

Tensor Contractions with Extended BLAS Kernels on CPU and GPU

Yang Shi

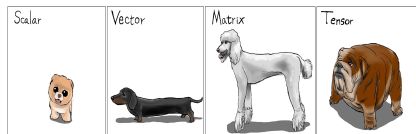
Electrical Engineering and Computer Science
University of California, Irvine

Joint work with U.N. Niranjan, Animashree Anandkumar and Cris Cecka

Southern California Machine Learning Symposium

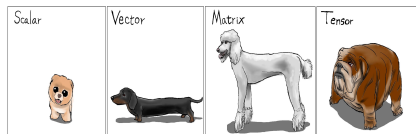
November 18, 2016

Tensor Contraction-Motivation



	2		
2	3	4 7 8 1	1 4 2 5
	5	2 5 1 6	2 3 0 4 9 2 0
	9	3 3 9 8	1 7 0 8 2 9 3
Scalar	Vector	Matrix	Tensor

Tensor Contraction-Motivation

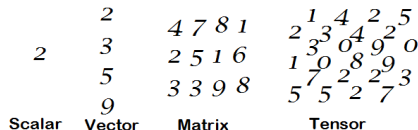
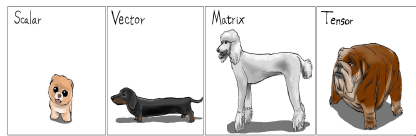


	2			
	3	4 7 8 1		
2	5	2 5 1 6		
	9	3 3 9 8		
Scalar	Vector	Matrix		
			1 4 2 5	
			2 3 0 4 2 5	
			1 3 0 4 9 2 0	
			7 0 8 2 9 3	
			5 5 2 2 7 3	
			Tensor	

Why we need tensor?

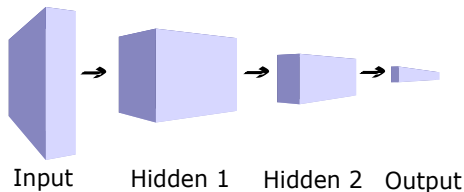
- 1 Analysis of high dimension data

Tensor Contraction-Motivation

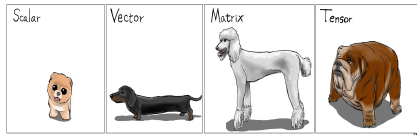


Why we need tensor?

- 1 Analysis of high dimension data



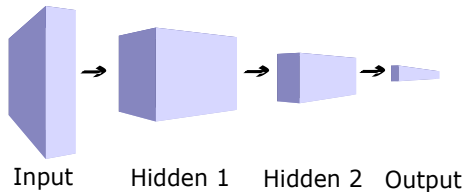
Tensor Contraction-Motivation



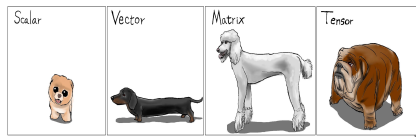
	2	4 7 8 1	¹ 4 ² ⁵ 2 ³ 3 ⁰ 4 ⁹ 0
2	3	2 5 1 6	1 ⁰ 3 ⁸ 2 ⁹ 0
	5	3 3 9 8	5 ⁷ 5 ² 2 ⁹ 3
	9		
Scalar	Vector	Matrix	Tensor

Why we need tensor?

- 1 Analysis of high dimension data
- 2 Multi-dimension relationship of low dimension data



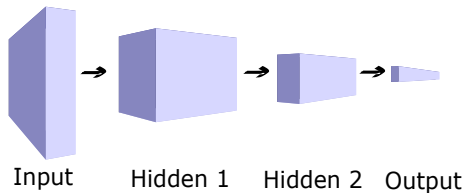
Tensor Contraction-Motivation



	2		
2	3	4 7 8 1	1 4 2 5
	5	2 5 1 6	2 3 0 4 2 5
	9	3 3 9 8	1 3 0 8 9 0
Scalar	Vector	Matrix	Tensor

Why we need tensor?

