



Benchmarking the Linear Algebra Awareness of TensorFlow and PyTorch

Authors | **A. Sankaran¹, NA. Alashti¹, C. Psarras¹, P. Bientinesi²**

Speaker | **Aravind Sankaran**

¹ RWTH Aachen University, Germany

² Umeå University, Sweden

State of the art

- Linear algebra operations form major performance bottlenecks.

State of the art

- Linear algebra operations form major performance bottlenecks.

1973: *Proposal for standard linear algebra subprograms:*

“[...] a fairly small number of basic operations which are generally responsible for a significant percentage of the total execution time” - Hanson, Krogh, Lawson

- Linear algebra operations form major performance bottlenecks.

1973: *Proposal for standard linear algebra subprograms:*

“[...] a fairly small number of basic operations which are generally responsible for a significant percentage of the total execution time” - **Hanson, Krogh, Lawson**

- Lots of mature numerical linear algebra libraries

2022: *Linear algebra software landscape*

PETSc, Trilinos, ...

ScaLAPACK, PLAPACK, Elemental, ...

LAPACK, Plasma, SuperMatrix, Magma, ...

BLAS-1, BLAS-2, BLAS-3, ATLAS, BTO-BLAS, BLIS, ...

State of the art

- Linear algebra operations form major performance bottlenecks.

1973: *Proposal for standard linear algebra subprograms:*

[..] a fairly small number of basic operations which are generally responsible for a significant percentage of the total execution time” - Hanson, Krogh, Lawson

- Lots of mature numerical linear algebra libraries

2022: *Linear algebra software landscape*

PETSc, Trilinos, ...

ScaLAPACK, PLAPACK, Elemental, ...

LAPACK, Plasma, SuperMatrix, Magma, ...

BLAS-1, BLAS-2, BLAS-3, ATLAS, BTO-BLAS, BLIS, ...

- Is that enough?

Kernel development: Typical Approach - ?

1. Identify the bottleneck

e.g., linear system, eigen problem, least-squares problem, SVD

2. Encapsulate into a function

e.g., dgesv, dsyevd, dgelsx, dgesvd

3. Optimize

2022: Linear algebra software landscape

PETSc, Trilinos, ...

ScaLAPACK, PLAPACK, Elemental, ...

LAPACK, Plasma, SuperMatrix, Magma, ...

BLAS-1, BLAS-2, BLAS-3, ATLAS, BTO-BLAS, BLIS, ...

Kernel development: Typical Approach - ?

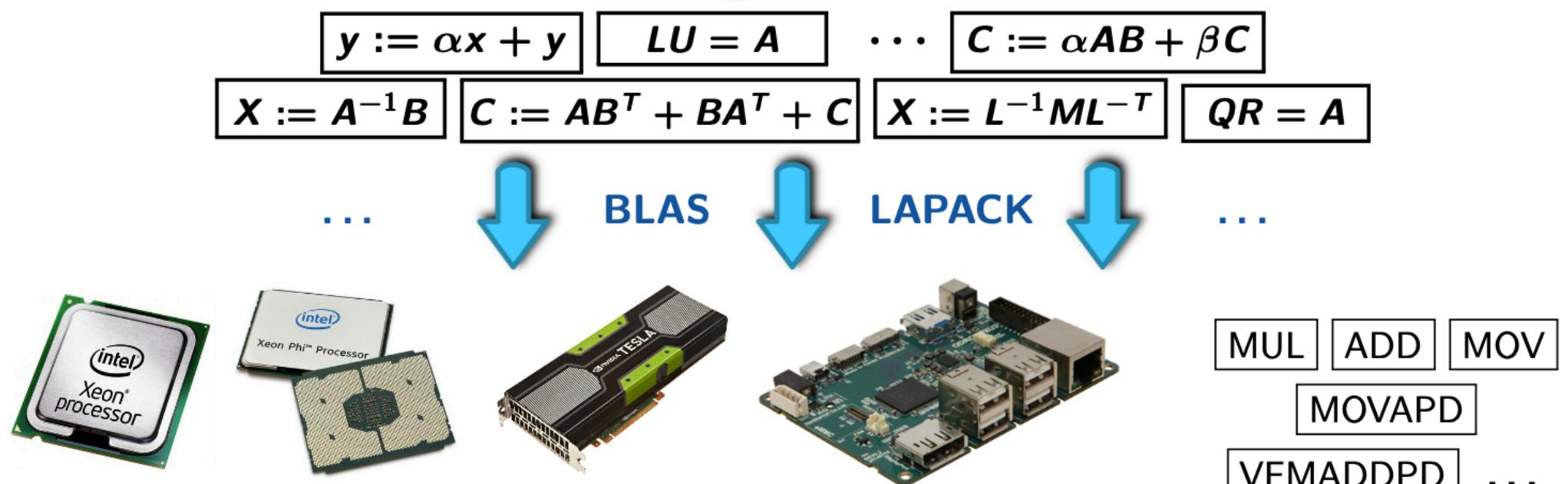
1. Identify the bottleneck

e.g., linear system, eigen problem, least-squares problem, SVD

2. Encapsulate into a function

e.g., dgesv, dsyevd, dgelsx, dgesvd

3. Optimize



2022: Linear algebra software landscape

PETSc, Trilinos, ...
ScaLAPACK, PLAPACK, Elemental, ...
LAPACK, Plasma, SuperMatrix, Magma, ...
BLAS-1, BLAS-2, BLAS-3, ATLAS, BTO-BLAS, BLIS, ...

Where do TensorFlow and PyTorch fit in?

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_\dagger := PCP^T + Q \\ K := C_\dagger H^T (H C_\dagger H^T)^{-1} \end{cases}$$

$$\begin{aligned} \Lambda &:= S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} &:= X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta \end{aligned}$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$



?

$$\begin{array}{ccccccc} y := \alpha x + y & LU = A & \dots & C := \alpha AB + \beta C \\ X := A^{-1}B & C := AB^T + BA^T + C & X := L^{-1}ML^{-T} & QR = A \end{array}$$

...



BLAS



LAPACK

...



Kernels are highly optimized.

Where do TensorFlow and PyTorch fit in?

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_\dagger := PCP^T + Q \\ K := C_\dagger H^T (H C_\dagger H^T)^{-1} \end{cases}$$

$$\begin{aligned} \Lambda &:= S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} &:= X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta \end{aligned}$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$

TensorFlow

PyTorch

$$\begin{array}{c} y := \alpha x + y \quad LU = A \quad \dots \quad C := \alpha AB + \beta C \\ X := A^{-1}B \quad C := AB^T + BA^T + C \quad X := L^{-1}ML^{-T} \quad QR = A \end{array}$$

...

