

# TSM2: Optimizing Tall-and-Skinny Matrix-Matrix Multiplication on GPUs

**Jieyang Chen**, Nan Xiong, Xin Liang, Dingwen Tao\*, Sihuan Li, Kaiming Ouyang,  
Kai Zhao, Nathan DeBardeleben\*\*, Qiang Guan\*\*\*, Zizhong Chen

*University of California, Riverside*

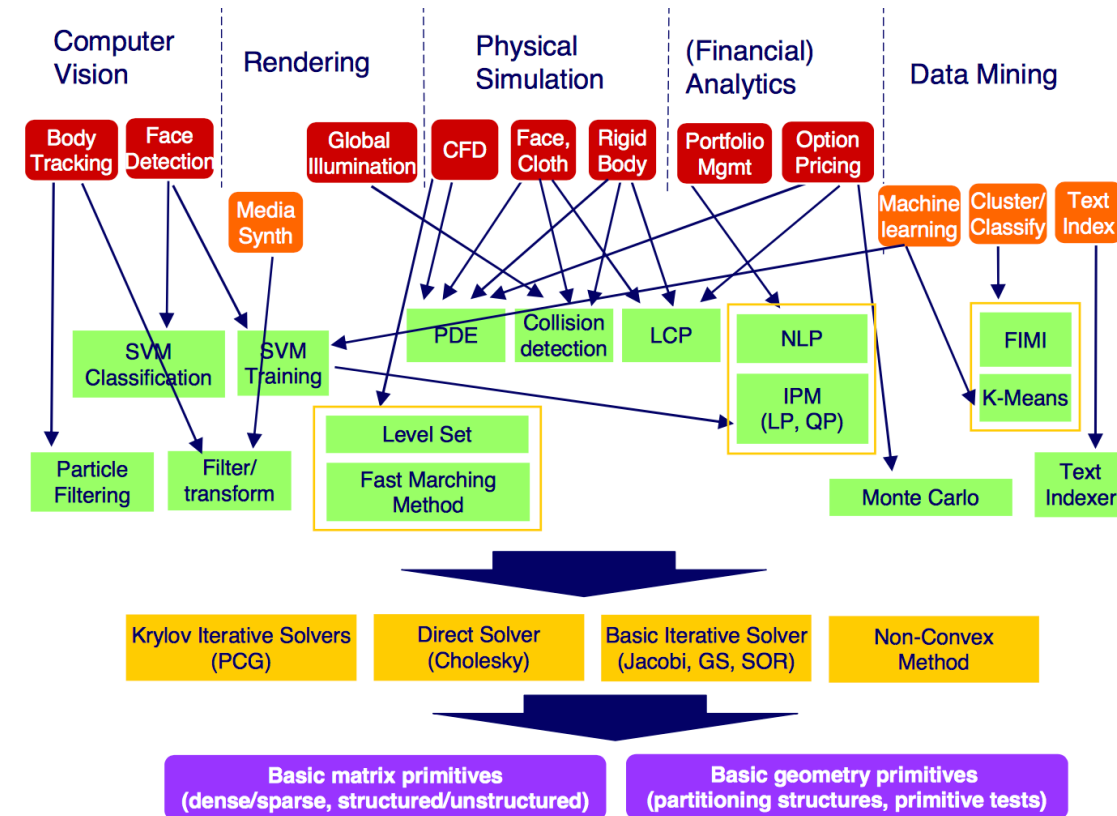
*\*\*Los Alamos National Laboratory*

*\*University of Alabama*

*\*\*\*Kent State University*

# Linear algebra kernels are widely used

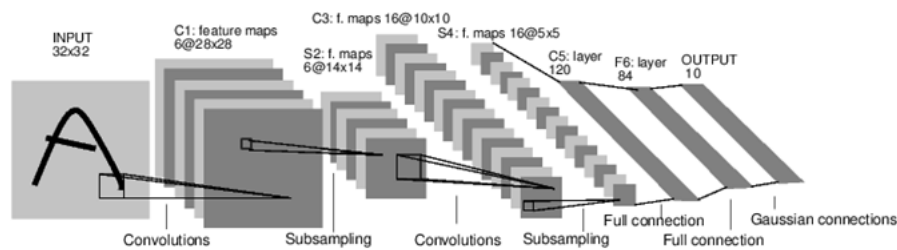
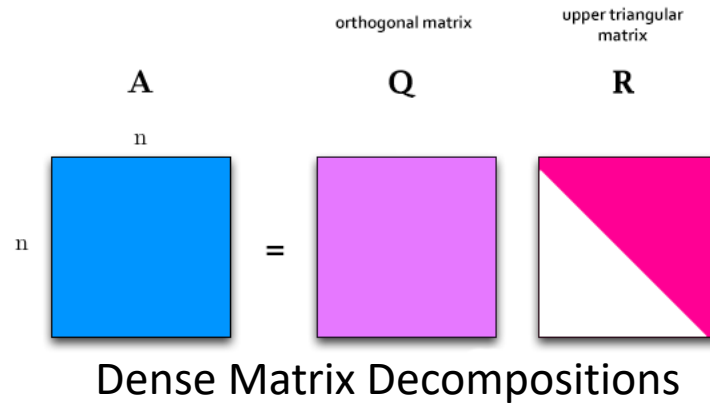
- Linear algebra kernels have been widely used.
  - *E.g., scientific simulation, big data analytics, machine learning, etc.*
- Matrix-matrix multiplication (**GEMM**)
  - One of the most fundamental computation kernel that is used to build up other kernels
  - Core computation of many applications.
  - Cost most of the computation time of applications



