# Tensor Contraction with Extended BLAS Kernels on CPU and GPU
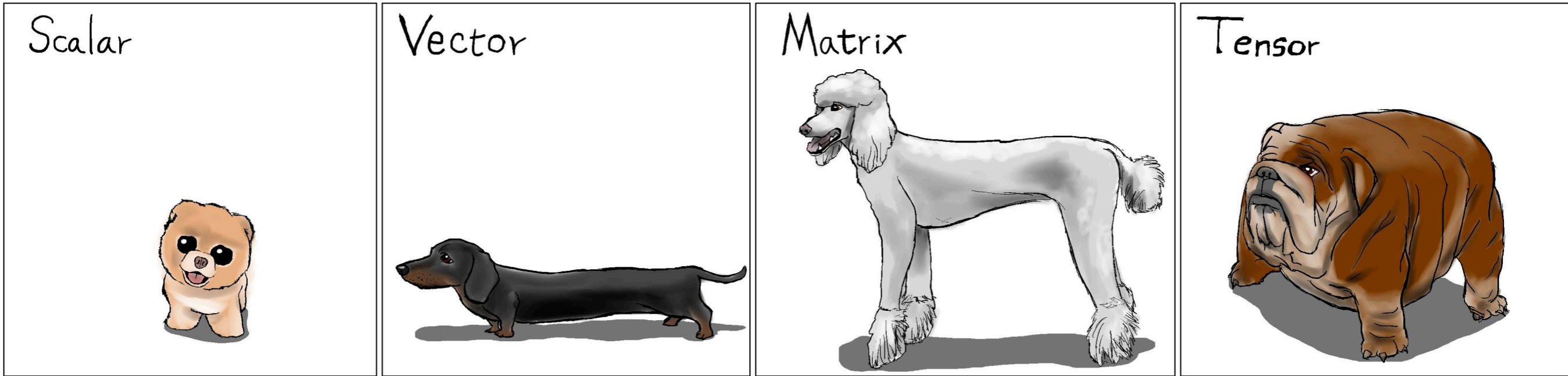
Yang Shi
University of California, Irvine,  EECS

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

**SIAM-ALA18**

# Tensor Contraction-Motivation



| Scalar | Vector | Matrix | Tensor |

2

2  
3  
5  
9

4 7 8 1  
2 5 1 6  
3 3 9 8

1 4 2 5  
2 3 3 0 4 9 2 0  
1 0 8 2 9  
7 2 2 9 3  
5 5 2 7

**Scalar    Vector    Matrix    Tensor**

# Tensor Contraction-Motivation

**Why we need tensor?**

# Tensor Contraction-Motivation
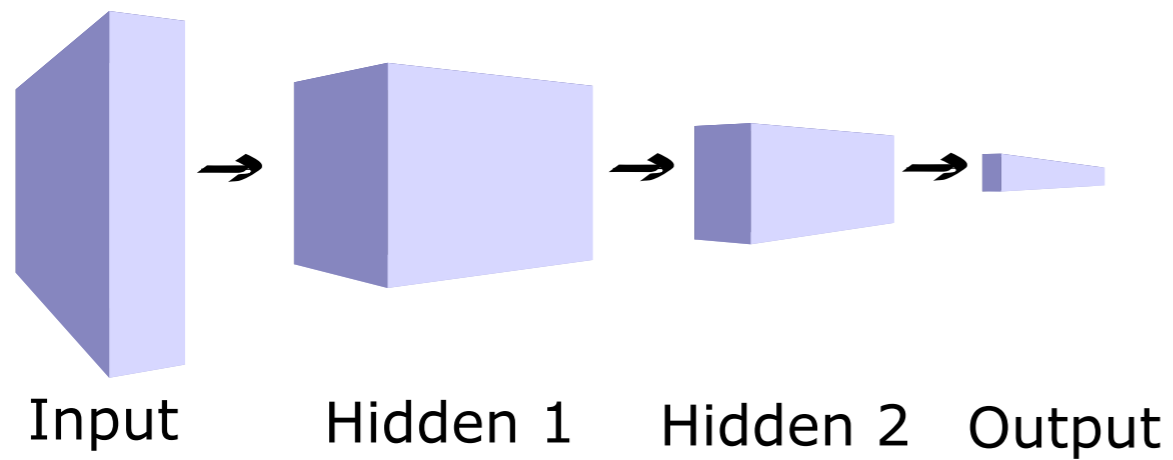
## Why we need tensor?

Modern data is inherently multi-dimensional

# Tensor Contraction-Motivation

## Why we need tensor?

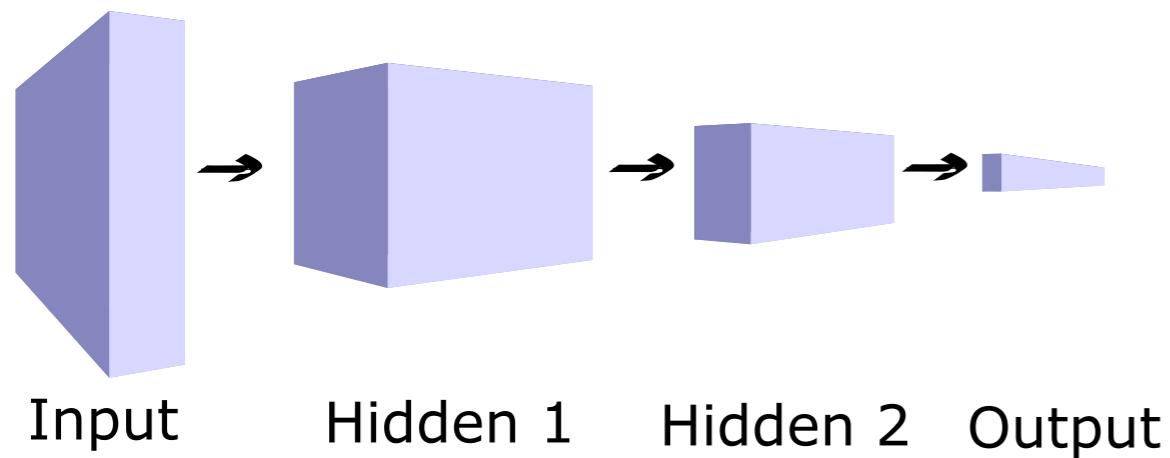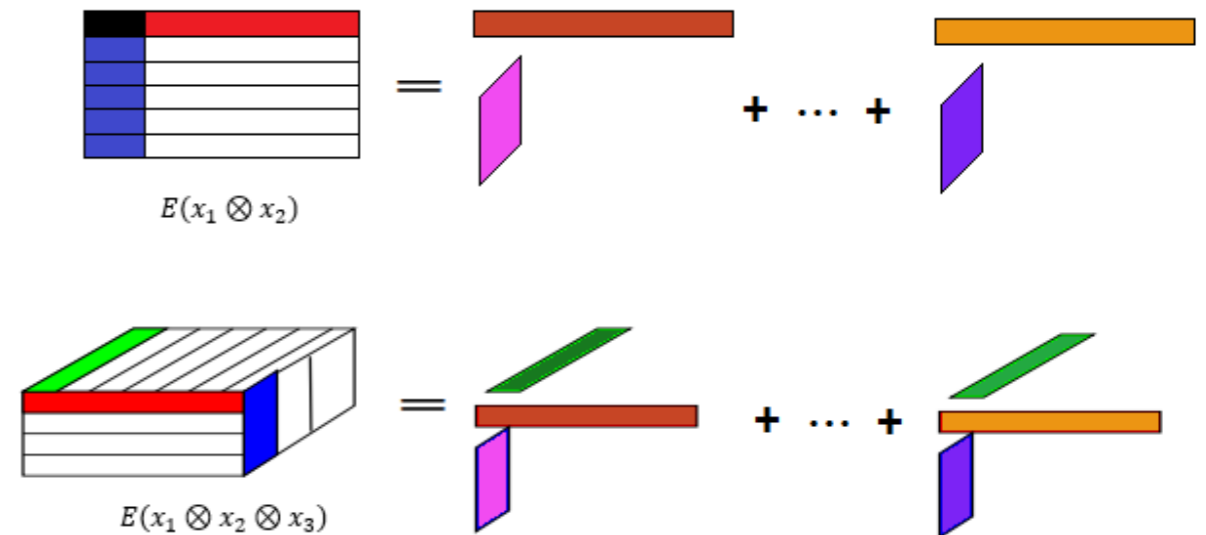Modern data is inherently multi-dimensional