BLAS

LAPACK

...

Kernels are highly optimized.



MUL ADD MOV
MOVAPD
VFMADDPD ...

Does the frameworks make best use of the kernels?

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_\dagger := PCP^T + Q \\ K := C_\dagger H^T (H C_\dagger H^T)^{-1} \end{cases}$$

$$\begin{aligned} \Lambda &:= S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} &:= X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta \end{aligned}$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$

Does the mapping make best use of the available kernels?

TensorFlow

PyTorch

$$\begin{array}{c} y := \alpha x + y \quad LU = A \quad \dots \quad C := \alpha AB + \beta C \\ X := A^{-1}B \quad C := AB^T + BA^T + C \quad X := L^{-1}ML^{-T} \quad QR = A \end{array}$$

...

BLAS

LAPACK

...

Kernels are highly optimized.



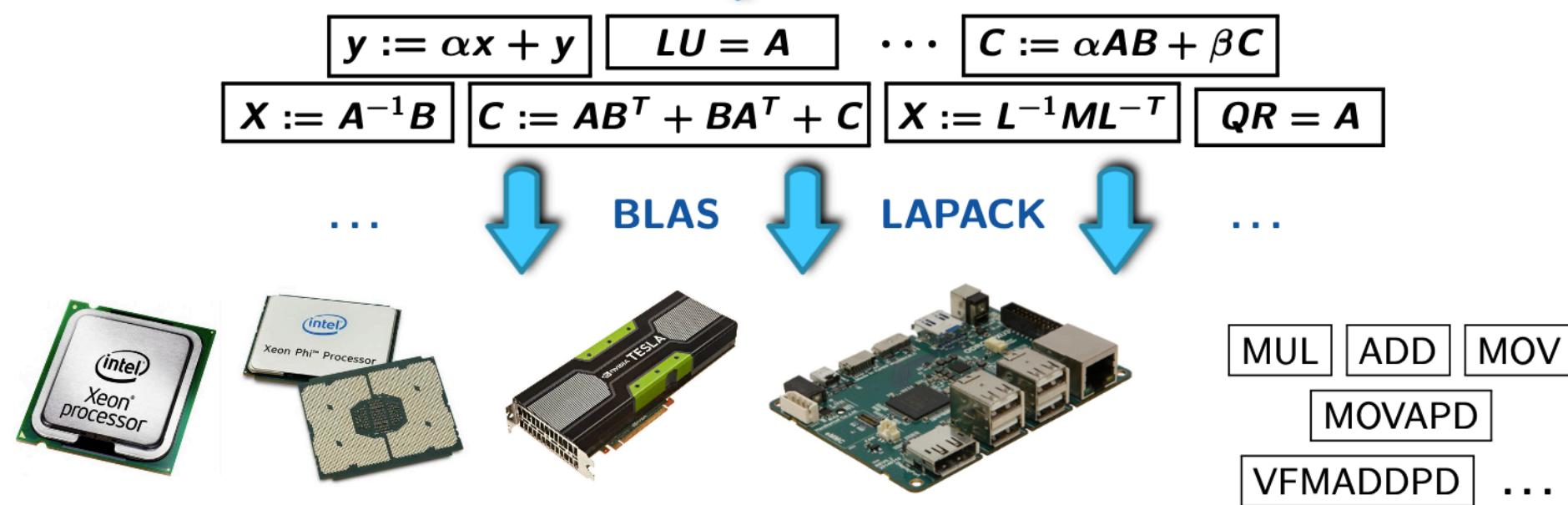
MUL ADD MOV
MOVAPD
VFMADDPD ...

Does the frameworks make best use of the kernels?

Linear Algebra Expression

$$\mathbf{y} := H^T \mathbf{y} + (I_n - H^T H) \mathbf{x}$$

where $H \in \mathbb{R}^{n \times n}$ $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$



Does the frameworks make best use of the kernels?

Linear Algebra Expression

$$\mathbf{y} := H^T \mathbf{y} + (I_n - H^T H) \mathbf{x}$$

where $H \in \mathbb{R}^{n \times n}$ $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

Operation	Kernel	FLOPs
$\mathbf{t}_1 := H^T \mathbf{y}$	gemv	$2n^2$
$T_2 := H^T H$	gemm	$2n^3$
$T_3 := I_n - T_2$	axpy	n
$\mathbf{t}_4 := T_3 \mathbf{x}$	gemv	$2n^2$
$\mathbf{y} := \mathbf{t}_1 + \mathbf{t}_4$	axpy	n

Total FLOPs $\approx 2n^3 + 4n^2 + 2n$

Exec time $\approx 0.43\text{sec}$

$n=3000$

executed on a single core of Intel Xeon CPU

$$\begin{array}{cccc} \mathbf{y} := \alpha \mathbf{x} + \mathbf{y} & LU = A & \dots & C := \alpha AB + \beta C \\ X := A^{-1}B & C := AB^T + BA^T + C & X := L^{-1}ML^{-T} & QR = A \end{array}$$

... ↓ BLAS ↓ LAPACK ...



Does the frameworks make best use of the kernels?

Linear Algebra Expression

Variant 1

$$\mathbf{y} := H^T \mathbf{y} + (I_n - H^T H) \mathbf{x}$$

$$\mathbf{y} := H^T \mathbf{y} + (I_n - H^T H) \mathbf{x}$$

where $H \in \mathbb{R}^{n \times n}$ $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

Variant 2

$$\mathbf{y} := H^T \mathbf{y} + \mathbf{x} - H^T (H \mathbf{x})$$

Operation	Kernel	FLOPs
$t_1 := H^T \mathbf{y}$	gemv	$2n^2$
$T_2 := H^T H$	gemm	$2n^3$
$T_3 := I_n - T_2$	axpy	n
$t_4 := T_3 \mathbf{x}$	gemv	$2n^2$
$\mathbf{y} := t_1 + t_4$	axpy	n