(Source: Berkeley Dwarfs Report)

# Input shape of GEMM can varies from application to application

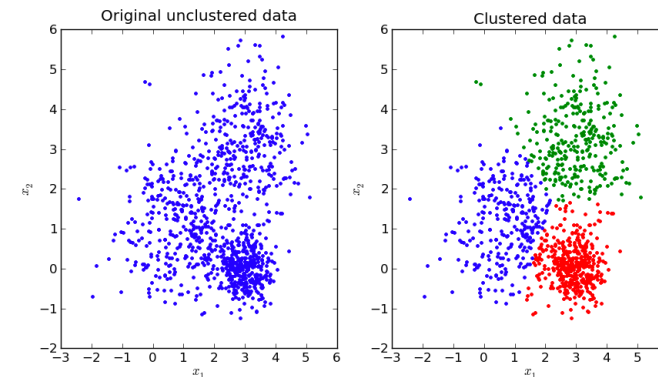
## Relative regular shape input



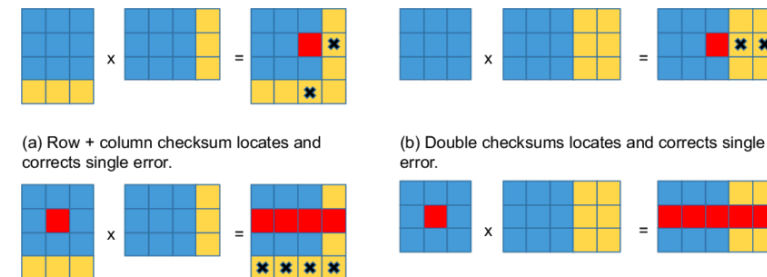
A Full Convolutional Neural Network (LeNet)

## Deep Neural Networks

## Tall-and-skinny shape input



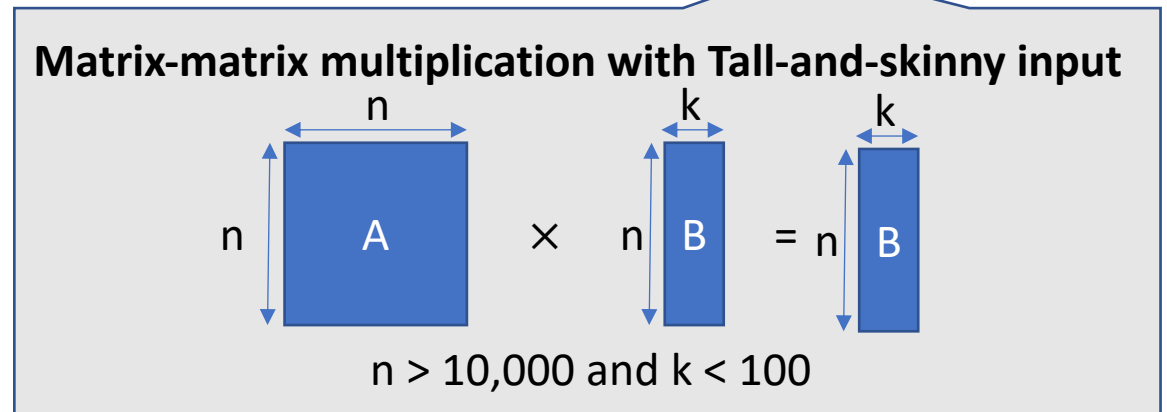
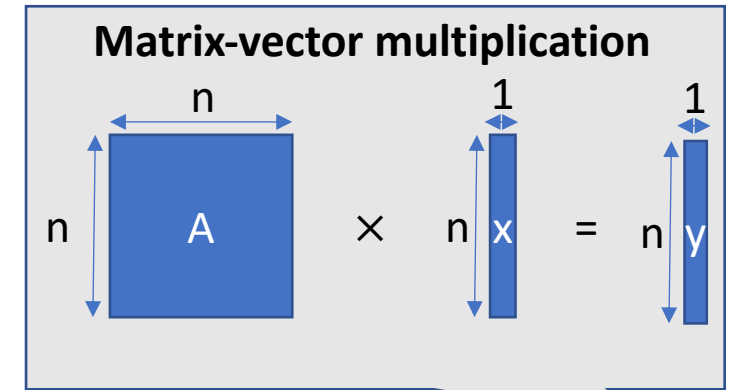
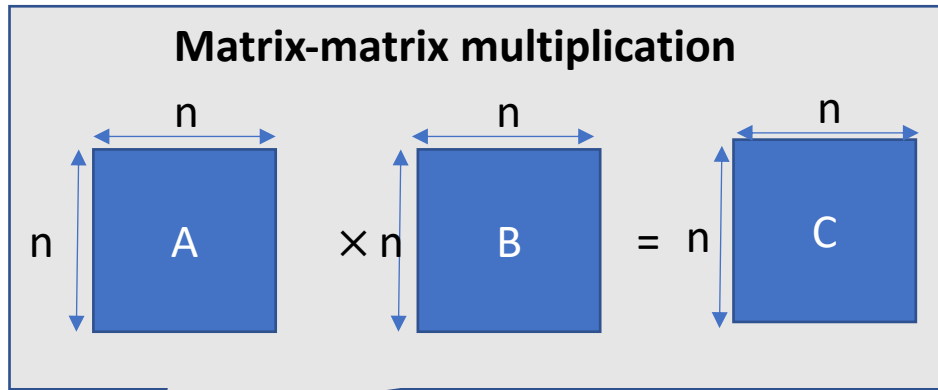
## K-means