- 1 Analysis of high dimension data
- 2 Multi-dimension relationship of low dimension data



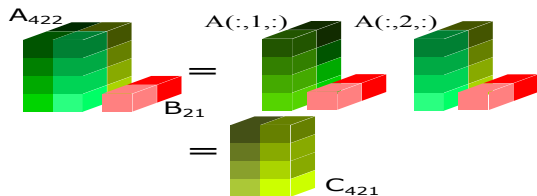
$$E(x_1 \otimes x_2) = \text{[diagram of a 2D tensor slice]} + \dots + \text{[diagram of a 2D tensor slice]}$$

$$E(x_1 \otimes x_2 \otimes x_3) = \text{[diagram of a 3D tensor slice]} + \dots + \text{[diagram of a 3D tensor slice]}$$

Tensor Contraction-Motivation

What is tensor contraction

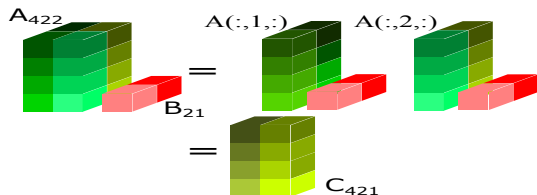
$$C_C = \alpha A_A B_B + \beta C_C$$



Tensor Contraction-Motivation

What is tensor contraction

$$C_C = \alpha A_A B_B + \beta C_C$$



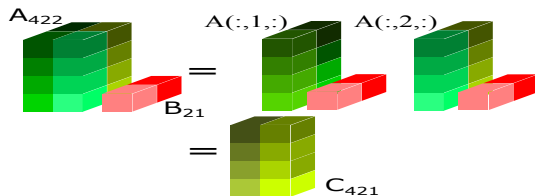
Why we need tensor contraction?

¹Picture from M. Alex O. Vasilescu

Tensor Contraction-Motivation

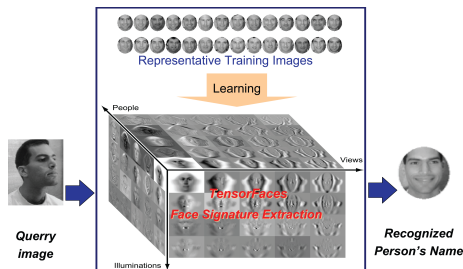
What is tensor contraction

$$C_C = \alpha A_A B_B + \beta C_C$$



Why we need tensor contraction?

- 1 Machine learning: Tensor decomposition



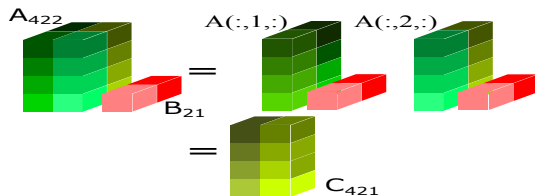
Tensor faces¹

¹Picture from M. Alex O. Vasilescu

Tensor Contraction-Motivation

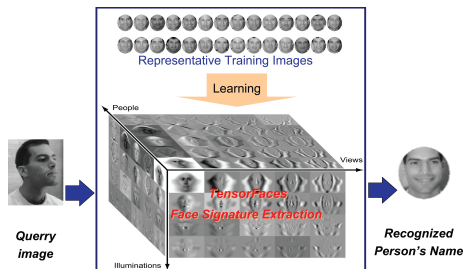
What is tensor contraction

$$C_C = \alpha A_A B_B + \beta C_C$$



Why we need tensor contraction?

- 1 Machine learning: Tensor decomposition
- 2 Physics
- 3 Chemistry



Tensor faces¹

¹Picture from M. Alex O. Vasilescu

Tensor Contraction-Motivation

What do we have?

What do we have?

Tensor computation libraries

- 1 Arbitrary/restricted tensor operation of any order and dimension
- 2 Such as: Matlab Tensortoolbox, BTAS, FTensor, Cyclops

What do we have?