Total FLOPs $\approx 2n^3 + 4n^2 + 2n$

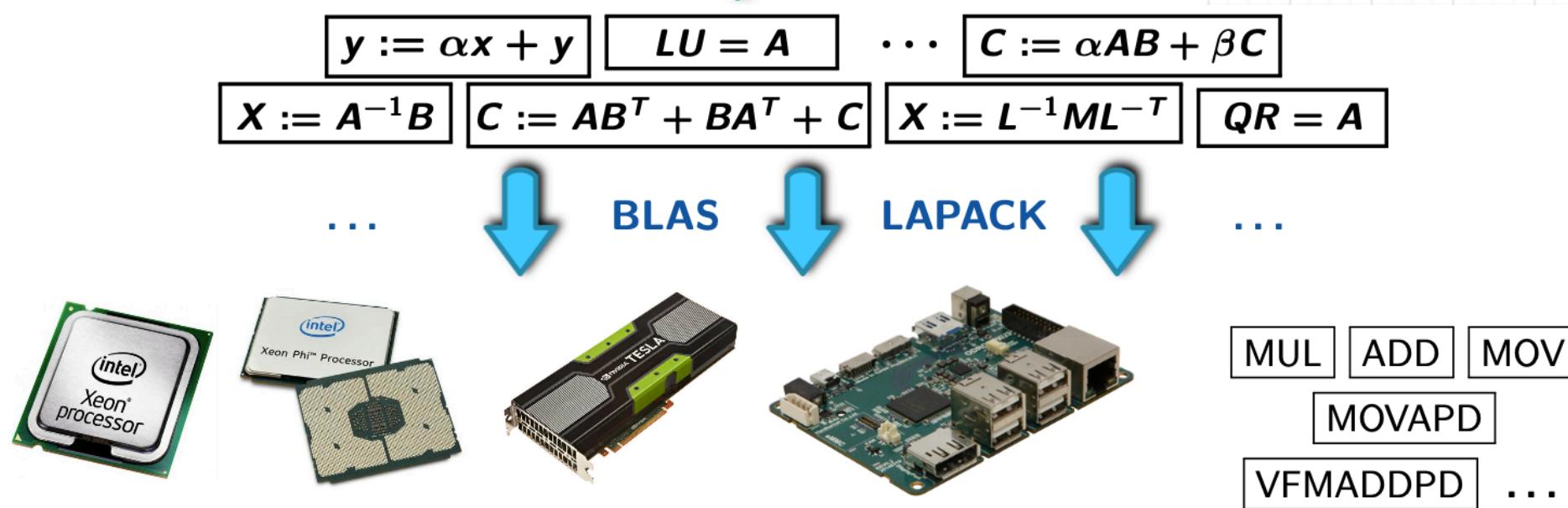
Exec time $\approx 0.43sec$

n=3000
executed on a single core of Intel Xeon CPU

Operation	Kernel	FLOPs
$t_1 := H^T \mathbf{y}$	gemv	$2n^2$
$t_2 := H \mathbf{x}$	gemv	$2n^2$
$t_3 := H^T t_2$	gemv	$2n^2$
$t_4 := t_1 - t_3$	axpy	n
$\mathbf{y} := t_4 + \mathbf{x}$	axpy	n

Total FLOPs $\approx 6n^2 + 2n$

Exec time $\approx 0.003sec$



Do we make best use of the kernels?

Linear Algebra Expression

Variant 1

$$\mathbf{y} := H^T \mathbf{y} + (I_n - H^T H) \mathbf{x}$$

$$\mathbf{y} := H^T \mathbf{y} + (I_n - H^T H) \mathbf{x}$$

where $H \in \mathbb{R}^{n \times n}$ $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

Variant 2

$$\mathbf{y} := H^T \mathbf{y} + \mathbf{x} - H^T (H \mathbf{x})$$

Operation	Kernel	FLOPs
$t_1 := H^T \mathbf{y}$	gemv	$2n^2$
$T_2 := H^T H$	gemm	$2n^3$
$T_3 := I_n - T_2$	axpy	n
$t_4 := T_3 \mathbf{x}$	gemv	$2n^2$
$\mathbf{y} := t_1 + t_4$	axpy	n

Total FLOPs $\approx 2n^3 + 4n^2 + 2n$

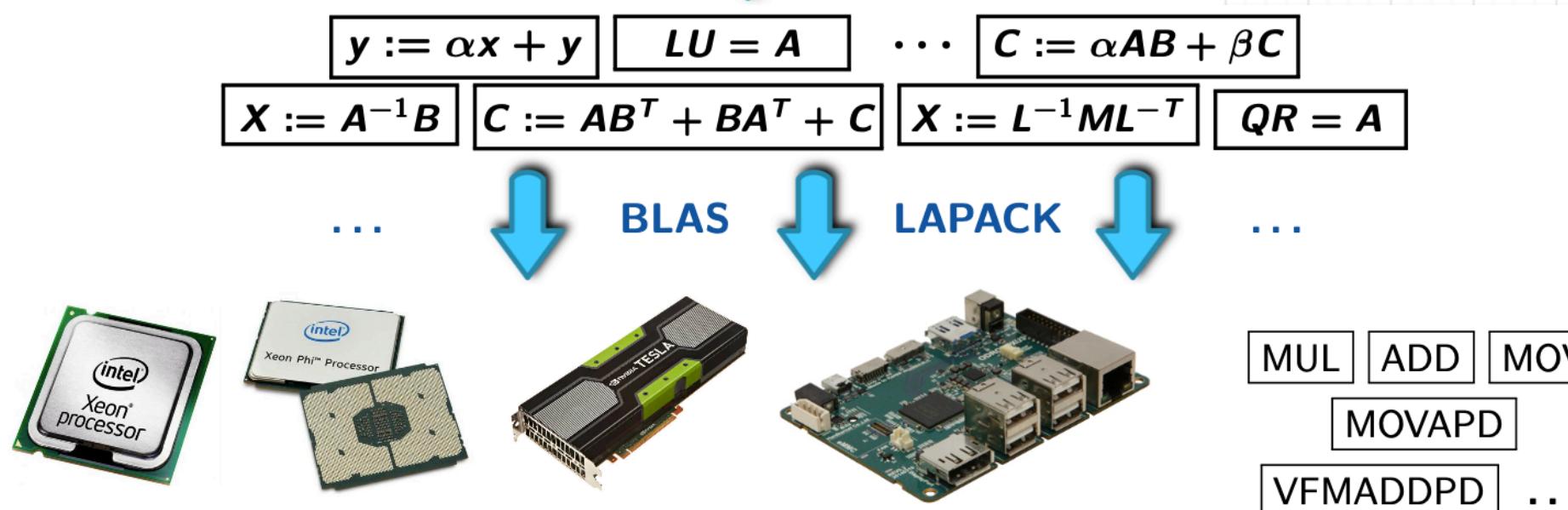
Exec time $\approx 0.43sec$

n=3000
executed on a single core of Intel Xeon CPU

Operation	Kernel	FLOPs
$t_1 := H^T \mathbf{y}$	gemv	$2n^2$
$t_2 := H \mathbf{x}$	gemv	$2n^2$
$t_3 := H^T t_2$	gemv	$2n^2$
$t_4 := t_1 - t_3$	axpy	n
$\mathbf{y} := t_4 + \mathbf{x}$	axpy	n

Total FLOPs $\approx 6n^2 + 2n$

Exec time $\approx 0.003sec$



Variant 2 is orders of magnitude faster than Variant 1!

