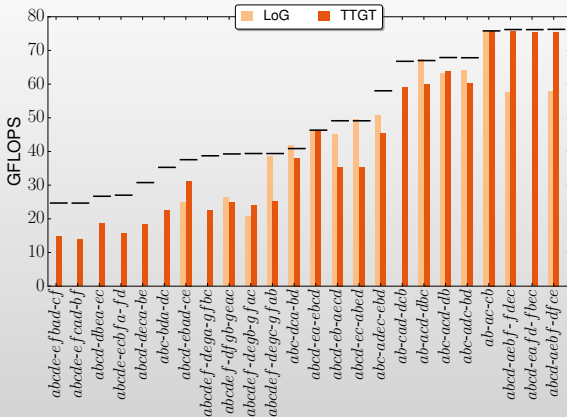


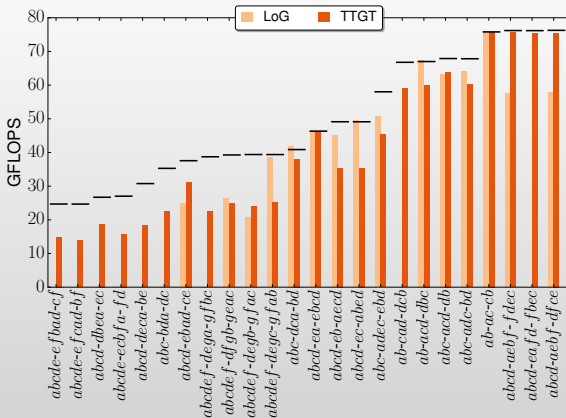
- Not all TCs can be implemented via LoG



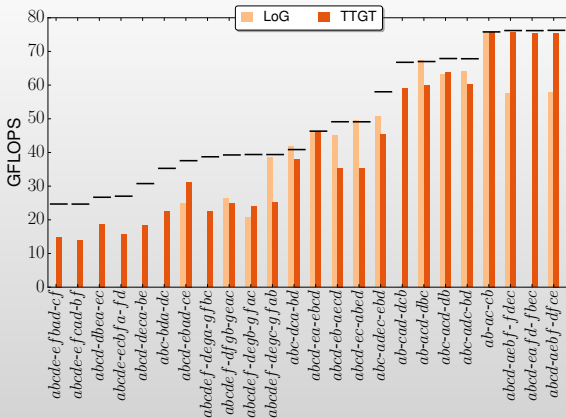
- Not all TCs can be implemented via LoG
- Mixed performance

Performance — Haswell (single core)