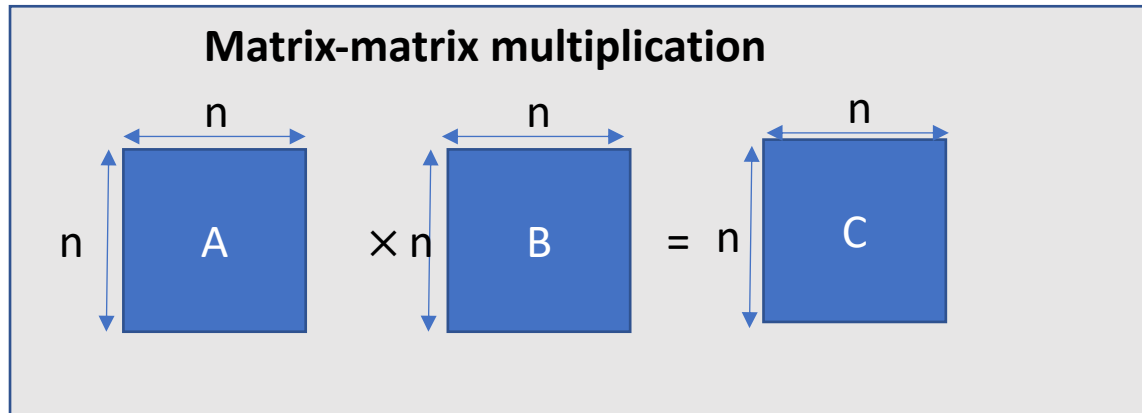
## Algorithm Based Fault Tolerance

# Two Kinds of Computations

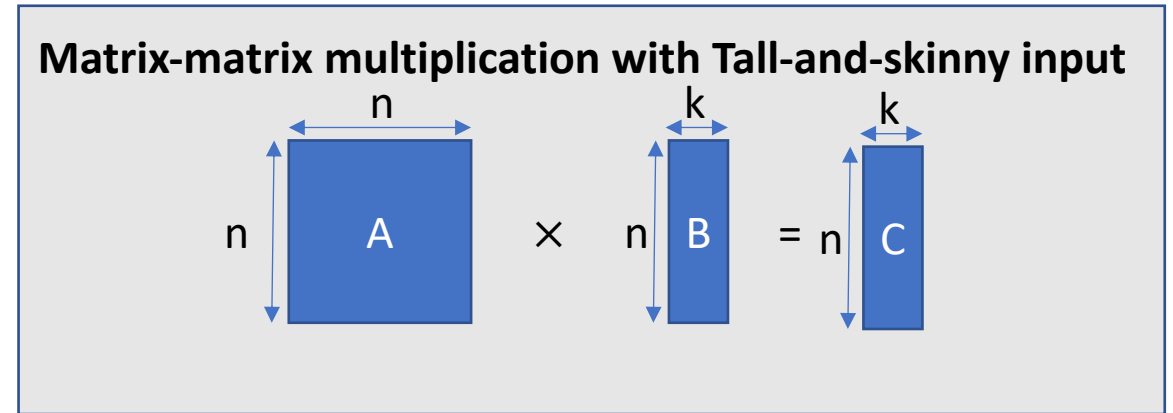
- Computation bound → Performance of application is bounded by **the computation power**.
- Memory bound → Performance of application is bounded by **the memory bandwidth**.