Linear Algebra Mapping Problem

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := PCP^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I$$
$$X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y$$

$$\dots \quad E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T$$



?

There can be hundreds of possible mappings!

$$y := \alpha x + y$$

$$LU = A$$

...

$$C := \alpha AB + \beta C$$

$$X := A^{-1}B$$

$$C := AB^T + BA^T + C$$

$$X := L^{-1}ML^{-T}$$

$$QR = A$$

...

BLAS

LAPACK

...



MUL
ADD
MOV
MOVAPD
VFMADDPD
...

Linear Algebra Mapping Problem

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := PCP^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I$$
$$X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$



Linear Algebra Mapping Problem
(LAMP)

$$\begin{array}{cccc} y := \alpha x + y & LU = A & \dots & C := \alpha AB + \beta C \\ X := A^{-1}B & C := AB^T + BA^T + C & X := L^{-1}ML^{-T} & QR = A \end{array}$$

...



BLAS



LAPACK



...



MUL ADD MOV
MOVAPD
VFMADDPD
...

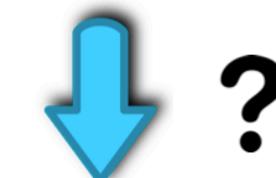
Linear Algebra Mapping Problem

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - Hx_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := PCP^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\begin{aligned} \Lambda &:= S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} &:= X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta \end{aligned}$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$



?

$$y := \alpha x + y \quad LU = A \quad \dots \quad C := \alpha AB + \beta C$$

$$X := A^{-1}B \quad C := AB^T + BA^T + C \quad X := L^{-1}ML^{-T} \quad QR = A$$

...



BLAS



LAPACK



...



$$\begin{array}{c} \text{MUL} \quad \text{ADD} \quad \text{MOV} \\ \text{MOVAPD} \\ \text{VFMADDPD} \quad \dots \end{array}$$

Linear Algebra Mapping Problem (LAMP)

- ▶ \mathcal{E} : a sequence of explicit assignments $var_i := EXP_i;$
- ▶ \mathcal{K} : a set of available computational building blocks e.g., BLAS, LAPACK, ...
- ▶ \mathcal{M} : a cost function defined over \mathcal{K}^+ #FLOPs, exec. time, #mem.ops, stability

LAMP:

Find a sequence of calls to building blocks in \mathcal{K} , optimal according to \mathcal{M} , that computes all the assignments in \mathcal{E} .

[arXiv:1911.09421]

C. Psarras, H. Barthels, P. Bientinesi,
“The Linear Algebra Mapping Problem. Current state of linear algebra languages and libraries”, ACM TOMS, 2022

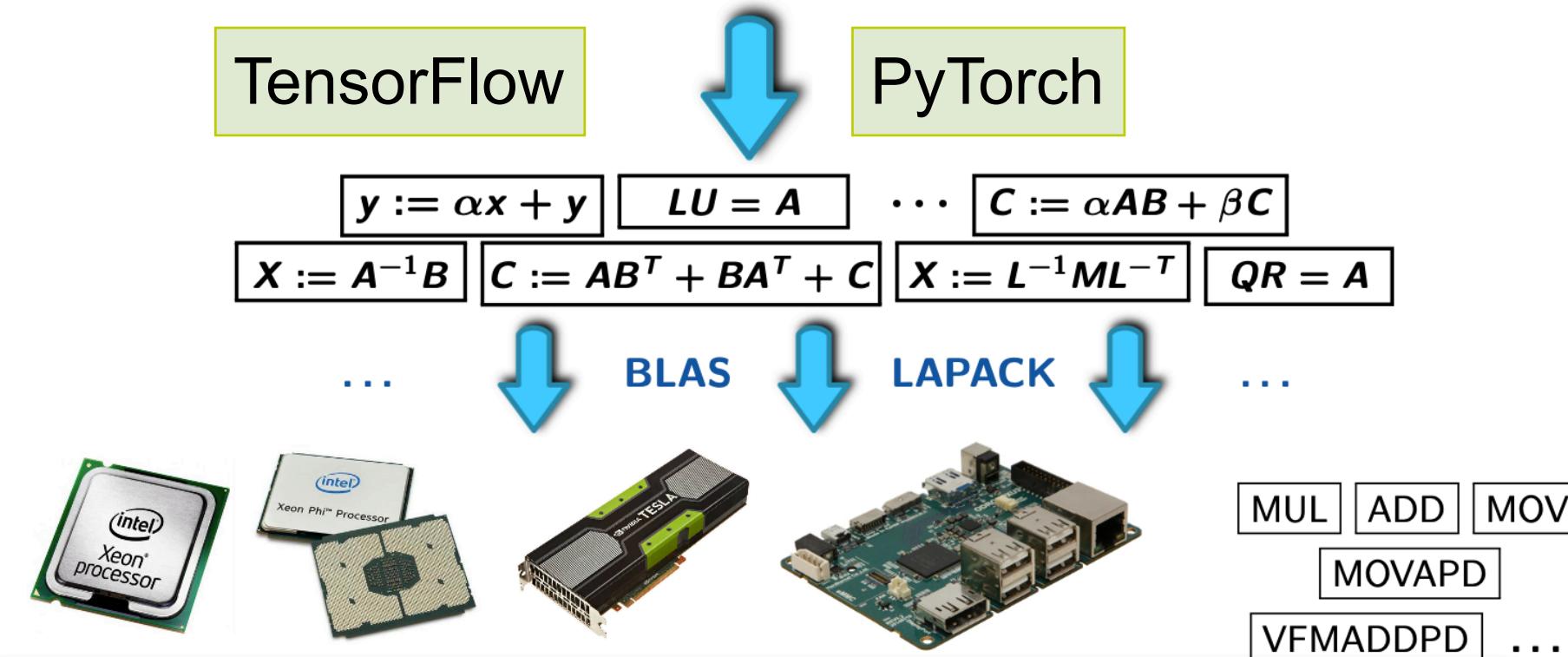
Linear Algebra Mapping Problem

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - Hx_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := PCP^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\begin{aligned} \Lambda &:= S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} &:= X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta \end{aligned}$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$



Linear Algebra Mapping Problem (LAMP)

- ▶ \mathcal{E} : a sequence of explicit assignments $var_i := EXP_i$
- ▶ \mathcal{K} : a set of available computational building blocks e.g., BLAS, LAPACK, ...
- ▶ \mathcal{M} : a cost function defined over \mathcal{K}^+ #FLOPs, exec. time, #mem.ops, stability

LAMP:

Find a sequence of calls to building blocks in \mathcal{K} , optimal according to \mathcal{M} , that computes all the assignments in \mathcal{E} .