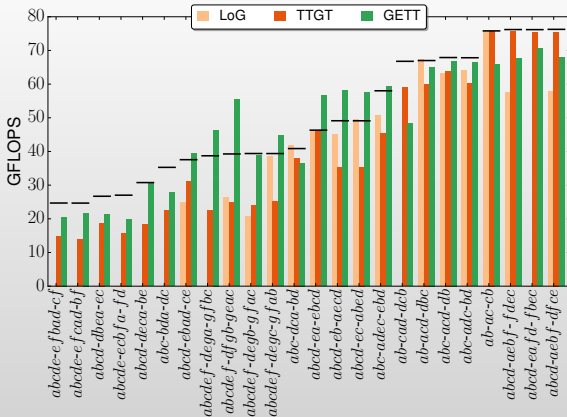


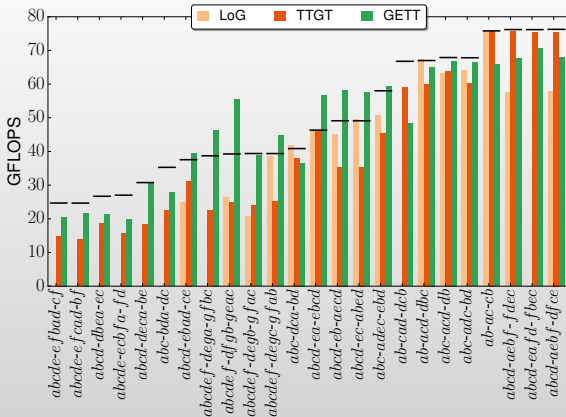
- TTGT: good for compute-bound TCs



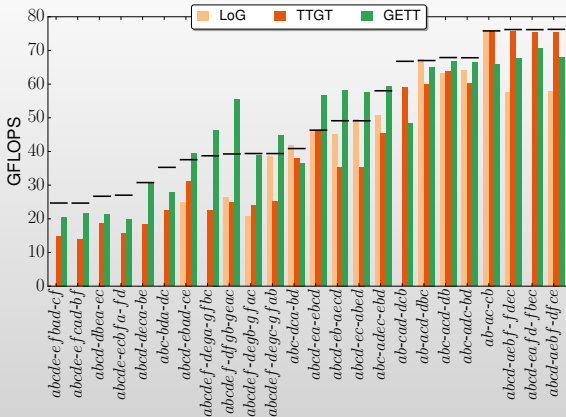
- TTGT: good for compute-bound TCs
- TTGT: bad for bandwidth-bound TCs

Performance — Haswell (single core)