**Neural Networks**



Input     Hidden 1     Hidden 2    Output

# Tensor Contraction-Motivation

## Why we need tensor?

Modern data is inherently multi-dimensional

**Neural Networks**



Input  Hidden 1  Hidden 2  Output

**Method of Moment**



$E(x_1 \otimes x_2)$

$E(x_1 \otimes x_2 \otimes x_3)$

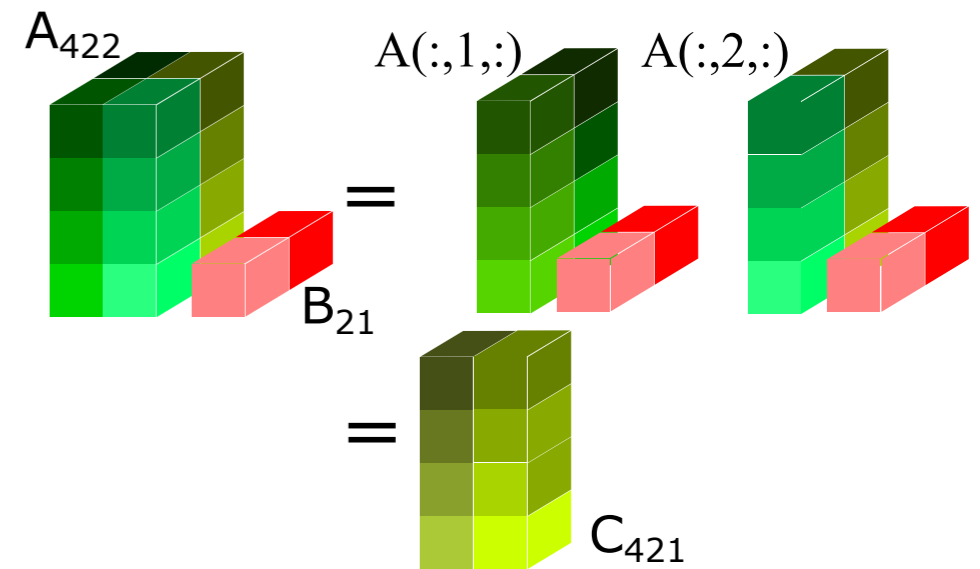# Tensor Contraction-Motivation

**What is tensor contraction?**

# Tensor Contraction-Motivation

**What is tensor contraction?**

$$C_{\mathcal{C}} = A_{\mathcal{A}}\, B_{\mathcal{B}}$$

# Tensor Contraction-Motivation

**What is tensor contraction?**

$$C_{\mathcal{C}} = A_{\mathcal{A}} \, B_{\mathcal{B}}$$



**Why do we need tensor contraction?**

# Tensor Contraction-Motivation

## What is tensor contraction?

$$C_{\mathcal{C}} = A_{\mathcal{A}} \, B_{\mathcal{B}}$$



## Why do we need tensor contraction?

- Physics
- Chemistry

# Tensor Contraction-Motivation

**Why do we need tensor contraction?**

# Tensor Contraction-Motivation

## Why do we need tensor contraction?

- **Deep Learning**

# Tensor Contraction-Motivation

## Why do we need tensor contraction?

- **Deep Learning**



- **Learning latent variable model with tensor decomposition**
  Example: Topic modeling

# Tensor Contraction-Motivation

## Why do we need tensor contraction?

- **Deep Learning**



- **Learning latent variable model with tensor decomposition**
  Example: Topic modeling

# Tensor Contraction-Motivation

## Why do we need tensor contraction?

- **Deep Learning**



- **Learning latent variable model with tensor decomposition**

  Example: Topic modeling

  h: Proportion of topics in a document

  $h = i$ with prob. $w_i$

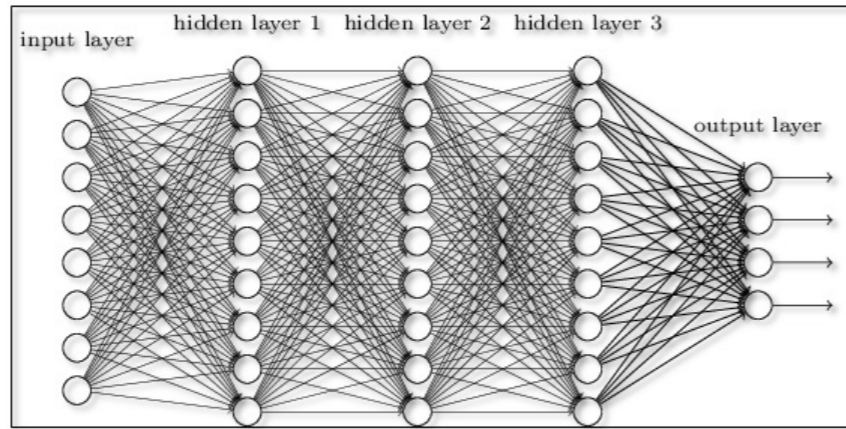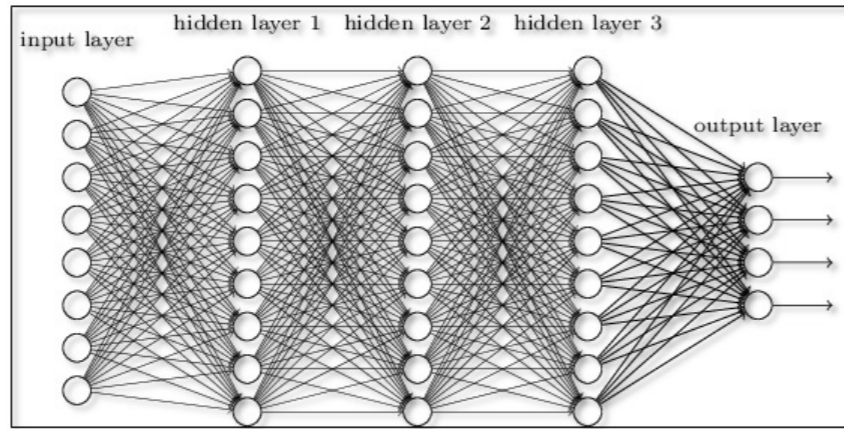  A: Topic-word matrix

  $A(i,j) = \mathcal{P}(x_m = i | y_m = j)$

# Tensor Contraction-Motivation

## Why do we need tensor contraction?

- **Deep Learning**



- **Learning latent variable model with tensor decomposition**

  Example: Topic modeling

  h: Proportion of topics in a document

  $$h = i \text{ with prob. } w_i$$

  A: Topic-word matrix

  $$A(i,j) = \mathcal{P}(x_m = i | y_m = j)$$

  Third order moment:

  $$M_3 = \mathbb{E}(x \otimes x \otimes x) = \sum_i w_i a_i \otimes a_i \otimes a_i$$

# Tensor Contraction-Motivation

**What do we have?**

# Tensor Contraction-Motivation

## What do we have?

**Tensor computation libraries:**

• Arbitrary/restricted tensor operations of any order and dimension

• Such as: Matlab Tensortoolbox, BTAS, FTensor, Cyclops

# Tensor Contraction-Motivation

## What do we have?

**Tensor computation libraries:**

• Arbitrary/restricted tensor operations of any order and dimension

• Such as: Matlab Tensortoolbox, BTAS, FTensor, Cyclops

**Efficient computing frame:**

• Static analysis solutions: loop reorganization, fusion

• Parallel and distributed computing system: BatchedGEMM functions in MKL 11.3, CuBLAS v4.1: compute many matrix-matrix multiplies at once.

# Tensor Contraction-Motivation

## What do we have?

**Tensor computation libraries:**

• Arbitrary/restricted tensor operations of any order and dimension

• Such as: Matlab Tensortoolbox, BTAS, FTensor, Cyclops

**Efficient computing frame:**

• Static analysis solutions: loop reorganization, fusion

• Parallel and distributed computing system: BatchedGEMM functions in MKL 11.3, CuBLAS v4.1: compute many matrix-matrix multiplies at once.