- ▶ Suboptimal solution → easy
- ▶ Optimality → NP complete

C. Psarras, H. Barthels, P. Bientinesi,

"The Linear Algebra Mapping Problem. Current state of linear algebra languages and libraries", ACM TOMS, 2022

[arXiv:1911.09421]

Linear Algebra Awareness Benchmark

In this work:

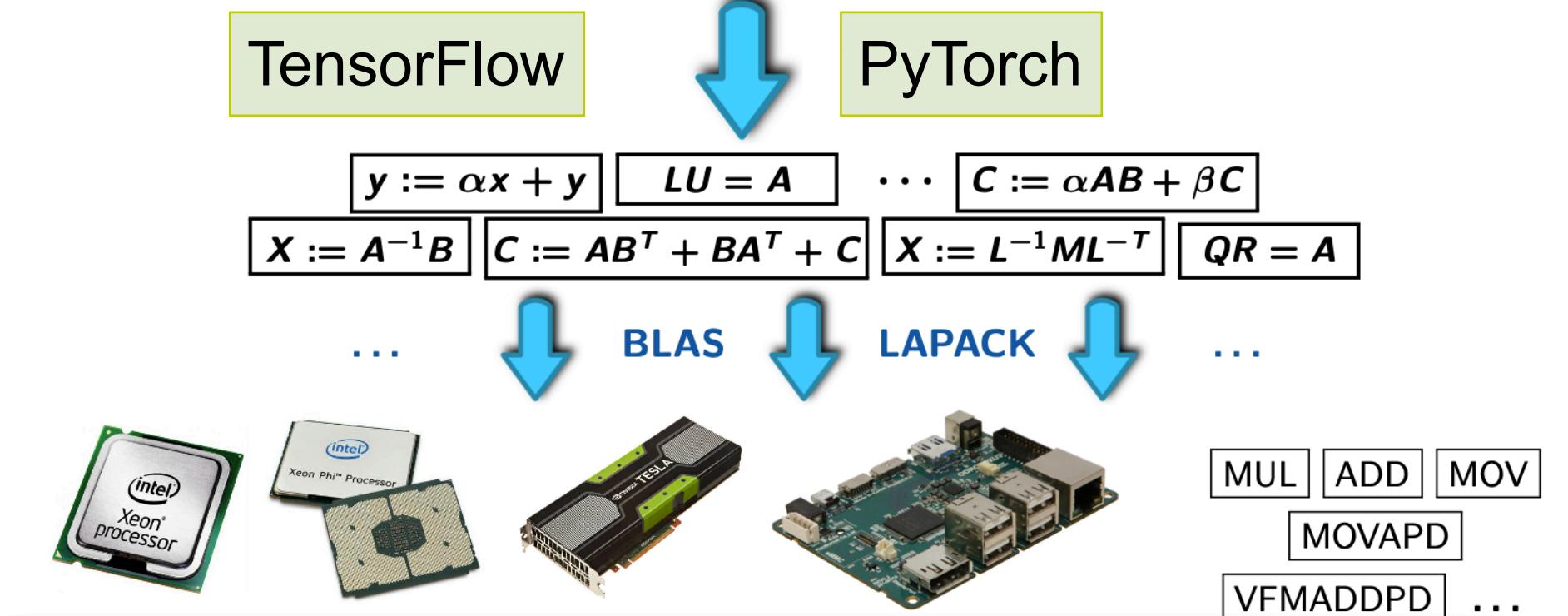
- We pin point the sub-optimals and expose optimization opportunities.
- Focus on ML frameworks - TensorFlow (TF) and PyTorch (PyT). Linear algebra operations are ubiquitous in machine learning.

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := PCP^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I$$
$$X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$



LAMP

- ▶ Suboptimal solution → easy
- ▶ Optimality → NP complete

Setting up TF and PyT

Given: $A, B \in R^{3000 \times 3000}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$

Setting up TF and PyT

Given: $A, B \in R^{3000 \times 3000}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$

Expression	MKL-C (TF / PyT)	Eager (TF / PyT)	Graph (TF / PyT)
$A^T B$	0.39	0.40 / 0.40	0.40 / 0.40
Invoke GEMM?			

Execution times in sec

Setting up TF and PyT

Given: $A, B \in R^{3000 \times 3000}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$

Expression	MKL-C (TF / PyT)	Eager (TF / PyT)	Graph (TF / PyT)
$A^T B$	0.39	0.40 / 0.40	0.40 / 0.40
Invoke GEMM?		✓	✓

Execution times in sec

Setting up TF and PyT

Given: $A, B \in R^{3000 \times 3000}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$

Expression: $Y = (A^T B)^T (A^T B)$

Mappings:

```
S1 = gemm(ATB);  
S2 = gemm(ATB);  
Y = gemm(S1TS2);
```

```
S = gemm(ATB);  
Y = gemm(STS);
```

Expression	MKL-C (TF / PyT)	Eager (TF / PyT)	Graph (TF / PyT)
$A^T B$	0.39	0.40 / 0.40	0.40 / 0.40
Invoke GEMM?		✓	✓

Setting up TF and PyT

Given: $A, B \in R^{3000 \times 3000}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$

Expression: $Y = (A^T B)^T (A^T B)$

Mappings:

```
S1 = gemm(ATB);  
S2 = gemm(ATB);  
Y = gemm(S1TS2);
```

```
S = gemm(ATB);  
Y = gemm(STS);
```