# Why tall-and-skinny behaves differently than regular shape input?



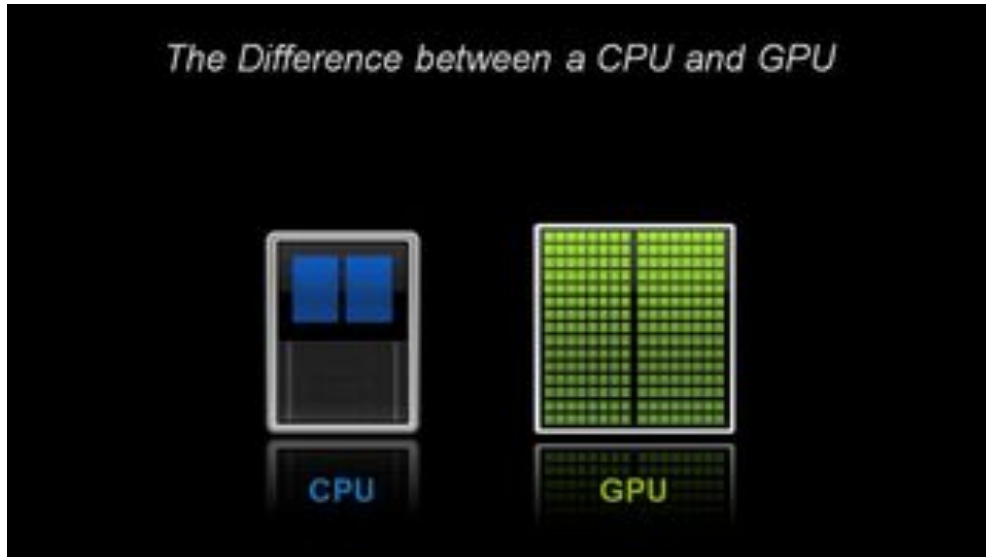
Input matrices size is  $O(n^2)$ .  
Computing time complexity is  $O(n^3)$ .  
Each element is used **n** times.



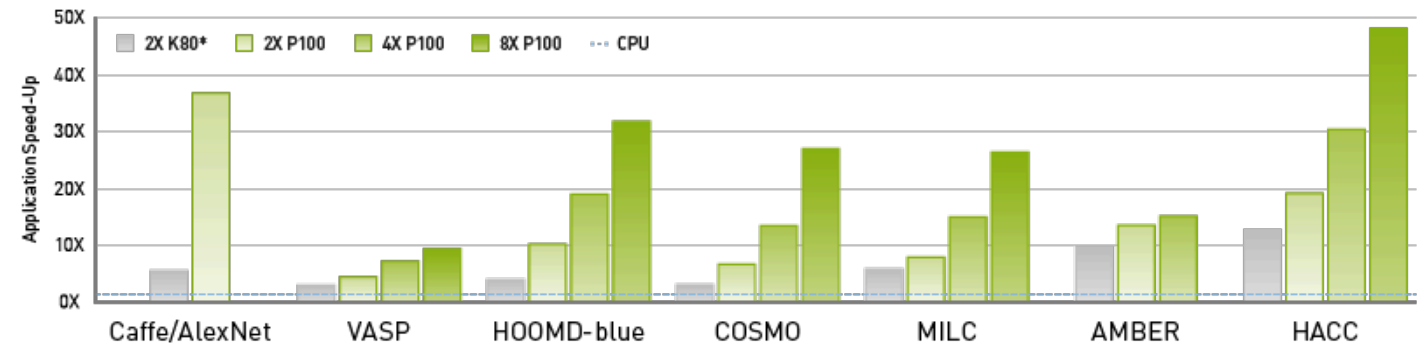
Input matrices size is  $O(n^2)$ .  
Computing time complexity is  $O(n^2k)$ .  
Each element is used **k** times on average.

- So for tall and skinny matrix input, depending on the k and the ratio between target GPU's peak computation power and peak memory throughput, it is usually memory bound.

# GPUs are widely used for accelerating applications



NVIDIA Tesla P100 Performance



Dual CPU server, Intel E5-2698 v3 @ 2.3 GHz, 256 GB system memory, Pre-production P100  
\* M40 for Caffe/AlexNet

- Good at parallelized computations.
- Higher computation power and memory throughput.
- Commonly used for accelerating matrix-related computations.

# cuBLAS library

---

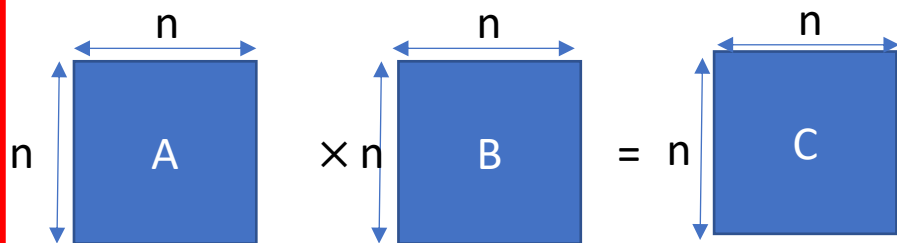


cuBLAS

- One of the most commonly used standard linear algebra libraries optimized for GPUs, which is developed by Nvidia.
- The core computing library of many big data and scientific computing applications.
- With deep optimization by Nvidia, the cuBLAS library is able to provide state-of-the-art performance in regular-shaped input matrix cases.
- But not fully optimized for tall-and-skinny matrix cases.