# Tensor Contraction-Motivation

**What are the limitations?**

# Tensor Contraction-Motivation

## What are the limitations?

- Explicit permutation takes long time in current tensor libraries:

# Tensor Contraction-Motivation

## What are the limitations?

- Explicit permutation takes long time in current tensor libraries:

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

# Tensor Contraction-Motivation

## What are the limitations?

- Explicit permutation takes long time in current tensor libraries:

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

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

# Tensor Contraction-Motivation

## What are the limitations?

- Explicit permutation takes long time in current tensor libraries:

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

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



Figure: The fraction of time spent in copies/transpositions when computing Cmnp = AmkBpkn . Lines are shown with 1, 2, 3, and 6 total transpositions performed on either the input or output. (Left) CPU. (Right) GPU.

# Overview

# Overview

- Propose tensor operation kernel: StridedBatchedGEMM

# Overview

- Propose tensor operation kernel: StridedBatchedGEMM
  - Library-based approaches that avoid memory movement

# Overview

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

# Overview

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

- Provide evaluation strategies for tensor contractions

# Overview

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

- Provide evaluation strategies for tensor contractions

- Apply to tensor decomposition

# Overview

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

- Provide evaluation strategies for tensor contractions

- Apply to tensor decomposition

- Introduce TensorLy: Tensor learning in python

# BLAS Operations

# BLAS Operations

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

# BLAS Operations

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

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

# BLAS Operations

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

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

Example:
GEMM(ORDER, TRANSA, TRANSB, M, N, K, $\alpha$, A, LDA, B, LDB, $\beta$, C, LDC)

# BLAS Operations

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

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

Example:
GEMM(ORDER, TRANSA, TRANSB, M, N, K, $\alpha$, A, LDA, B, LDB, $\beta$, C, LDC)

# Extended BLAS Operator

# Extended BLAS Operator

**Focusing: one-index contraction**

> ### Extended BLAS Kernel for tensor one-index contraction
>
> $$C = \alpha\, op(A)\, op(B) + \beta C$$

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

$$\boxed{\text{m n p}} \;\langle k \rangle$$

$$\Downarrow$$

$$\text{A\_\_B\_\_\_\_}$$

# Extended BLAS Operator

**Focusing: one-index contraction**

> **Extended BLAS Kernel for tensor one-index contraction**
>
> $$C = \alpha\, op(A)\, op(B) + \beta C$$

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

$$\boxed{m\ n\ p}\ \langle k \rangle$$

$$\Downarrow$$

$$A_{\_\_}B_{\_\_\_}$$

If fixing indices of C, there are total 3 x 2 x 3  x 2  x 1 = 36 cases.

# Extended BLAS Operator

# Extended BLAS Operator

| 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: Example: possible mappings to Level 3 BLAS routines

# Extended BLAS Operator

| 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: Example: possible mappings to Level 3 BLAS routines

StridedBatchedGEMM(ORDER, TRANSA, TRANSB, M, N, K, $\alpha$, A, LDA, LOA, B, LDB, LOB, $\beta$, C, LDC, LOC, P)

# Extended BLAS Operator

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

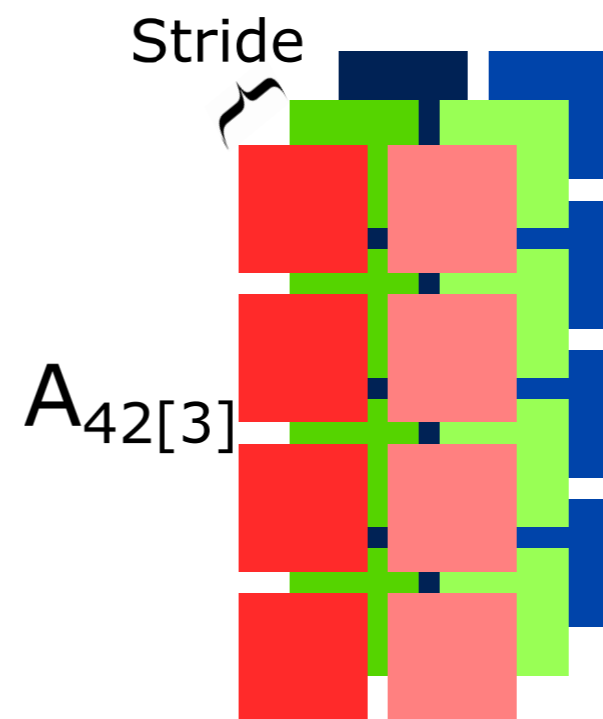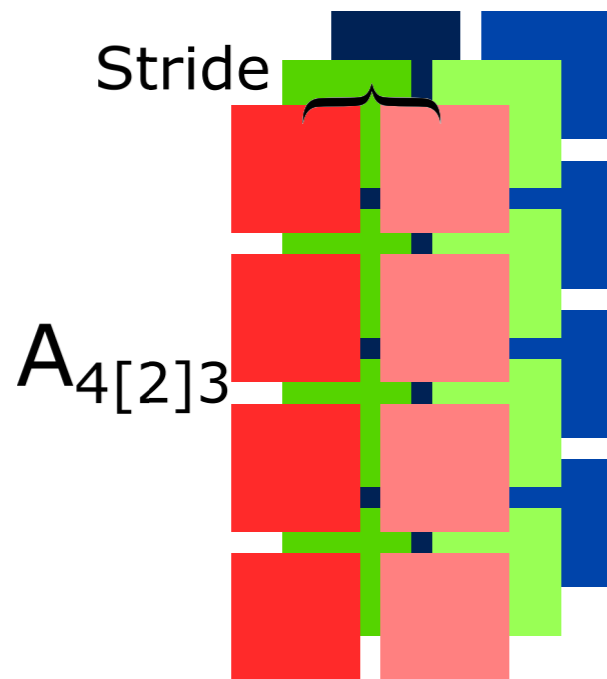Table: Example: possible mappings to Level 3 BLAS routines

StridedBatchedGEMM(ORDER, TRANSA, TRANSB, M, N, K, $\alpha$, A, LDA, LOA, B, LDB, LOB, $\beta$, C, LDC, LOC, P)