Expression	MKL-C (TF / PyT)	Eager (TF / PyT)	Graph (TF / PyT)
$A^T B$	0.39	0.40 / 0.40	0.40 / 0.40
Invoke GEMM?		✓	✓
Graph optimizations?			
Execution times in sec			

Setting up TF and PyT

Given: $A, B \in R^{3000 \times 3000}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$

Expression: $Y = (A^T B)^T (A^T B)$

Mappings:

```
S1 = gemm(ATB);  
S2 = gemm(ATB);  
Y = gemm(S1TS2);
```

Expression	MKL-C (TF / PyT)	Eager (TF / PyT)	Graph (TF / PyT)
$A^T B$	0.39	0.40 / 0.40	0.40 / 0.40
Invoke GEMM?		✓	✓
$(A^T B)^T (A^T B)$	-	1.25 / 1.27	0.78 / 0.80
Graph optimizations?		✗	✓

Execution times in sec

```
S = gemm(ATB);  
Y = gemm(STS);
```

Graph mode in TF and PyT

```

ret := (ATB)T(ATB)          A, B ∈ ℝn×n
import tensorflow as tf

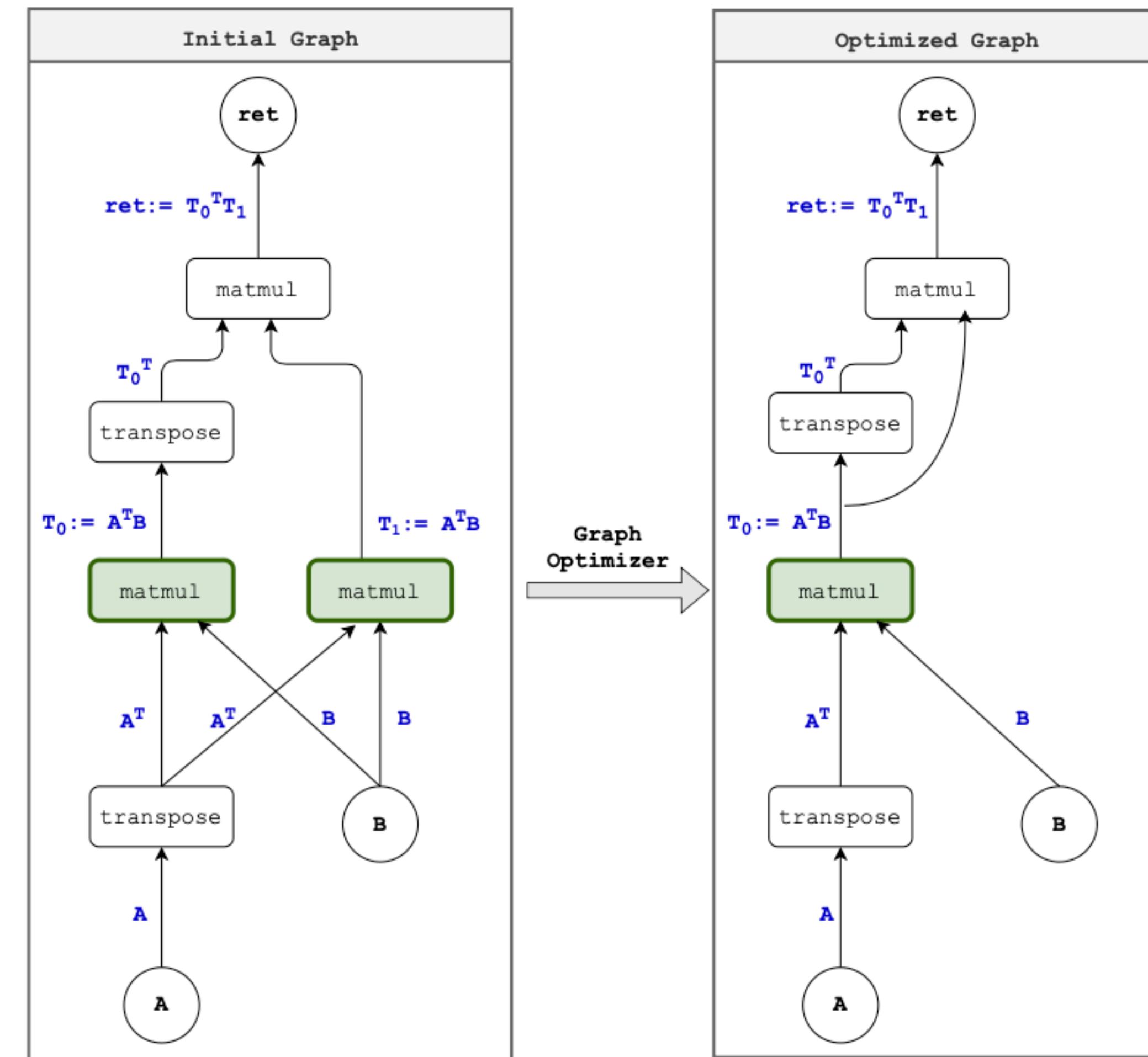
@tf.function
def transformation(A, B):
    AT = tf.transpose(A)
    ret = tf.transpose(AT@B) @ (AT@B)
    return ret
}

```

TensorFlow Code.

Expression	MKL-C	Eager (TF / PyT)	Graph (TF / PyT)
$A^T B$	0.39	0.40 / 0.40	0.40 / 0.40
$(A^T B)^T (A^T B)$	-	1.25 / 1.27	0.78 / 0.80

Execution time (in sec) for $n = 3000$



The Computational Graphs for $(A^T B)^T (A^T B)$.

Linear Algebra Awareness Benchmark

Conditions:

- Linked to optimized libraries (e.g., Intel MKL).
- Experiments are run in graph mode

Goal:

Evaluate the extent to which TF and PyT incorporate Linear algebra knowledge to speed up computations.