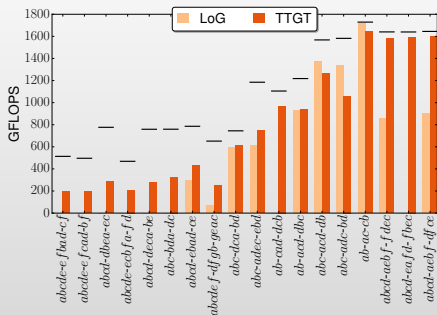




- GETT: excels for bandwidth-bound TCs

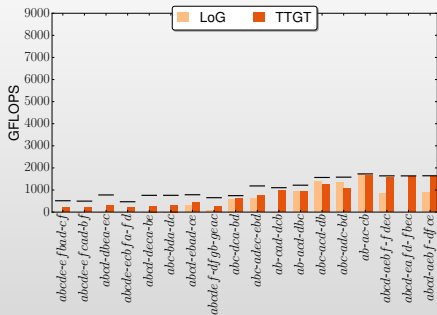


- GETT: excels for bandwidth-bound TCs
- GETT: good for compute-bound TCs

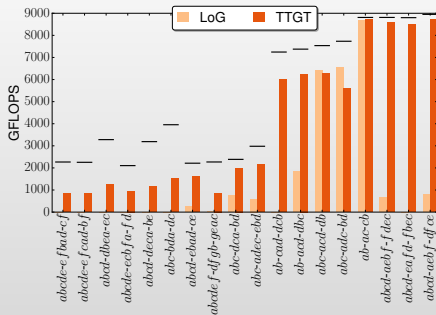


(a) 2x Intel Xeon E5-2680 v3

- Performance gap increases for bandwidth-bound TCs

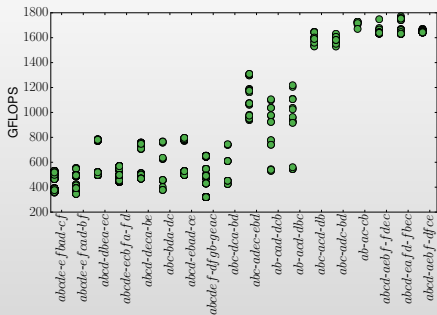


(a) 2x Intel Xeon E5-2680 v3

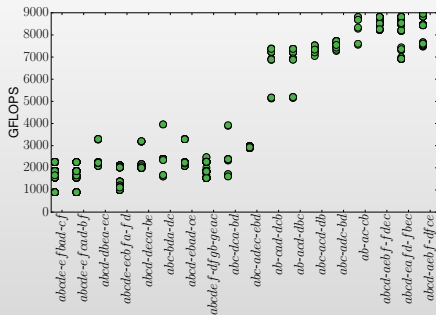


(b) NVIDIA Tesla P100

- Performance gap increases for bandwidth-bound TCs

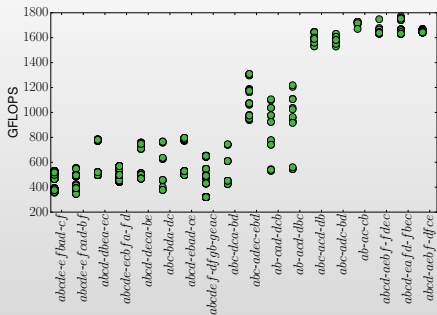


(a) 2×Intel Xeon E5-2680 v3

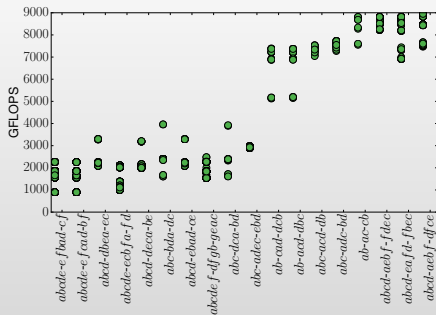


(b) NVIDIA Tesla P100

⁶Elmar Peise et al. "On the Performance Prediction of BLAS-based Tensor Contractions"



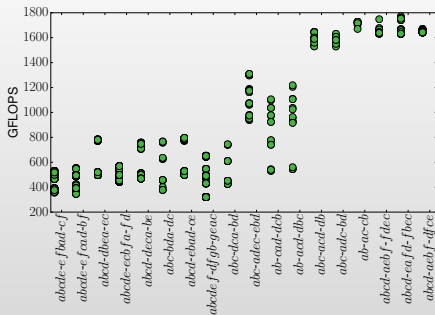
(a) 2×Intel Xeon E5-2680 v3



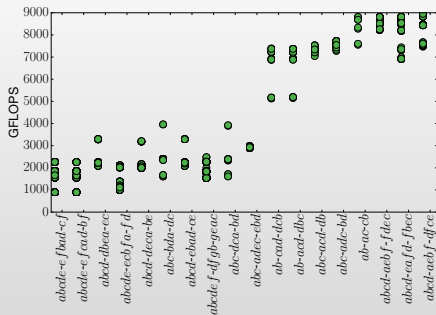
(b) NVIDIA Tesla P100

- Performance for equally-sized GEMMs varies greatly
 - For different settings: opA, opB, interchanged \mathcal{A} and \mathcal{B}

⁶Elmar Peise et al. "On the Performance Prediction of BLAS-based Tensor Contractions"



(a) 2×Intel Xeon E5-2680 v3



(b) NVIDIA Tesla P100

- Performance for equally-sized GEMMs varies greatly
 - For different settings: opA, opB, interchanged \mathcal{A} and \mathcal{B}
- Performance Model for TTGT and LoG:
 - Account for varying GEMM perf⁶

⁶Elmar Peise et al. "On the Performance Prediction of BLAS-based Tensor Contractions"

Conclusion

- A survey of different approaches to TCs has been presented
- GETT exhibits high performance across a wide range of TCs
- TCCG is available at <https://github.com/HPAC/tccg>

Conclusion

- A survey of different approaches to TCs has been presented
- GETT exhibits high performance across a wide range of TCs
- TCCG is available at <https://github.com/HPAC/tccg>