# Example

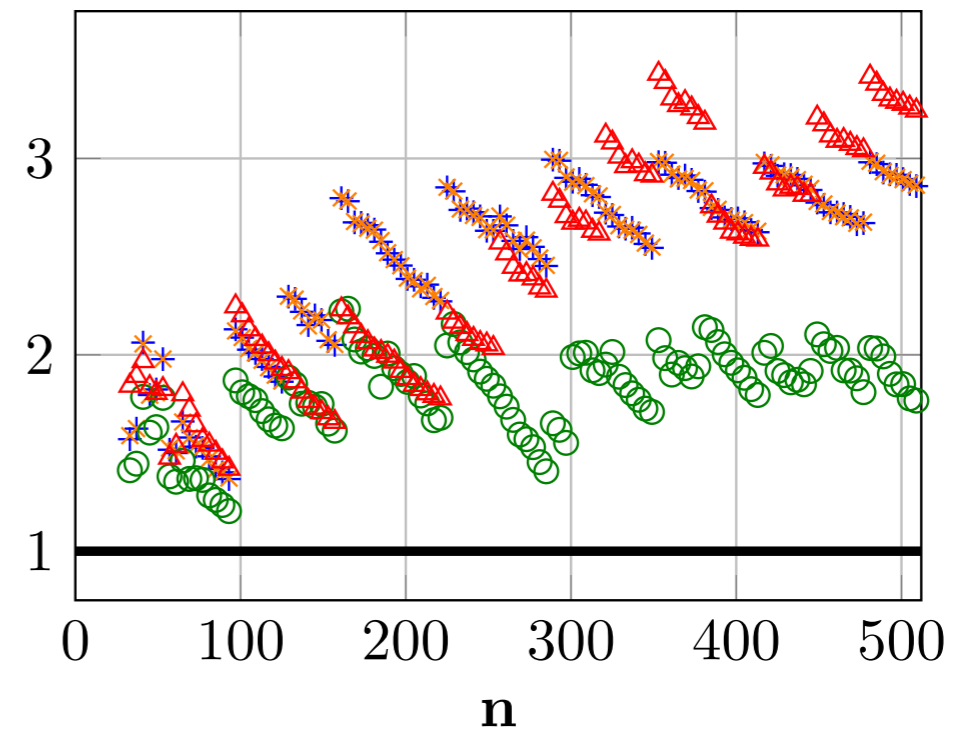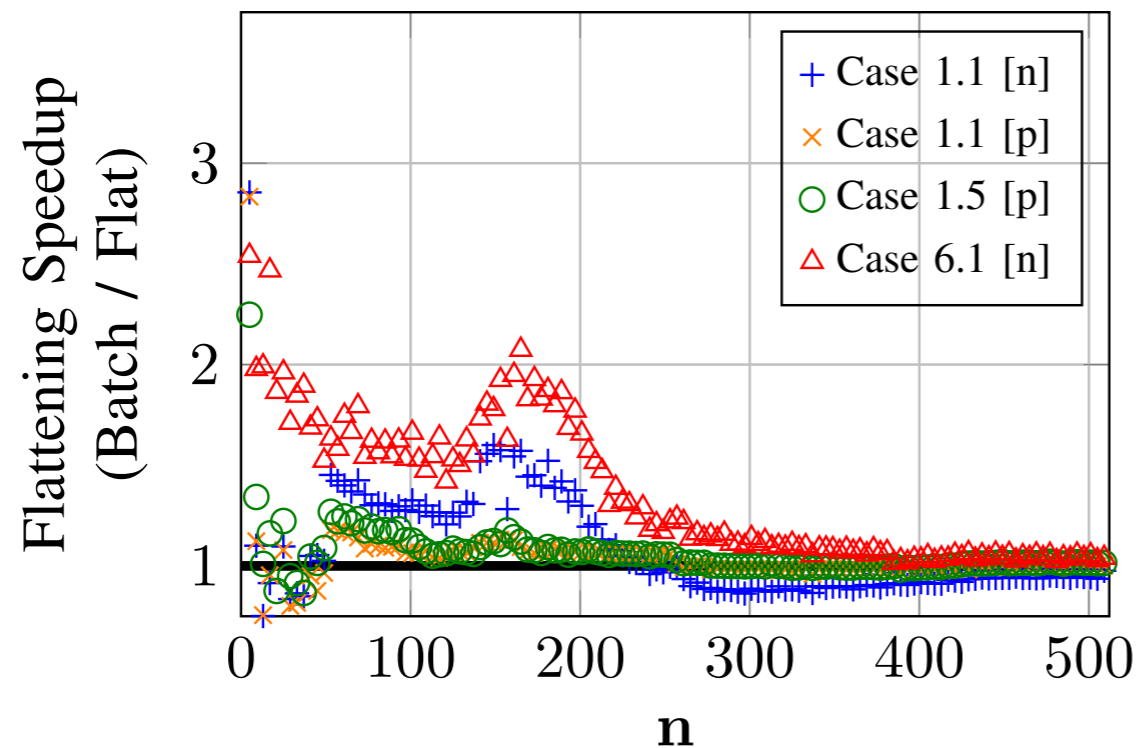| Case | Contraction | Kernel1 | Kernel2 | Kernel3 | Case | Contraction | Kernel1 | Kernel2 |
|---|---|---|---|---|---|---|---|---|
| 1.1 | $A_{mk}B_{knp}$ | $C_{m(np)} = A_{mk}B_{k(np)}$ | $C_{mn[p]} = A_{mk}B_{kn[p]}$ | $C_{m[n]p} = A_{mk}B_{k[n]p}$ | 4.1 | $A_{kn}B_{kmp}$ | $C_{mn[p]} = B_{km[p]}^\top A_{kn}$ | |
| 1.2 | $A_{mk}B_{kpn}$ | $C_{mn[p]} = A_{mk}B_{k[p]n}$ | $C_{m[n]p} = A_{mk}B_{kp[n]}$ | | 4.2 | $A_{kn}B_{kpm}$ | $C_{mn[p]} = B_{k[p]m}^\top A_{kn}$ | |
| 1.3 | $A_{mk}B_{nkp}$ | $C_{mn[p]} = A_{mk}B_{nk[p]}^\top$ | | | 4.3 | $A_{kn}B_{mkp}$ | $C_{mn[p]} = B_{mk[p]} A_{kn}$ | |
| 1.4 | $A_{mk}B_{pkn}$ | $C_{m[n]p} = A_{mk}B_{pk[n]}^\top$ | | | 4.4 | $A_{kn}B_{pkm}$ | $TRANS(A_{kn}^\top B_{pk[m]}^\top)$ | $C_{[m][n]p} = B_{pk[m]}A_{k[n]}$ |
| 1.5 | $A_{mk}B_{npk}$ | $C_{m(np)} = A_{mk}B_{(np)k}^\top$ | $C_{mn[p]} = A_{mk}B_{n[p]k}^\top$ | | 4.5 | $A_{kn}B_{mpk}$ | $C_{mn[p]} = B_{m[p]k} A_{kn}$ | |
| 1.6 | $A_{mk}B_{pnk}$ | $C_{m[n]p} = A_{mk}B_{p[n]k}^\top$ | | | 4.6 | $A_{kn}B_{pmk}$ | $TRANS(A_{kn}^\top B_{p[m]k}^\top)$ | $C_{[m][n]p} = B_{p[m]k}A_{k[n]}$ |
| 2.1 | $A_{km}B_{knp}$ | $C_{m(np)} = A_{km}^\top B_{k(np)}$ | $C_{mn[p]} = A_{km}^\top B_{kn[p]}$ | $C_{m[n]p} = A_{km}^\top B_{k[n]p}$ | 5.1 | $A_{pk}B_{kmn}$ | $C_{(mn)p} = B_{k(mn)}^\top A_{pk}^\top$ | $C_{m[n]p} = B_{km[n]}^\top A_{pk}^\top$ |
| 2.2 | $A_{km}B_{kpn}$ | $C_{mn[p]} = A_{km}^\top B_{k[p]n}$ | $C_{m[n]p} = A_{km}^\top B_{kp[n]}$ | | 5.2 | $A_{pk}B_{knm}$ | $C_{m[n]p} = B_{k[n]m}^\top A_{pk}^\top$ | |
| 2.3 | $A_{km}B_{nkp}$ | $C_{mn[p]} = A_{km}^\top B_{nk[p]}^\top$ | | | 5.3 | $A_{pk}B_{mkn}$ | $C_{m[n]p} = B_{mk[n]} A_{pk}^\top$ | |
| 2.4 | $A_{km}B_{pkn}$ | $C_{m[n]p} = A_{km}^\top B_{pk[n]}^\top$ | | | 5.4 | $A_{pk}B_{nkm}$ | $TRANS(B_{nk[m]} A_{pk}^\top)$ | $C_{[m]n[p]} = B_{nk[m]}A_{[p]k}$ |
| 2.5 | $A_{km}B_{npk}$ | $C_{m(np)} = A_{km}^\top B_{(np)k}^\top$ | $C_{mn[p]} = A_{km}^\top B_{n[p]k}^\top$ | | 5.5 | $A_{pk}B_{mnk}$ | $C_{(mn)p} = B_{(mn)k}A_{pk}^\top$ | $C_{m[n]p} = B_{m[n]k}A_{pk}^\top$ |
| 2.6 | $A_{km}B_{pnk}$ | $C_{m[n]p} = A_{km}^\top B_{p[n]k}^\top$ | | | 5.6 | $A_{pk}B_{nmk}$ | $TRANS(B_{n[m]k} A_{pk}^\top)$ | $C_{[m]n[p]} = B_{n[m]k}A_{[p]k}$ |
| 3.1 | $A_{nk}B_{kmp}$ | $C_{mn[p]} = B_{km[p]}^\top A_{nk}^\top$ | | | 6.1 | $A_{kp}B_{kmn}$ | $C_{(mn)p} = B_{k(mn)}^\top A_{kp}$ | $C_{m[n]p} = B_{km[n]}^\top A_{kp}$ |
| 3.2 | $A_{nk}B_{kpm}$ | $C_{mn[p]} = B_{k[p]m}^\top A_{nk}^\top$ | | | 6.2 | $A_{kp}B_{knm}$ | $C_{m[n]p} = B_{k[n]m}^\top A_{kp}$ | |
| 3.3 | $A_{nk}B_{mkp}$ | $C_{mn[p]} = B_{mk[p]} A_{nk}^\top$ | | | 6.3 | $A_{kp}B_{mkn}$ | $C_{m[n]p} = B_{mk[n]} A_{kp}$ | |
| 3.4 | $A_{nk}B_{pkm}$ | $TRANS(A_{nk}B_{pk[m]}^\top)$ | $C_{[m][n]p} = B_{pk[m]}A_{[n]k}$ | | 6.4 | $A_{kp}B_{nkm}$ | $TRANS(B_{nk[m]} A_{kp})$ | $C_{[m]n[p]} = B_{nk[m]}A_{k[p]}$ |
| 3.5 | $A_{nk}B_{mpk}$ | $C_{mn[p]} = B_{m[p]k} A_{nk}^\top$ | | | 6.5 | $A_{kp}B_{mnk}$ | $C_{(mn)p} = B_{(mn)k}A_{kp}$ | $C_{m[n]p} = B_{m[n]k}A_{kp}$ |
| 3.6 | $A_{nk}B_{pmk}$ | $TRANS(A_{nk}B_{p[m]k}^\top)$ | $C_{[m][n]p} = B_{p[m]k}A_{[n]k}$ | | 6.6 | $A_{kp}B_{nmk}$ | $TRANS(B_{n[m]k} A_{kp})$ | $C_{[m]n[p]} = B_{n[m]k}A_{k[p]}$ |