Experiments to evaluate Linear Algebra Awareness:

1. Common Sub-expression Elimination
2. Optimization of Matrix Chains
3. Exploiting Matrix Properties
4. Algebraic Manipulations
5. Code Motion

Experiment 1: Common Sub-expression Elimination

	Expression	TF	PyT
<i>Parenthesized:</i>	$(A^T B)^T (A^T B)$	0.78	0.80

Experiment 1: Common Sub-expression Elimination

Expression: $Y = (A^T B)^T A^T B$

Mappings:

```
S1 = gemm (ATB) ;
S2 = gemm (ATB) ;
Y = gemm (S1TS2) ;
```

```
S = gemm (ATB) ;
Y = gemm (STS) ;
```

```
S1 = gemm (ATB) ;
S2 = gemm (S1TAT) ;
Y = gemm (S2B) ;
```

Parenthesized:

Expression	TF	PyT
$(A^T B)^T (A^T B)$	0.78 ✓	0.80 ✓

Experiment 1: Common Sub-expression Elimination

Expression: $Y = (A^T B)^T A^T B$

Mappings:

```
S1 = gemm (ATB) ;
S2 = gemm (ATB) ;
Y = gemm (S1TS2) ;
```

```
S = gemm (ATB) ;
Y = gemm (STS) ;
```

```
S1 = gemm (ATB) ;
S2 = gemm (S1TAT) ;
Y = gemm (S2B) ;
```

Parenthesized:

Not Parenthesized:

Expression	TF	PyT
$(A^T B)^T (A^T B)$	0.78 ✓	0.80 ✓
$(A^T B)^T A^T B$	1.17	1.15

Experiment 1: Common Sub-expression Elimination

Expression: $Y = (A^T B)^T A^T B$

Mappings:

```
S1 = gemm (ATB) ;
S2 = gemm (ATB) ;
Y = gemm (S1TS2) ;
```

```
S = gemm (ATB) ;
Y = gemm (STS) ;
```

```
S1 = gemm (ATB) ;
S2 = gemm (S1TAT) ;
Y = gemm (S2B) ;
```

Parenthesized:

Not Parenthesized:

Expression	TF	PyT
$(A^T B)^T (A^T B)$	0.78 ✓	0.80 ✓
$(A^T B)^T A^T B$	1.17 ✗	1.15 ✗

Experiment 1: Common Sub-expression Elimination

Expression: $Y = (A^T B)^T A^T B$

Mappings:

```
S1 = gemm (ATB) ;
S2 = gemm (ATB) ;
Y = gemm (S1TS2) ;
```

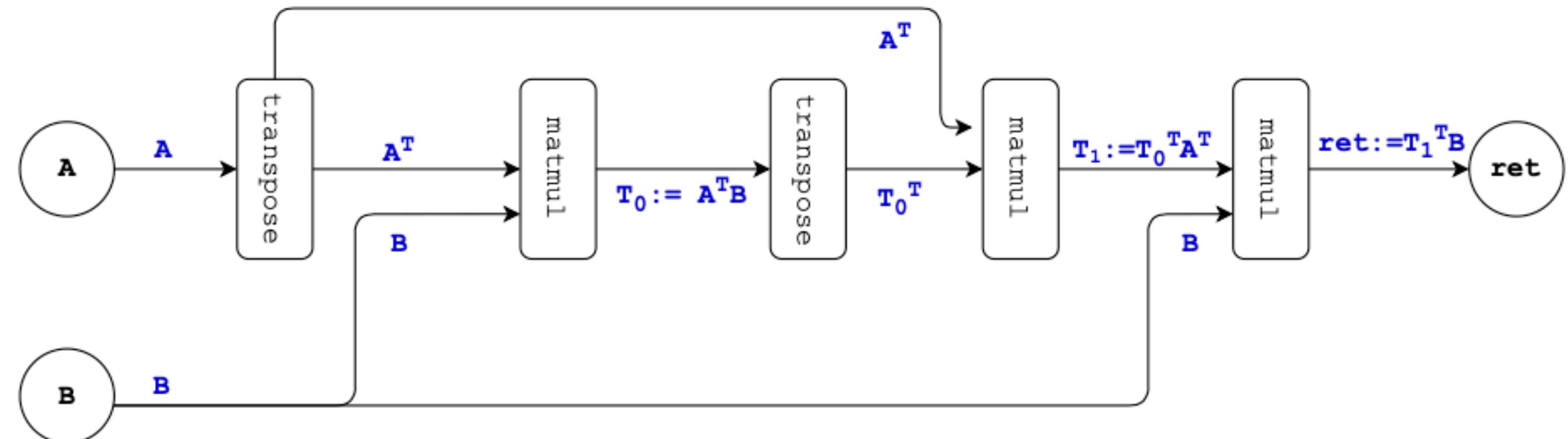
```
S = gemm (ATB) ;
Y = gemm (STS) ;
```

```
S1 = gemm (ATB) ;
S2 = gemm (S1TAT) ;
Y = gemm (S2B) ;
```

Parenthesized:

Expression	TF	PyT
$(A^T B)^T (A^T B)$	0.78	0.80
	✓	✓

Not Parenthesized:	Execution times in sec
$(A^T B)^T A^T B$	1.17
	✗
	✗



The Computational Graph for $(A^T B)^T A^T B$.

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul	matmul	multi_dot	
$H^T H \mathbf{x}$	0.40	0.41	0.006	
$H^T(H\mathbf{x})$	0.006	0.004	-	

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT		$H^T(H\mathbf{x})$
	matmul	matmul	matmul	multi_dot	
$H^T H\mathbf{x}$	0.40	0.41	0.006		
$H^T(H\mathbf{x})$	0.006	0.004	-		

$$(H^T H)\mathbf{x}$$

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul		matmul	multi_dot
$H^T H \mathbf{x}$	0.40		0.41	0.006
$H^T(H\mathbf{x})$	0.006		0.004	-

$$H^T(H\mathbf{x})$$

$$\begin{aligned} \mathbf{t}_0 &\leftarrow H\mathbf{x} & 2n^2 \\ \mathbf{t}_1 &\leftarrow H^T\mathbf{t}_0 & 2n^2 \end{aligned}$$

$$(H^T H)\mathbf{x}$$

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul	matmul	multi_dot	
$H^T H \mathbf{x}$	0.40	0.41	0.006	
$H^T(H\mathbf{x})$	0.006	0.004	-	

$$H^T(H\mathbf{x})$$

$$\mathbf{t}_0 \leftarrow H\mathbf{x} \quad 2n^2$$

$$\mathbf{t}_1 \leftarrow H^T\mathbf{t}_0 \quad 2n^2$$

$$(H^T H)\mathbf{x}$$

$$T_0 \leftarrow H^T H \quad 2n^3$$

$$\mathbf{t}_1 \leftarrow T_0\mathbf{x} \quad 2n^2$$

Experiment 2: Optimization of Matrix Chains

Expression	TF			PyT		
	matmul	matmul	multi_dot	matmul	multi_dot	
$H^T H \mathbf{x}$	0.40	0.41	0.006			
$H^T(H\mathbf{x})$	0.006	0.004	-			