Future Work

- Implement TC library based on GETT
- Parallelize GETT

Conclusion

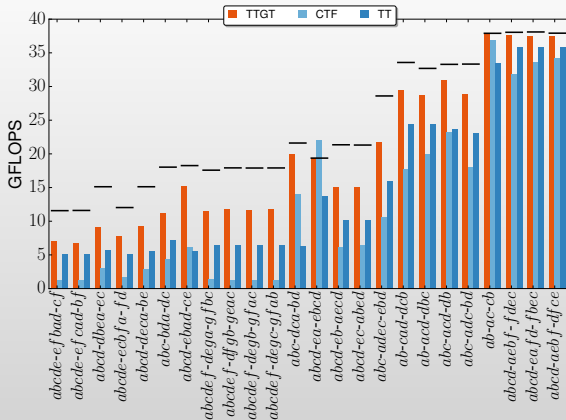
- A survey of different approaches to TCs has been presented
- GETT exhibits high performance across a wide range of TCs
- TCCG is available at <https://github.com/HPAC/tccg>

Future Work

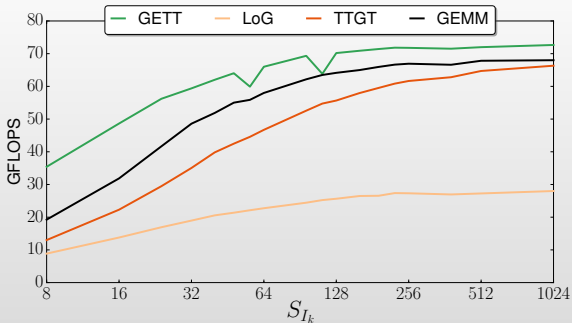
- Implement TC library based on GETT
- Parallelize GETT

Thank you for your attention.

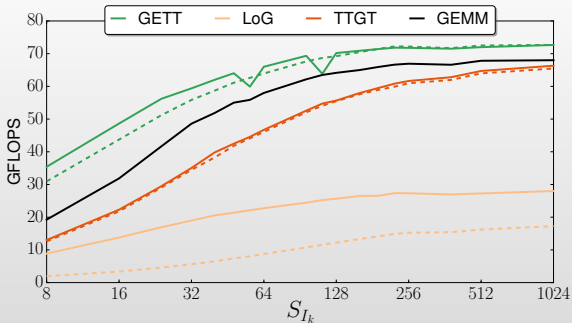
- Systems:
 - Intel *Xeon E5-2680 v3* CPU (Haswell)
 - NVIDIA *Tesla P100* GPU (Pascal)
- Compilers:
 - *icpc 16.0.1 20151021*
 - *nvcc v8.0.44*
- Benchmark
 - Collection of 48 TCs
 - Compiled from four publications
 - Each TC is at least 200 MiB
- Correctness checked against naive loop-based implementation



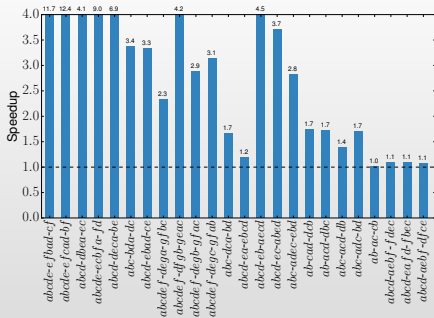
- TTGT faster than CTF everywhere.
- TTGT good in compute-bound regime
- TTGT bad in bandwidth-bound regime



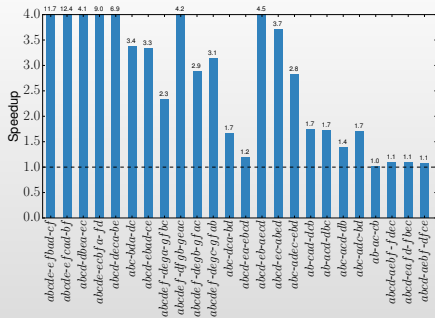
- GETT especially good in bandwidth-bound regime
 - GETT still attains up to 91.3% of peak floating-point performance
- TTGT poor in bandwidth-bound regime