# Poor Performance on Current State-of-the-Art Design:

## Regular-sized matrix multiplication

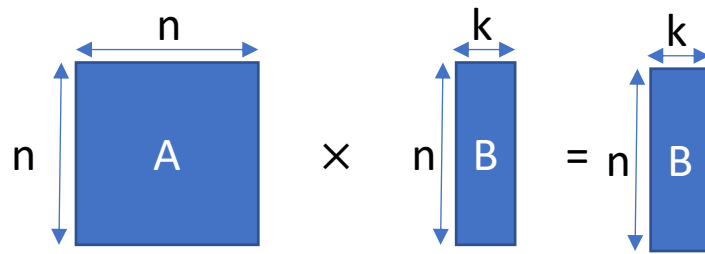


With large  $n$ ,  $k$  in similar magnitude

Computation bound

Regular size:  
80%-90% of  
the peak  
computation  
power

## Tall-and-skinny matrix multiplication



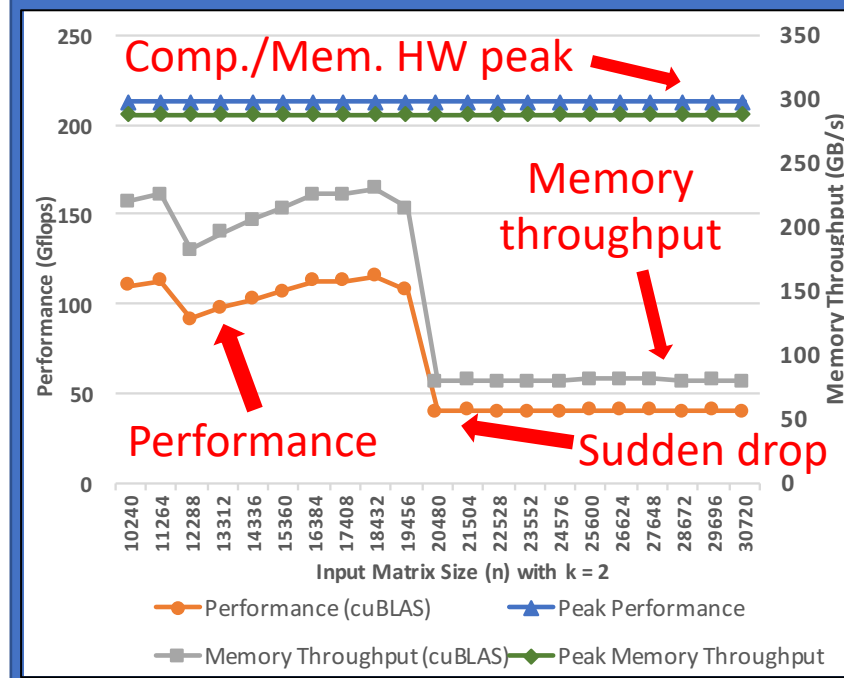
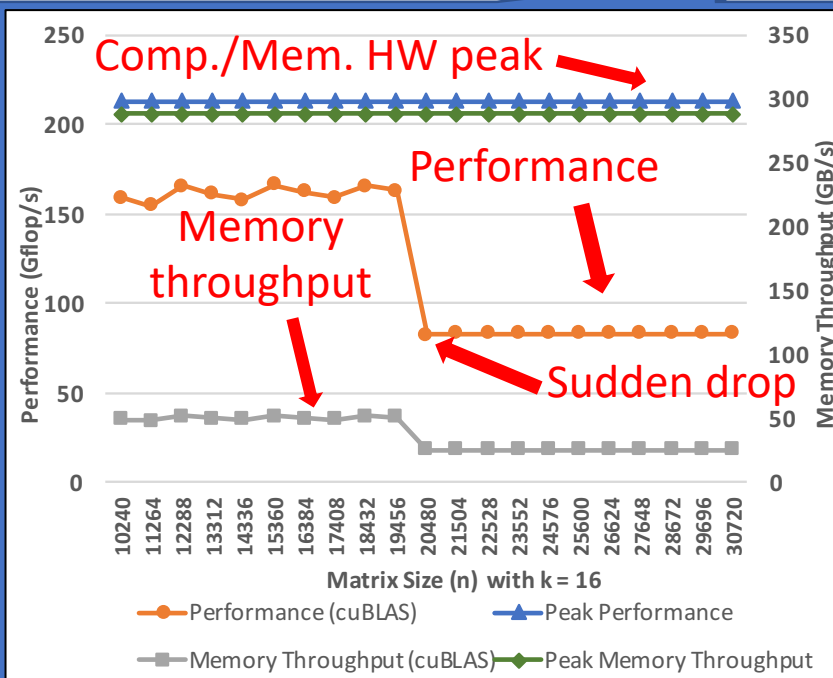
$n \gg k$

memory bound

Current state-of-the-art design only  
optimized for computation bound case

Low GPU utilization:

- $K = 2$ :
  - 49.9% memory band.
  - 37.9% peak comp. power
- $K=16$ :
  - 31.1% memory band.
  - 56.6% peak comp. power