Tensor computation libraries

- ① Arbitrary/restricted tensor operation of any order and dimension
- ② Such as: Matlab Tensortoolbox, BTAS, FTensor, Cyclops

Efficient computing frame

- ① Static analysis solutions: loop reorganization/fusion
- ② Parallel and distributed computing system:
BatchedGEMM functions in MKL 11.3 β , CuBLAS v4.1

Tensor Contraction-Motivation

What are the limitations?

Tensor Contraction-Motivation

What are the limitations?

- 1 Explicit permutation takes long time in current tensor libraries:

Tensor Contraction-Motivation

What are the limitations?

- ① Explicit permutation takes long time in current tensor libraries:

Consider $C_{mnp} = A_{km}B_{pkn}$.

- $A_{km} \rightarrow A_{mk}$.
- $B_{pkn} \rightarrow B_{kpn}$.
- $C_{mnp} \rightarrow C_{mpn}$.
- $C_{mpn} = A_{mk}B_{kpn}$.
- $C_{mpn} \rightarrow C_{mnp}$.

Tensor Contraction-Motivation

What are the limitations?

- ① Explicit permutation takes long time in current tensor libraries:

Consider $C_{mnp} = A_{km}B_{pkn}$.

- $A_{km} \rightarrow A_{mk}$.
- $B_{pkn} \rightarrow B_{kpn}$.
- $C_{mnp} \rightarrow C_{mpn}$.
- $C_{mpn} = A_{mk}B_{kpn}$.
- $C_{mpn} \rightarrow C_{mnp}$.

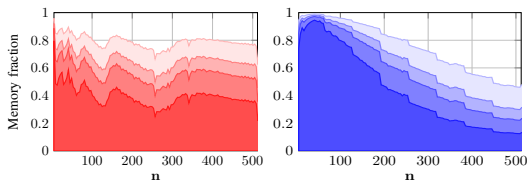


Figure: The fraction of time spent in copies/transpositions when computing $C_{mnp} = A_{mk}B_{pkn}$. Lines are shown with 1, 2, 3, and 6 total transpositions performed on either the input or output. (Left) CPU. (Right) GPU.

- **Proposed tensor operation kernel: StridedBatchedGEMM**

- **Proposed tensor operation kernel: StridedBatchedGEMM**
 - ① Library-based approaches that avoid memory movement
 - ② Constant-strided BatchedGEMM that has more optimization opportunities

- **Proposed tensor operation kernel: StridedBatchedGEMM**
 - ① Library-based approaches that avoid memory movement
 - ② Constant-strided BatchedGEMM that has more optimization opportunities
- **Provided evaluation strategies for tensor contractions**

- **Proposed tensor operation kernel: StridedBatchedGEMM**
 - ① Library-based approaches that avoid memory movement
 - ② Constant-strided BatchedGEMM that has more optimization opportunities
- **Provided evaluation strategies for tensor contractions**
- **Applied to tensor decomposition**

BLAS Operators

BLAS(Basic Linear Algebra Subprograms): Low-level routines for performing common linear algebra operations.

Level 1	Level 2	Level 3
$y \leftarrow \alpha x + y$	$y \leftarrow \alpha op(A)x + \beta y$	$C \leftarrow \alpha op(A)op(B) + \beta C$

Table: BLAS operators

BLAS Operators

BLAS(Basic Linear Algebra Subprograms): Low-level routines for performing common linear algebra operations.

Level 1	Level 2	Level 3
$y \leftarrow \alpha x + y$	$y \leftarrow \alpha op(A)x + \beta y$	$C \leftarrow \alpha op(A)op(B) + \beta C$

Table: BLAS operators

Example:

GEMM(ORDER, TRANSA, TRANSB, M, N, K, α , A, LDA, B, LDB, β , C, LDC)

BLAS Operators

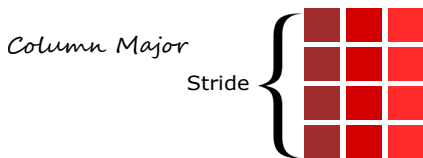
BLAS(Basic Linear Algebra Subprograms): Low-level routines for performing common linear algebra operations.