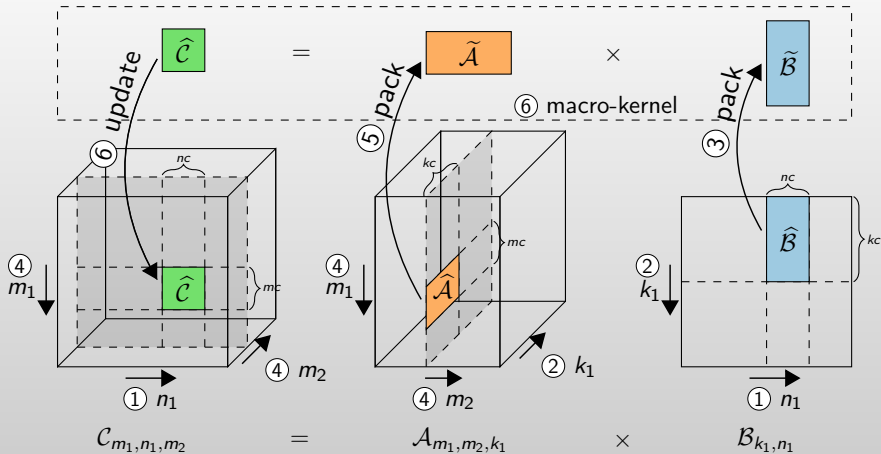
- GETT especially good in bandwidth-bound regime
 - GETT still attains up to 91.3% of peak floating-point performance
- TTGT poor in bandwidth-bound regime
- LoG performance can become arbitrarily bad
- GETT and TTGT barely affected by higher dimensions

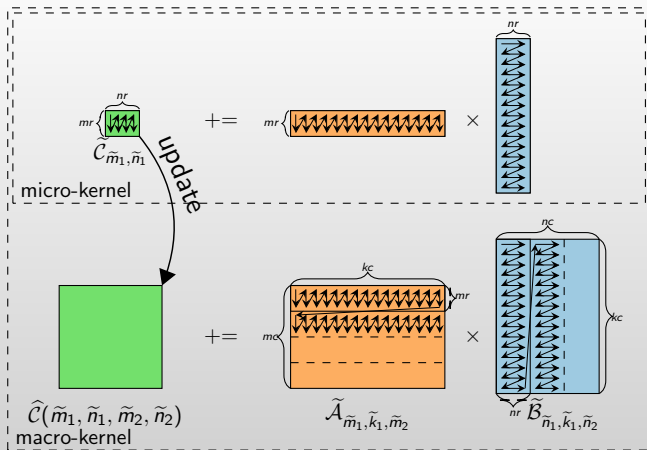


(a) Single-Precision.

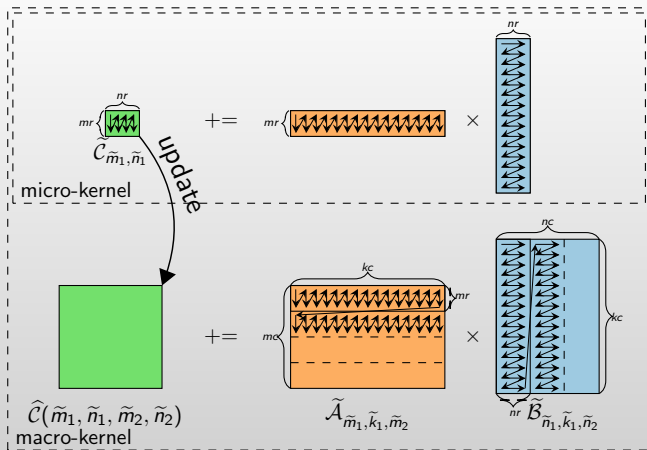


(b) Double-Precision.