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul	matmul	multi_dot	
$H^T H \mathbf{x}$	0.40	0.41	0.006	
$H^T(H\mathbf{x})$	0.006	0.004	-	
Right to left?	✗	✗	✓	

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul	matmul	matmul	multi_dot
$H^T H \mathbf{x}$	0.40	0.41	0.006	
$H^T(H\mathbf{x})$	0.006	0.004	-	
Right to left?	✗	✗	✓	
<hr/>				
$\mathbf{y}^T H^T H$	0.006	0.005	0.005	
$(\mathbf{y}^T H^T)H$	0.006	0.005	-	
Left to right?	✓	✓	✓	

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul	matmul	matmul	multi_dot
$H^T H \mathbf{x}$	0.40	0.41	0.006	
$H^T(H\mathbf{x})$	0.006	0.004	-	
Right to left?	✗	✗	✓	
$\mathbf{y}^T H^T H$	0.006	0.005	0.005	
$(\mathbf{y}^T H^T)H$	0.006	0.005	-	
Left to right?	✓	✓	✓	
$H^T \mathbf{y} \mathbf{x}^T H$	0.41	0.40	0.01	
$(H^T \mathbf{y})(\mathbf{x}^T H)$	0.01	0.01	-	
Mixed?	✗	✗	✓	

Experiment 2: Optimization of Matrix Chains

Expression	TF		PyT	
	matmul	matmul	matmul	multi_dot
$H^T H \mathbf{x}$	0.40	0.41	0.006	
$H^T(H\mathbf{x})$	0.006	0.004	-	
Right to left?	✗	✗	✓	
$\mathbf{y}^T H^T H$	0.006	0.005	0.005	
$(\mathbf{y}^T H^T)H$	0.006	0.005	-	
Left to right?	✓	✓	✓	
$H^T \mathbf{y} \mathbf{x}^T H$	0.41	0.40	0.01	
$(H^T \mathbf{y})(\mathbf{x}^T H)$	0.01	0.01	-	
Mixed?	✗	✗	✓	

- Evaluate all possible parenthesizations and choose the one with minimum FLOP?

- Number of possible parenthesizations for a matrix chain of length m is given by the $(m-1)^{th}$ Catalan number:

$$C_{m-1} = \frac{(2m)!}{(m+1)!m!}.$$

- Evaluate optimal order at least for $m < K$?

Experiment 3: Exploiting Matrix Properties

Given: $A, B \in R^{n \times n}$

Expression: $Y = A^T B$

Mapping: $\boxed{Y = \text{gemm}(A^T B)} \quad 2n^3$

If A is triangular (L)

$\boxed{Y = \text{trmm}(L^T B)} \quad n^3$

If $B = A$

$\boxed{Y = \text{syrk}(A^T A)} \quad n^3$

Experiment 3: Exploiting Matrix Properties

Given: $A, B \in R^{n \times n}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$ $2n^3$

If A is triangular (L)

$Y = \text{trmm}(L^T B)$ n^3

If $B = A$

$Y = \text{syrk}(A^T A)$ n^3

Expr	SciPy		TF		PyT	
	BLAS	matmul	optim	matmul	optim	
AB	0.40	0.40	-	0.41	-	
LB	0.24	0.40	n.a	0.40	n.a	X
AA^T	0.24	0.41	n.a	0.39	n.a	X

Execution times in sec

Experiment 3: Exploiting Matrix Properties

Given: $A, B \in R^{n \times n}$

Expression: $Y = A^T B$

Mapping: $Y = \text{gemm}(A^T B)$ $2n^3$

If A is triangular (L)

$Y = \text{trmm}(L^T B)$ n^3

If $B = A$

$Y = \text{syrk}(A^T A)$ n^3

Expr	SciPy		TF		PyT	
	BLAS	matmul	optim	matmul	optim	
AB	0.40	0.40	-	0.41	-	
LB	0.24	0.40	n.a	0.40	n.a	X
AA^T	0.24	0.41	n.a	0.39	n.a	X

- If matrix properties can be propagated through the computational graph, it would be possible to detect algebraic simplifications that can speed up computation.

E.g. if it can be known that Q is orthogonal, then $Q^T Q = I$; matrix multiplication can be avoided.

Experiment 4: Algebraic Manipulations

Distributive Law:

	TF		PyT		
	LHS	RHS	LHS	RHS	
$AB + AC = A(B + C)$	0.78	040	0.81	0.41	✗
$A\mathbf{x} - H^T(H\mathbf{x}) = (A - H^T H)\mathbf{x}$	0.01	0.42	0.01	0.41	✗

Experiment 4: Algebraic Manipulations

Blocked Operands:

$$A_B := \begin{bmatrix} A_1 & 0 \\ 0 & A_2 \end{bmatrix} \in \mathbb{R}^{n \times n} \quad B_B := \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$A_B B_B = \begin{bmatrix} (A_1 B_1) \\ (A_2 B_2) \end{bmatrix}$$

	TF		PyT	
	LHS	RHS	LHS	RHS
Blocked matrices	0.40	0.20	0.40	0.20 X

Experiment 5: Code Motion

Loop-Invariant Code Motion:

Naive

```
V = [v1, v2, v3]  
for i in range(3):  
    Y[:, :, i] = A @ B + V[i] @ V[i]T
```

Recommended

```
V = [v1, v2, v3]  
tmp = A @ B  
for i in range(3):  
    Y[:, :, i] = tmp + V[i] @ V[i]T
```

Property	TF		PyT	
	Naive	Reco	Naive	Reco
Loop-inv code motion	0.42	0.42	0.42	0.41



Experiment 5: Code Motion

Partial operand access:

Naive

```
#Sum  
Y = (A+B) [2, 2]  
  
#Product  
Y = (A@B) [2, 2]
```

Recommended

```
#Sum  
Y = A[2, 2] + B[2, 2]  
  
#Product  
Y = dot(A[2, :], B[:, 2])
```

Property

	TF		PyT	
	Naive	Reco	Naive	Reco
Partial-op access (sum)	0.011	6e-4	0.018	2e-3
Partial-op access (product)	0.39	2e-3	0.40	3e-3

Summary

- Evaluated the extent to which TF and PyT use Linear Algebra knowledge to speed up computations.
- We exposed some of the sub-optimalities presented guidelines to help both the framework developers and end-users to achieve better performance.

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k(z_k - Hx_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := PCP^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U(I + U^T Q^{-1} U)^{-1} U^T$$