## TSM2: redesigned matrix-matrix multiplication for tall-and-skinny input

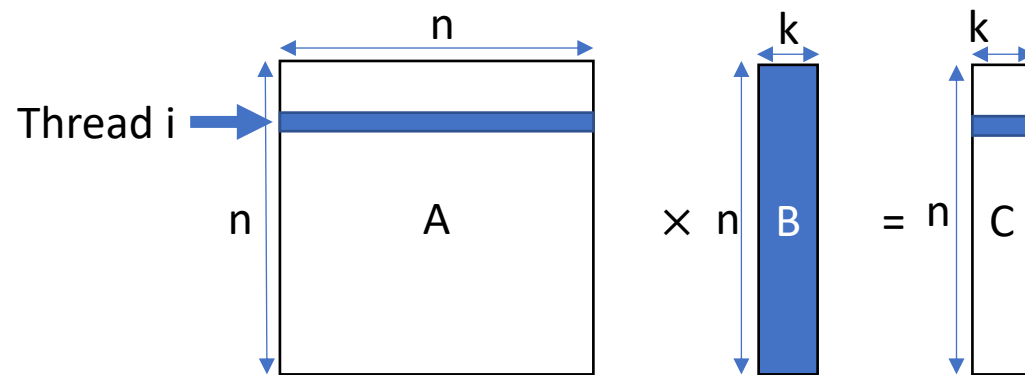
---

- Several factors are considered:
  - 1) Total number of global memory accesses.
  - 2) Efficiency on global memory throughput.
  - 3) Parallelism of overall workload.
  - 4) On-chip memory utilization.
  - 5) Streaming Multiprocessor (SM) utilization.