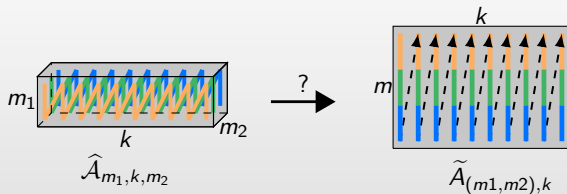


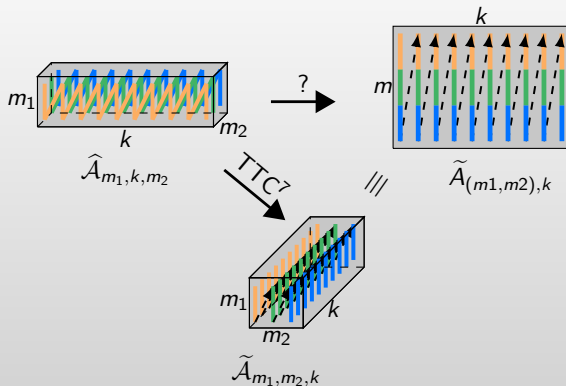


- Blocking for L3, L2, L1 cache as well as registers



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

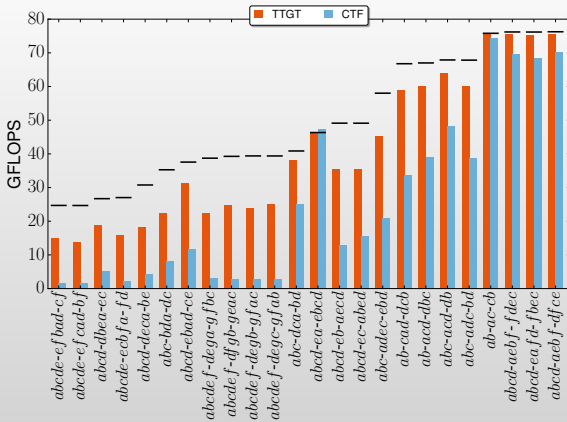


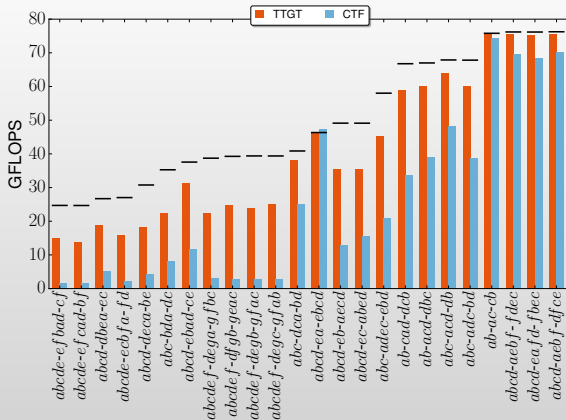


- Preserve stride-1 index
 ⇒ Efficient packing routines

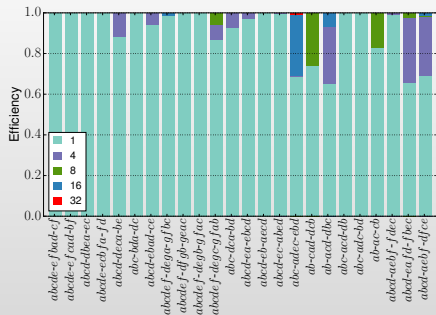
⁷Paul Springer et al. "TTC: A high-performance Compiler for Tensor Transpositions"

- 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

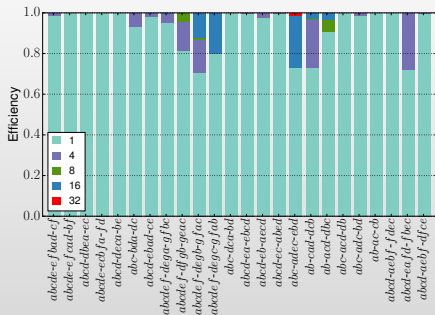




- TTGT good in compute-bound regime
- TTGT bad in bandwidth-bound regime
- TTGT faster than CTF everywhere.



(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%