Level 1	Level 2	Level 3
$y \leftarrow \alpha x + y$	$y \leftarrow \alpha op(A)x + \beta y$	$C \leftarrow \alpha op(A)op(B) + \beta C$

Table: BLAS operators

Example:

GEMM(ORDER, TRANSA, TRANSB, M, N, K, α , A, LDA, B, LDB, β , C, LDC)



BLAS Operators

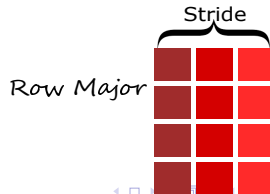
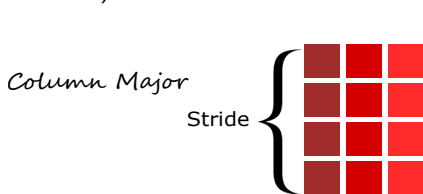
BLAS(Basic Linear Algebra Subprograms): Low-level routines for performing common linear algebra operations.

Level 1	Level 2	Level 3
$y \leftarrow \alpha x + y$	$y \leftarrow \alpha op(A)x + \beta y$	$C \leftarrow \alpha op(A)op(B) + \beta C$

Table: BLAS operators

Example:

GEMM(ORDER, TRANSA, TRANSB, M, N, K, α , A, LDA, B, LDB, β , C, LDC)



Tensor Contraction-Extended BLAS Operators

Extended BLAS Kernel for tensor one-index contraction

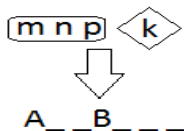
$$C = \alpha op(A)op(B) + \beta C$$

Tensor Contraction-Extended BLAS Operators

Extended BLAS Kernel for tensor one-index contraction

$$C = \alpha op(A)op(B) + \beta C$$

$$C_{mnp} = A_{**} \times B_{***}$$



If fixing C, there are total $3 \times 2 \times 3 \times 2 \times 1 = 36$ cases.

Tensor Contraction-Extended BLAS Operators

Case	Contraction	Kernel1	Kernel2	Kernel3
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	$C_{m[n]p} = A_{mk}B_{k[n]p}$

Table: Case 1.1 single index contraction operations between a second-order tensor and a third-order tensor and possible mappings to Level 3 BLAS routines

Tensor Contraction-Extended BLAS Operators

Case	Contraction	Kernel1	Kernel2	Kernel3
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	$C_{m[n]p} = A_{mk}B_{k[n]p}$

Table: Case 1.1 single index contraction operations between a second-order tensor and a third-order tensor and possible mappings to Level 3 BLAS routines

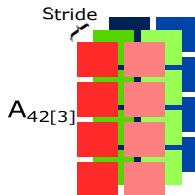
StridedBatchedGEMM(Order, TRANSA, TRANSB, M, N, K, α , A, LDA, LOA, B, LDB, LOB, β , C, LDC, LOC, P)

Tensor Contraction-Extended BLAS Operators

Case	Contraction	Kernel1	Kernel2	Kernel3
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	$C_{m[n]p} = A_{mk}B_{k[n]p}$

Table: Case 1.1 single index contraction operations between a second-order tensor and a third-order tensor and possible mappings to Level 3 BLAS routines

StridedBatchedGEMM(Order, TRANSA, TRANSB, M, N, K, α , A, LDA, LOA, B, LDB, LOB, β , C, LDC, LOC, P)

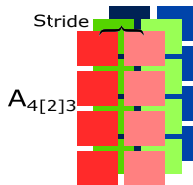
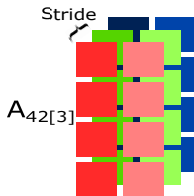


Tensor Contraction-Extended BLAS Operators

Case	Contraction	Kernel1	Kernel2	Kernel3
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	$C_{m[n]p} = A_{mk}B_{k[n]p}$

Table: Case 1.1 single index contraction operations between a second-order tensor and a third-order tensor and possible mappings to Level 3 BLAS routines

StridedBatchedGEMM(Order, TRANSA, TRANSB, M, N, K, α , A, LDA, LOA, B, LDB, LOB, β , C, LDC, LOC, P)