# Algorithm design: how to fit the workload into the programming model of CUDA(Continued)

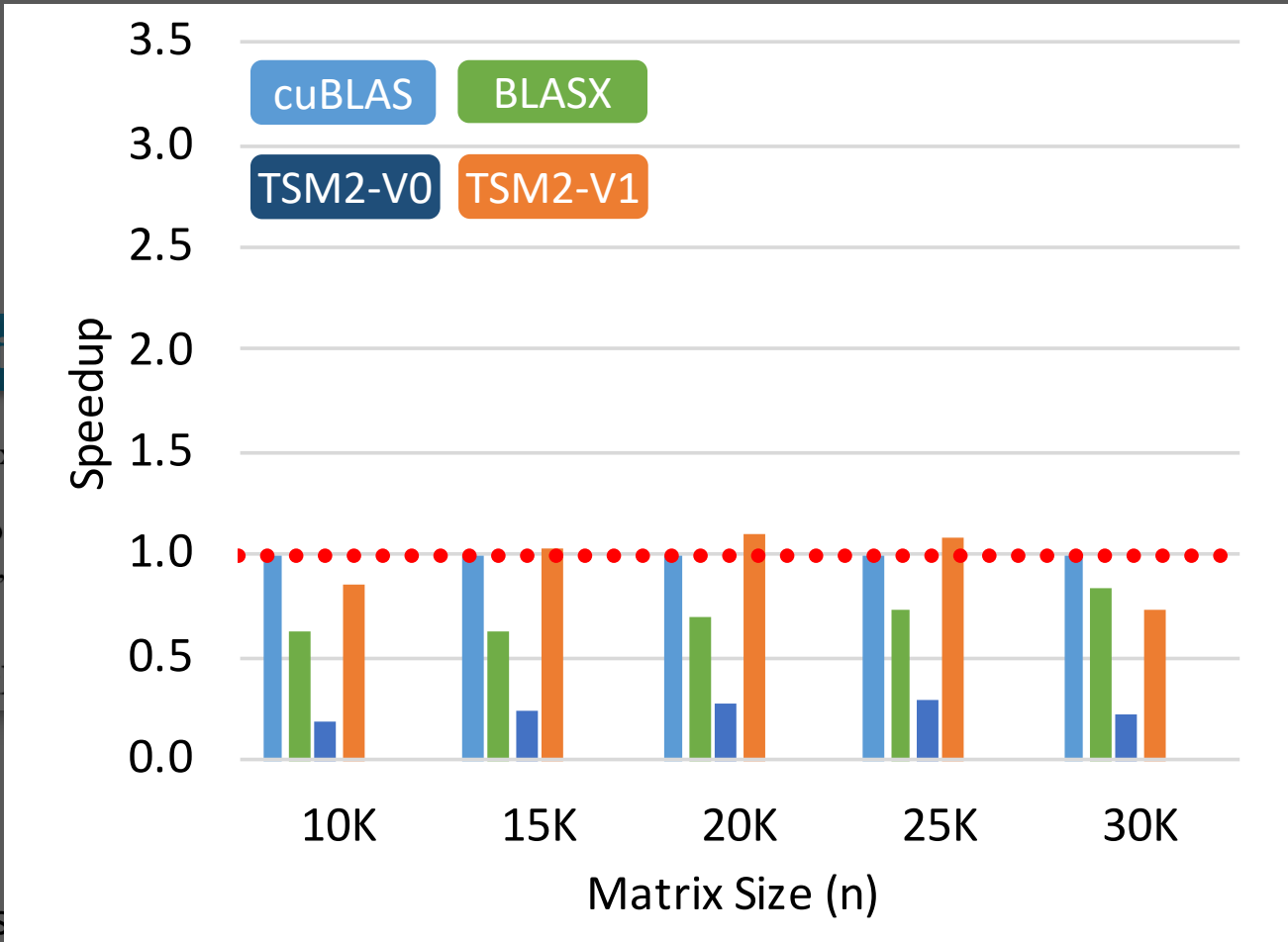
---

- We divide the workload by assigning  $n$  rows of matrix  $A$  to  $n$  different threads. Each vector-matrix multiplication is assigned to one thread.
  - i. To ensure high parallelism and high Streaming Multiprocessor occupancy.
  - ii. To ensure minimum number of memory access in favor of matrix  $A$ .
  - iii. To enable high memory accesses efficiency.



# Redesigning matrix-matrix multiplication for tall-and-skinny input

- Rethinking algorithm design – aiming to reduce total number of memory access
  - Inner product

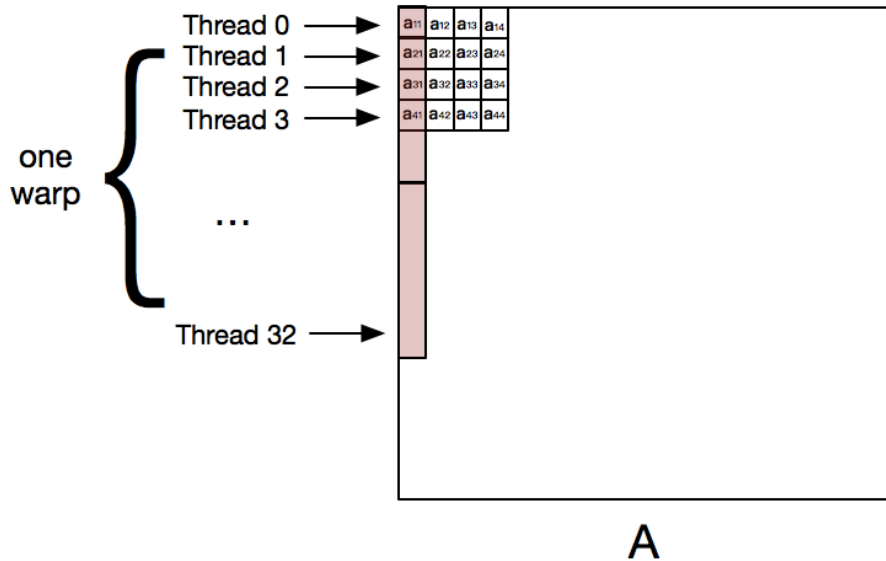


```
Version 1.0
Require: input matrix A (n x k)
Require: output matrix C (n x k)
1: for i = 1 to k do
2:   for j = 1 to n do
3:     C[thread_id, j] = 0
4:   end for
5: end for
Algorithm 1: Work
```

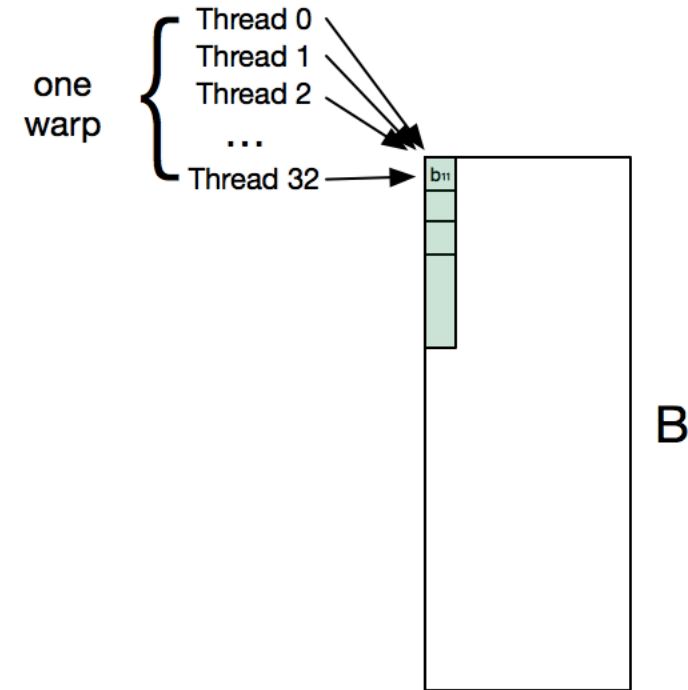
- Memory access to each element of A: 1 time
- Memory access to each element of B: n times
- Total number of accesses:  $2kn^2$
- Total number of accesses:  $(k+1)n^2$

# Global memory access efficiency analysis

- Global memory access efficiency per transaction = useful data/cache line size
  - Affect overall application memory access efficiency
  - Determined by the memory access pattern and the algorithm
  - Can be challenging to improve without modifying the algorithm design
- For outer product GEMM:



$$\frac{128 \text{ bytes}}{128 \text{ bytes}} = 100\% \text{ or } \frac{32 \text{ bytes}}{32 \text{ bytes}} = 100\%$$



$$\frac{8 \text{ bytes}}{128 \text{ bytes}} = 6.25\% \text{ or } \frac{8 \text{ bytes}}{32 \text{ bytes}} = 25\%$$