# Example

| Case | Contraction | Kernel1 | Kernel2 | Kernel3 |
|---|---|---|---|---|
| 1.1 | $A_{mk}B_{knp}$ | $C_{m(np)} = A_{mk}B_{k(np)}$ | $C_{mn[p]} = A_{mk}B_{kn[p]}$ | $C_{m[n]p} = A_{mk}B_{k[n]p}$ |
| 1.2 | $A_{mk}B_{kpn}$ | $C_{mn[p]} = A_{mk}B_{k[p]n}$ | $C_{m[n]p} = A_{mk}B_{kp[n]}$ | |
| 1.3 | $A_{mk}B_{nkp}$ | $C_{mn[p]} = A_{mk}B_{nk[p]}^\top$ | | |
| 1.4 | $A_{mk}B_{pkn}$ | $C_{m[n]p} = A_{mk}B_{pk[n]}^\top$ | | |
| 1.5 | $A_{mk}B_{npk}$ | $C_{m(np)} = A_{mk}B_{(np)k}^\top$ | $C_{mn[p]} = A_{mk}B_{n[p]k}^\top$ | |
| 1.6 | $A_{mk}B_{pnk}$ | $C_{m[n]p} = A_{mk}B_{p[n]k}^\top$ | | |
| 2.1 | $A_{km}B_{knp}$ | $C_{m(np)} = A_{km}^\top B_{k(np)}$ | $C_{mn[p]} = A_{km}^\top B_{kn[p]}$ | $C_{m[n]p} = A_{km}^\top B_{k[n]p}$ |
| 2.2 | $A_{km}B_{kpn}$ | $C_{mn[p]} = A_{km}^\top B_{k[p]n}$ | $C_{m[n]p} = A_{km}^\top B_{kp[n]}$ | |
| 2.3 | $A_{km}B_{nkp}$ | $C_{mn[p]} = A_{km}^\top B_{nk[p]}^\top$ | | |
| 2.4 | $A_{km}B_{pkn}$ | $C_{m[n]p} = A_{km}^\top B_{pk[n]}^\top$ | | |
| 2.5 | $A_{km}B_{npk}$ | $C_{m(np)} = A_{km}^\top B_{(np)k}^\top$ | $C_{mn[p]} = A_{km}^\top B_{n[p]k}^\top$ | |
| 2.6 | $A_{km}B_{pnk}$ | $C_{m[n]p} = A_{km}^\top B_{p[n]k}^\top$ | | |
| 3.1 | $A_{nk}B_{kmp}$ | $C_{mn[p]} = B_{km[p]}^\top A_{nk}^\top$ | | |
| 3.2 | $A_{nk}B_{kpm}$ | $C_{mn[p]} = B_{k[p]m}^\top A_{nk}^\top$ | | |
| 3.3 | $A_{nk}B_{mkp}$ | $C_{mn[p]} = B_{mk[p]} A_{nk}^\top$ | | |
| 3.4 | $A_{nk}B_{pkm}$ | $TRANS(A_{nk}B_{pk[m]}^\top)$ | $C_{[m][n]p} = B_{pk[m]}A_{[n]k}$ | |
| 3.5 | $A_{nk}B_{mpk}$ | $C_{mn[p]} = B_{m[p]k}A_{nk}^\top$ | | |
| 3.6 | $A_{nk}B_{pmk}$ | $TRANS(A_{nk}B_{p[m]k}^\top)$ | $C_{[m][n]p} = B_{p[m]k}A_{[n]k}$ | |

| Case | Contraction | Kernel1 | Kernel2 |
|---|---|---|---|
| 4.1 | $A_{kn}B_{kmp}$ | $C_{mn[p]} = B_{km[p]}^\top A_{kn}$ | |
| 4.2 | $A_{kn}B_{kpm}$ | $C_{mn[p]} = B_{k[p]m}^\top A_{kn}$ | |
| 4.3 | $A_{kn}B_{mkp}$ | $C_{mn[p]} = B_{mk[p]} A_{kn}$ | |
| 4.4 | $A_{kn}B_{pkm}$ | $TRANS(A_{kn}^\top B_{pk[m]}^\top)$ | $C_{[m][n]p} = B_{pk[m]}A_{k[n]}$ |
| 4.5 | $A_{kn}B_{mpk}$ | $C_{mn[p]} = B_{m[p]k}A_{kn}$ | |
| 4.6 | $A_{kn}B_{pmk}$ | $TRANS(A_{kn}^\top B_{p[m]k}^\top)$ | $C_{[m][n]p} = B_{p[m]k}A_{k[n]}$ |
| 5.1 | $A_{pk}B_{kmn}$ | $C_{(mn)p} = B_{k(mn)}^\top A_{pk}^\top$ | $C_{m[n]p} = B_{km[n]}^\top A_{pk}^\top$ |
| 5.2 | $A_{pk}B_{knm}$ | $C_{m[n]p} = B_{k[n]m}^\top A_{pk}^\top$ | |
| 5.3 | $A_{pk}B_{mkn}$ | $C_{m[n]p} = B_{mk[n]} A_{pk}^\top$ | |
| 5.4 | $A_{pk}B_{nkm}$ | $TRANS(B_{nk[m]}A_{pk}^\top)$ | $C_{[m]n[p]} = B_{nk[m]}A_{[p]k}$ |
| 5.5 | $A_{pk}B_{mnk}$ | $C_{(mn)p} = B_{(mn)k}A_{pk}^\top$ | $C_{m[n]p} = B_{m[n]k}A_{pk}^\top$ |
| 5.6 | $A_{pk}B_{nmk}$ | $TRANS(B_{n[m]k}A_{pk}^\top)$ | $C_{[m]n[p]} = B_{n[m]k}A_{[p]k}$ |
| 6.1 | $A_{kp}B_{kmn}$ | $C_{(mn)p} = B_{k(mn)}^\top A_{kp}$ | $C_{m[n]p} = B_{km[n]}^\top A_{kp}$ |
| 6.2 | $A_{kp}B_{knm}$ | $C_{m[n]p} = B_{k[n]m}^\top A_{kp}$ | |
| 6.3 | $A_{kp}B_{mkn}$ | $C_{m[n]p} = B_{mk[n]} A_{kp}$ | |
| 6.4 | $A_{kp}B_{nkm}$ | $TRANS(B_{nk[m]}A_{kp})$ | $C_{[m]n[p]} = B_{nk[m]}A_{k[p]}$ |
| 6.5 | $A_{kp}B_{mnk}$ | $C_{(mn)p} = B_{(mn)k}A_{kp}$ | $C_{m[n]p} = B_{m[n]k}A_{kp}$ |
| 6.6 | $A_{kp}B_{nmk}$ | $TRANS(B_{n[m]k}A_{kp})$ | $C_{[m]n[p]} = B_{n[m]k}A_{k[p]}$ |