Flatten V.S. SBGEMM

Batching in last mode V.S. Batching in earlier mode

Mixed mode batching on input/output tensors

Flatten V.S. SBGEMM

- ① A single large GEMM is more efficient
- ② Flatten modes whenever possible

Batching in last mode V.S. Batching in earlier mode

Mixed mode batching on input/output tensors

Flatten V.S. SBGEMM

- ① A single large GEMM is more efficient
- ② Flatten modes whenever possible

Batching in last mode V.S. Batching in earlier mode

Mixed mode batching on input/output tensors

Flatten V.S. SBGEMM

- ① A single large GEMM is more efficient
- ② Flatten modes whenever possible

Batching in last mode V.S. Batching in earlier mode

- ① On CPU: batching in the third mode is advantageous when the size of the tensor is small
- ② On GPU: there is no discernible preference in the choice of batching mode

Mixed mode batching on input/output tensors

Flatten V.S. SBGEMM

- ① A single large GEMM is more efficient
- ② Flatten modes whenever possible

Batching in last mode V.S. Batching in earlier mode

- ① On CPU: batching in the third mode is advantageous when the size of the tensor is small
- ② On GPU: there is no discernible preference in the choice of batching mode

Mixed mode batching on input/output tensors

Flatten V.S. SBGEMM

- ① A single large GEMM is more efficient
- ② Flatten modes whenever possible

Batching in last mode V.S. Batching in earlier mode

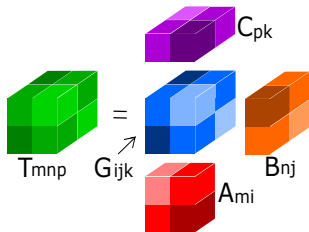
- ① On CPU: batching in the third mode is advantageous when the size of the tensor is small
- ② On GPU: there is no discernible preference in the choice of batching mode

Mixed mode batching on input/output tensors

- ① On CPU: mode of the output tensor is more important than the batching mode of the input tensor.

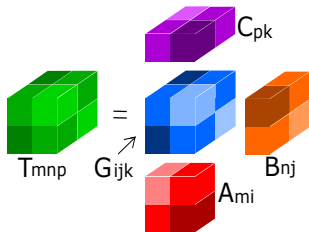
Applications: Tucker Decomposition

$$T_{mnp} = G_{ijk} A_{mi} B_{nj} C_{pk}$$



Applications: Tucker Decomposition

$$T_{mnp} = G_{ijk} A_{mi} B_{nj} C_{pk}$$



Main steps in the algorithm

- $Y_{mjk} = T_{mnp} B_{nj}^t C_{pk}^t$
- $Y_{ink} = T_{mnp} A_{mi}^{t+1} C_{pk}^t$
- $Y_{ijp} = T_{mnp} B_{nj}^{t+1} A_{mi}^{t+1}$

Applications: Tucker Decomposition

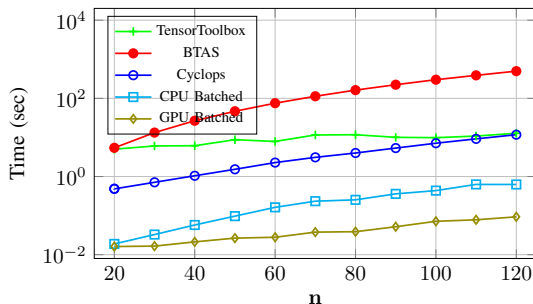


Figure: Performance on Tucker decomposition.

Conclusion

- StridedBatchedGEMM for generalized tensor contractions.
- Avoid explicit transpositions or permutations.
- 10x(GPU) and 2x(CPU) speedup on small and moderate sized tensors.
- Available in CuBLAS 8.0

Conclusion

- StridedBatchedGEMM for generalized tensor contractions.
- Avoid explicit transpositions or permutations.
- 10x(GPU) and 2x(CPU) speedup on small and moderate sized tensors.
- Available in CuBLAS 8.0
- Future works: Multi-index contractions or sparse tensor algebra.

Thank you!

More information: <http://shiyangdaisy1.wixsite.com/ucieecs>