Table: List of 36 possible single mode contraction operations between a second-order tensor and a third-order tensor and possible mappings to Level-3 BLAS routines

# Analysis

## Flatten v.s. SBGEMM

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



Prefer flatten than SBGEMM

# Analysis

## Batching in last mode v.s. middle mode

| Case | Contraction | Kernel1 | Kernel2 | Kernel3 |
|------|-------------|---------|---------|---------|
| 1.1 | $A_{mk}B_{knp}$ | $C_{m(np)} = A_{mk}B_{k(np)}$ | $C_{mn[p]} = A_{mk}B_{kn[p]}$ | $C_{m[n]p} = A_{mk}B_{k[n]p}$ |
| 2.1 | $A_{km}B_{knp}$ | $C_{m(np)} = A_{km}^{\top}B_{k(np)}$ | $C_{mn[p]} = A_{km}^{\top}B_{kn[p]}$ | $C_{m[n]p} = A_{km}^{\top}B_{k[n]p}$ |



On CPU, it's better to batch in last mode when tensor size is small/moderate
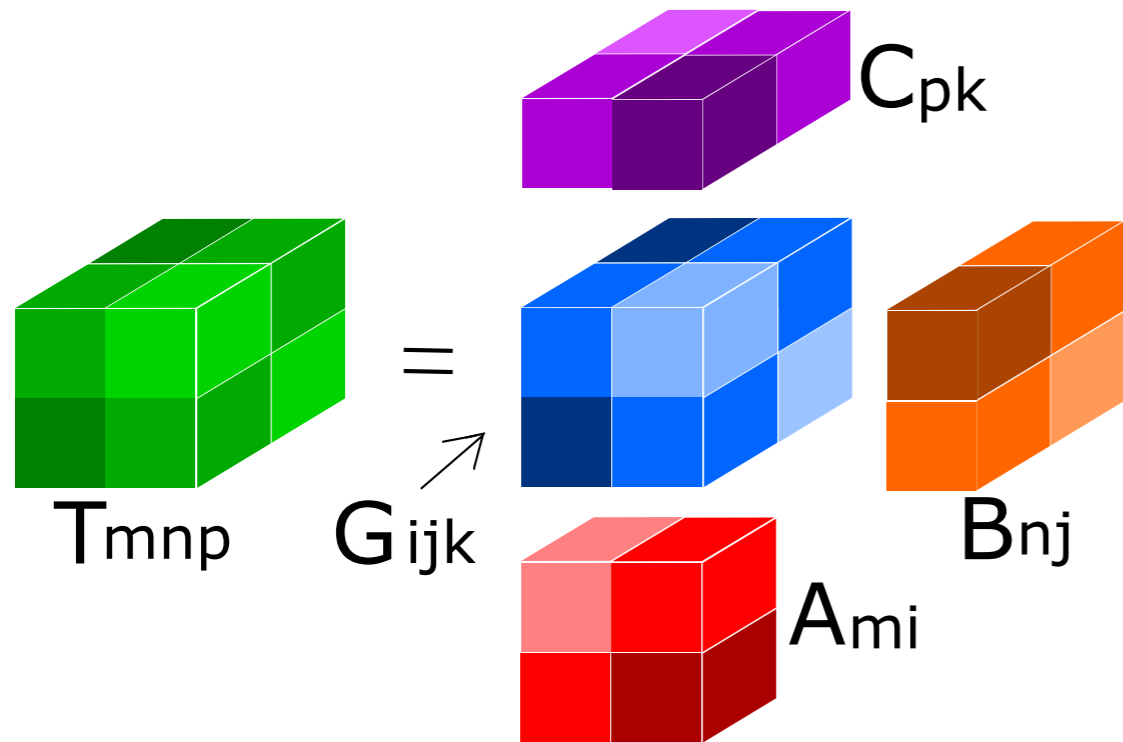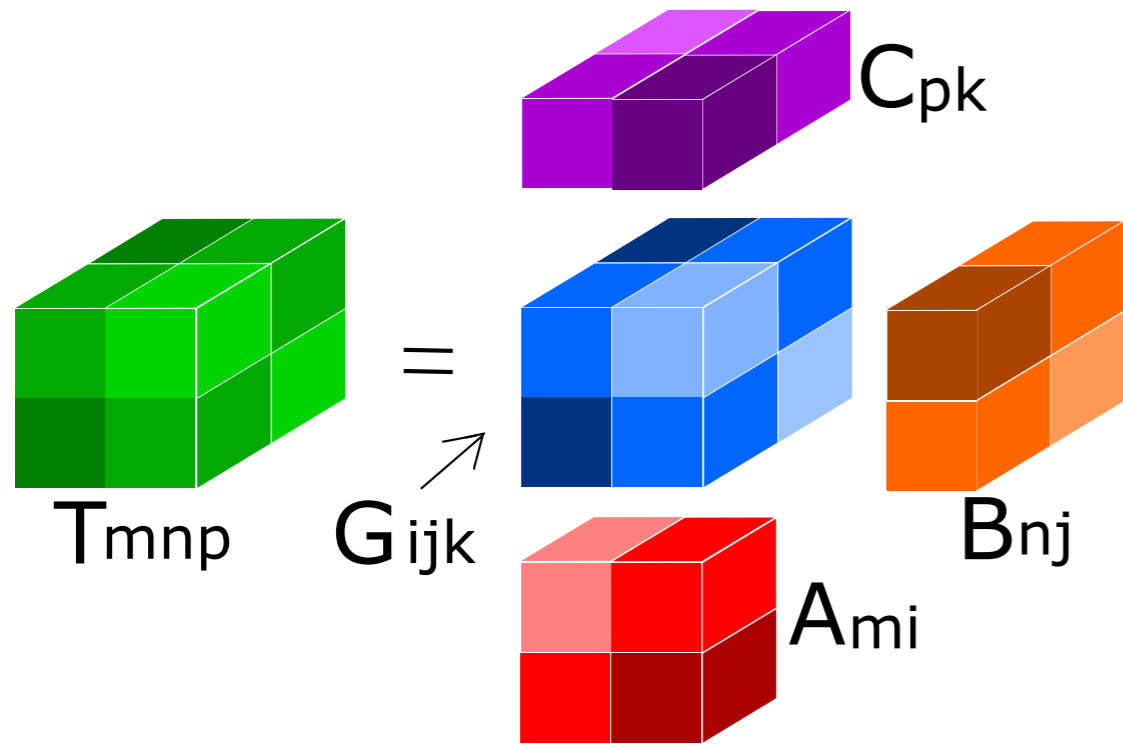
# Application: Tucker Decomposition

# Application: Tucker Decomposition

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

# Application: Tucker Decomposition

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

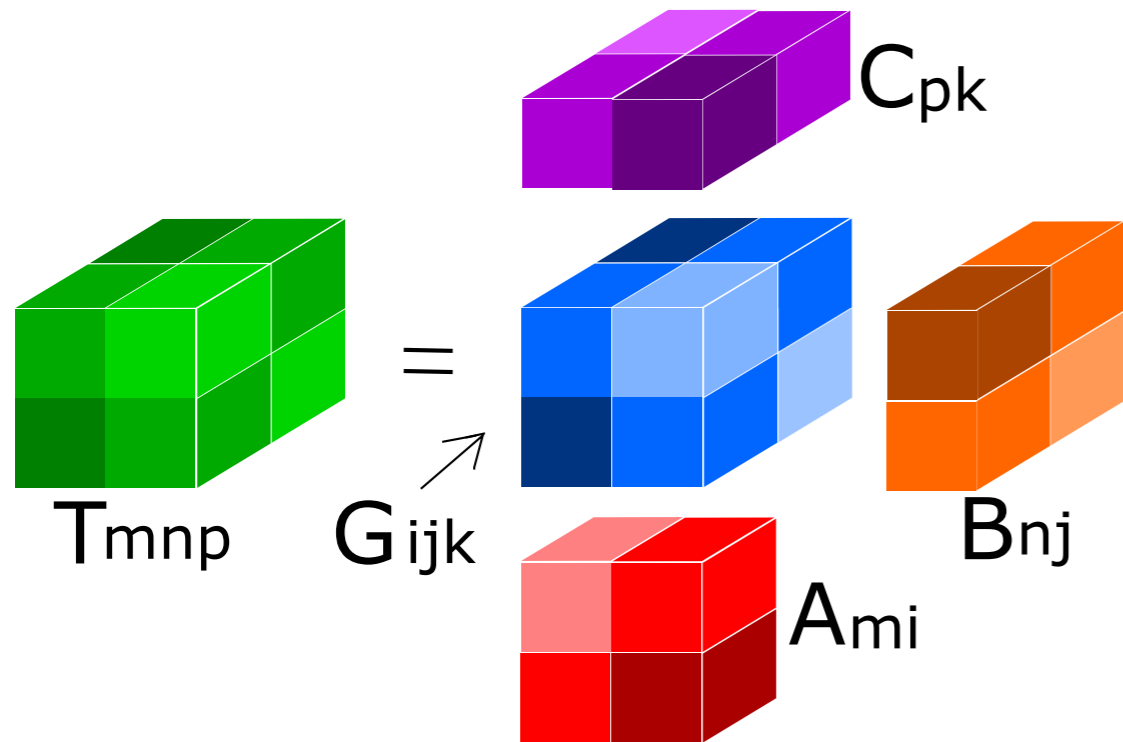# Application: Tucker Decomposition

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



Main Steps:

# Application: Tucker Decomposition

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



Main Steps:

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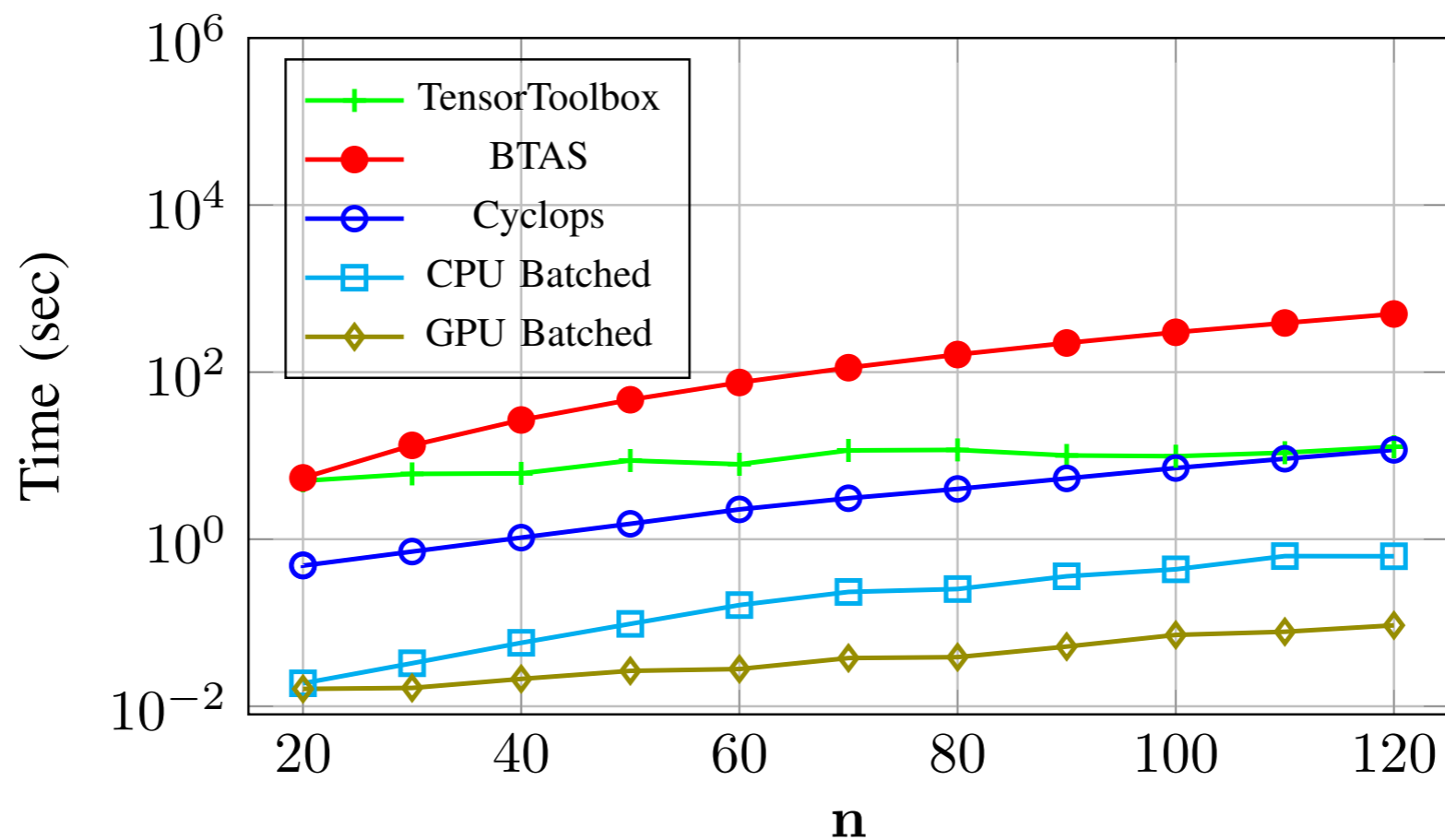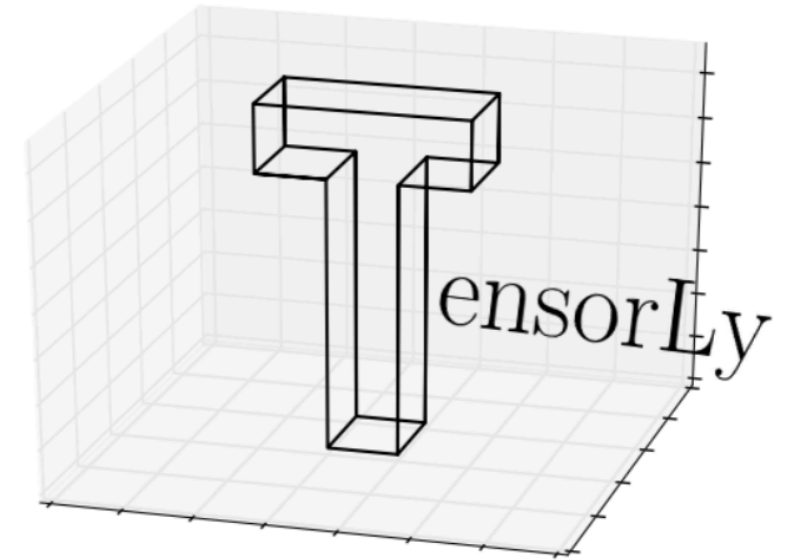
# Application: Tucker Decomposition



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.

# Introduction of TensorLy

by Jean Kossaifi, Imperial College London
   Yannis Panagakis, Imperial College London
   Anima Anandkumar, Caltech

# Introduction of TensorLy
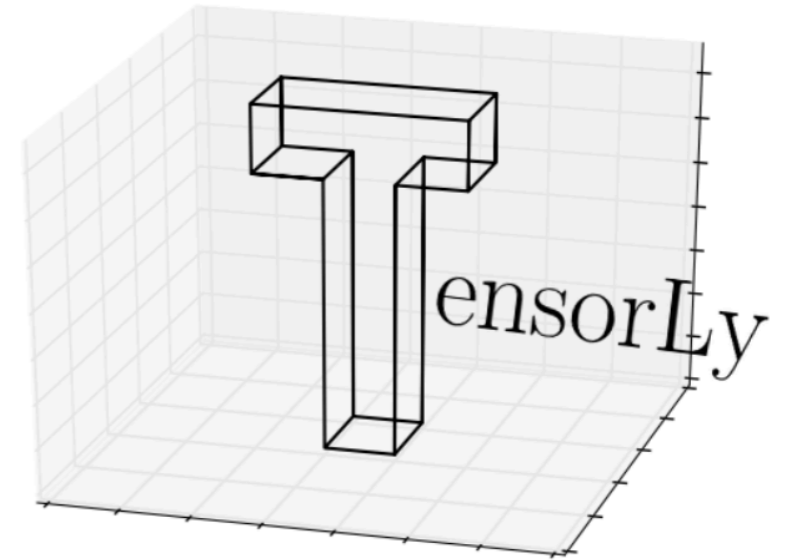
by Jean Kossaifi, Imperial College London
Yannis Panagakis, Imperial College London
Anima Anandkumar, Caltech

- **Open source**

  Homepage: http://tensorly.org/dev/

  Github: https://github.com/tensorly/tensorly

  Suitable for academic / industrial applications

# Introduction of TensorLy

**by Jean Kossaifi, Imperial College London**
**Yannis Panagakis, Imperial College London**
**Anima Anandkumar, Caltech**

- **Open source**

  Homepage: http://tensorly.org/dev/
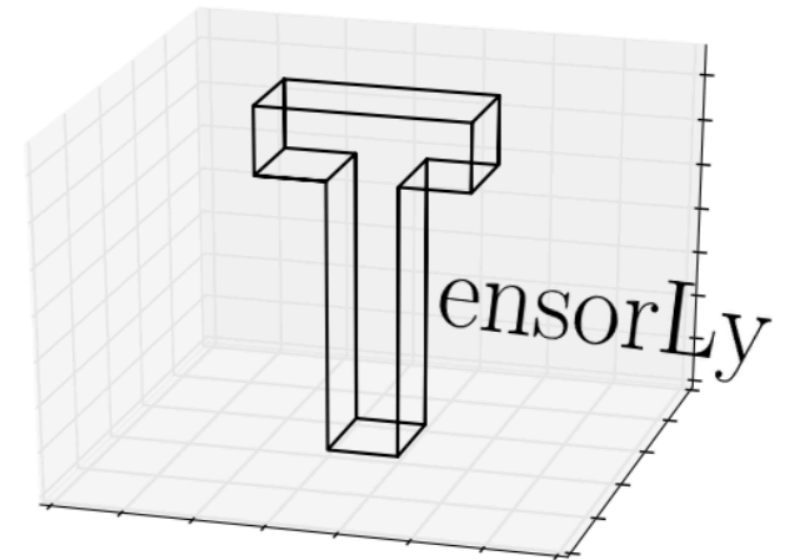
  Github: https://github.com/tensorly/tensorly

  Suitable for academic / industrial applications

- **Reliability and easy to use**

  Depends only on NumPy, SciPy [Optionally Matplotlib, MXNet and PyTorch]

  Exhaustive documentation, Unit-testing for all functions

  Fast

# User-friendly API

**TensorLy**

| Tensor decomposition | Tensor regression | Deep learning |
|---|---|---|

**Basic tensor operations**

**Unified backend**

NumPy · SciPy · mxnet · PYTORCH · TensorFlow

# TensorLy Operators

- Kronecker
- Khatri-rao
- Hadamard products
- Tensor unfolding/folding/vectorization
- N-mode product

| | |
|---|---|
| `khatri_rao` (matrices[, skip_matrix, reverse]) | Khatri-Rao product of a list of matrices |
| `kronecker` (matrices[, skip_matrix, reverse]) | Kronecker product of a list of matrices |
| `mode_dot` (tensor, matrix_or_vector, mode) | n-mode product of a tensor and a matrix or vector at the specified mode |
| `multi_mode_dot` (tensor, matrix_or_vec_list[, ...]) | n-mode product of a tensor and several matrices or vectors over several modes |
| `proximal.soft_thresholding` (tensor, threshold) | Soft-thresholding operator |
| `proximal.svd_thresholding` (matrix, threshold) | Singular value thresholding operator |
| `proximal.procrustes` (matrix) | Procrustes operator |
| `inner` (tensor1, tensor2[, n_modes]) | Generalised inner products between tensors |

- CANONICAL-POLYADIC (CP)
- Non-negative CP Tucker (HO-SVD)
- Non-negative Tucker
- Robust Tensor PCA

| | |
|---|---|
| `parafac` (tensor, rank[, n_iter_max, init, ...]) | CANDECOMP/PARAFAC decomposition via alternating least squares (ALS) |
| `non_negative_parafac` (tensor, rank[, ...]) | Non-negative CP decomposition |
| `tucker` (tensor[, rank, ranks, n_iter_max, ...]) | Tucker decomposition via Higher Order Orthogonal Iteration (HOI) |
| `partial_tucker` (tensor, modes[, rank, ...]) | Partial tucker decomposition via Higher Order Orthogonal Iteration (HOI) |
| `non_negative_tucker` (tensor, rank[, ...]) | Non-negative Tucker decomposition |
| `robust_pca` (X[, mask, tol, reg_E, reg_J, ...]) | Robust Tensor PCA via ALM with support for missing values |

# TensorLy Example

```python
from tensorly.decomposition import parafac

factors = parafac(image, rank=50, init='random')
cp_reconstruction = tl.kruskal_to_tensor(factors)


from tensorly.decomposition import tucker

core, factors = tucker(image, ranks=(50, 50, 3), init='random')
tucker_reconstruction = tl.tucker_to_tensor(core, factors)
```
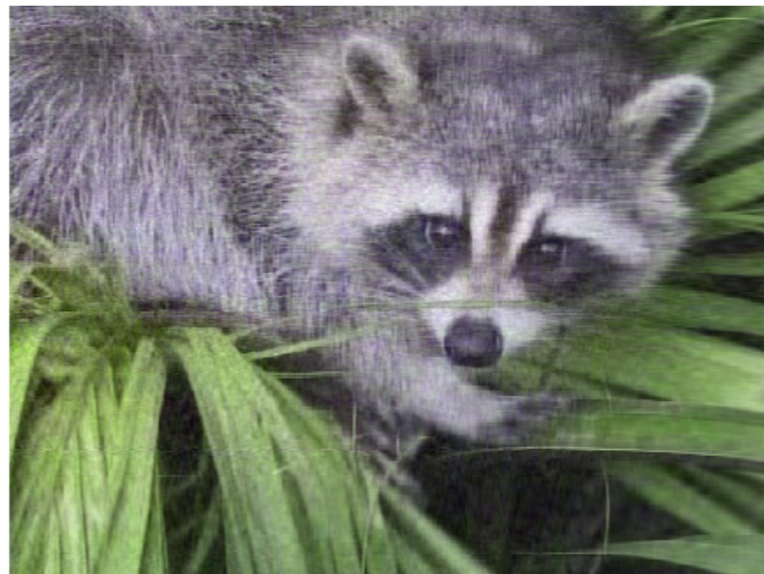


original     CP     Tucker

# TensorLy Backend

**tl.set_backend('numpy') # or 'mxnet' or 'pytorch'**

```
import tensorly as tl

T = tl.tensor([[1, 2, 3], [4, 5, 6]])
tl.tenalg.kronecker([T, T])
tl.clip(T, a_min=2, a_max=5)

tl.set_backend('mxnet')
T = tl.tensor([[1, 2, 3], [4, 5, 6]])

tl.set_backend('pytorch')
T = tl.tensor([[1, 2, 3], [4, 5, 6]])
```

NumPy ndarray

MXNet NDArray

PyTorch FloatTensor

# TensorLy Example

**Back-propagate through tensor operations with PyTorch**

```python
import tensorly as tl
from tensorly.random import tucker_tensor

tl.set_backend('pytorch')
core, factors = tucker_tensor((5, 5, 5),
                              rank=(3, 3, 3))

core = Variable(core, requires_grad=True)
factors = [Variable(f, requires_grad=True) for f in factors]

optimiser = torch.optim.Adam([core]+factors, lr=lr)

for i in range(1, n_iter):
    optimiser.zero_grad()
    rec = tucker_to_tensor(core, factors)
    loss = (rec - tensor).pow(2).sum()
    for f in factors:
        loss = loss + 0.01*f.pow(2).sum()

    loss.backward()
    optimiser.step()
```

PyTorch FloatTensor

We can attach gradients

Penalty on the factors

# Contribute to TensorLy

**Contributions welcome!**

- If you have a cool tensor method you want to add

- If you spot a bug

# Thank you!

# Questions?