# Improving global memory access efficiency

## Version 2: Outer Product + Shared Mem.

**Require:** input matrix A ( $n \times n$ ) and B ( $n \times k$ )

**Require:** output matrix C ( $n \times k$ )

1:  $t_1 \leftarrow \text{tile\_size\_B}$ ,  $t_2 \leftarrow \text{tile\_size\_A}$

2: Register:  $A_1, A_2, \dots, A_{t_3}$

3: Register:  $C_1, C_2, \dots, C_{t_2}$

4: Shared Memory: currB with size  $t_1$

5: Threads per thread block  $\leftarrow t_1$

6: Total thread blocks  $\leftarrow n/t_1$

7: for  $p = 1$  to  $k$  with step size  $= t_2$  do

8:  $C_1 \leftarrow C[\text{thread\_id}, p]$

9:  $C_2 \leftarrow C[\text{thread\_id}, p+1]$

10:  $\ddot{C}_{t_2} \leftarrow C[\text{thread\_id}, p+t_2-1]$

11: for  $j = 0$  to  $n$  with step size  $= t_1$  do

12: *\* Load a tile of B into shared memory*

13: ThreadsSynchronization()

14:  $\text{currB}[\text{thread\_id}, 1] \leftarrow B[j+1, j+t_1+1]$

15:  $\text{currB}[\text{thread\_id}, 2] \leftarrow B[j+1, j+t_1+1]$

16:  $\dots$

17:  $\text{currB}[\text{thread\_id}, t_2] \leftarrow B[j+1, j+t_1+1]$

18: ThreadsSynchronization()

19: for  $l = j$  to  $j+t_1$  with step size  $= t_2$  do

20: *\* Load a tile of A into registers*

21:  $A_1 \leftarrow A[\text{thread\_id}, l]$

22:  $A_2 \leftarrow A[\text{thread\_id}, l+2]$

23:  $\dots$

24:  $\ddot{A}_{t_2} \leftarrow A[\text{thread\_id}, l+t_2-1]$

25:  $C_1 += A_{[1..t_3]} \times \text{currB}[[l, j+1, j+t_1+1]]$

26:  $C_2 += A_{[1..t_3]} \times \text{currB}[[l, j+1, j+t_1+1]]$

27:  $\dots$

28:  $\ddot{C}_{t_2} += A_{[1..t_3]} \times \text{currB}[[l, j+1, j+t_1+1]]$

29: end for

30: end for

31:  $C[\text{thread\_id}, p] \leftarrow C_1$

32:  $C[\text{thread\_id}, p+1] \leftarrow C_2$

33:  $\dots$

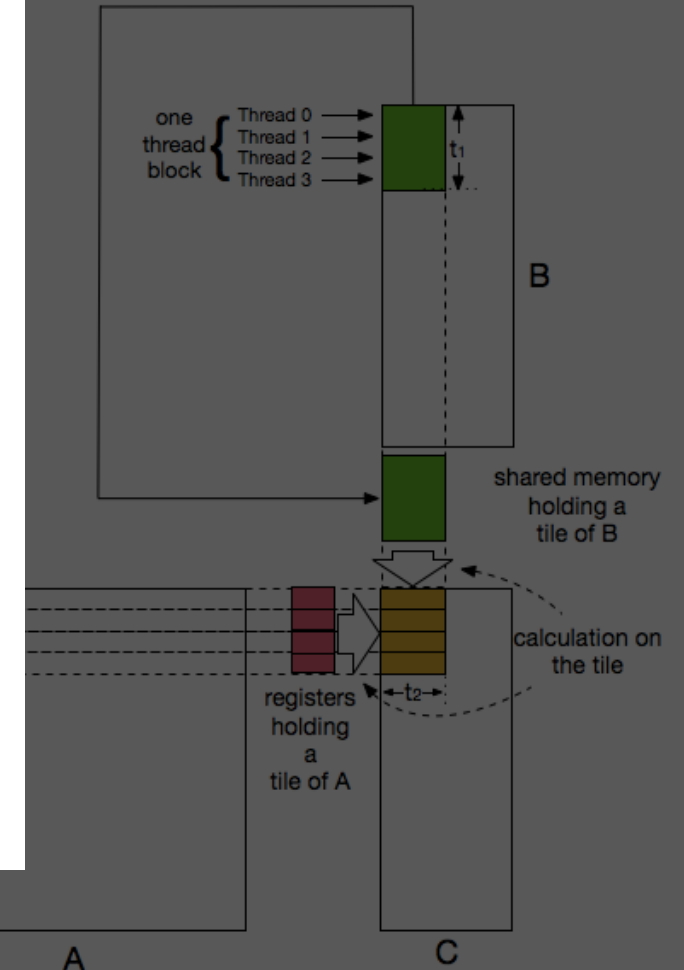
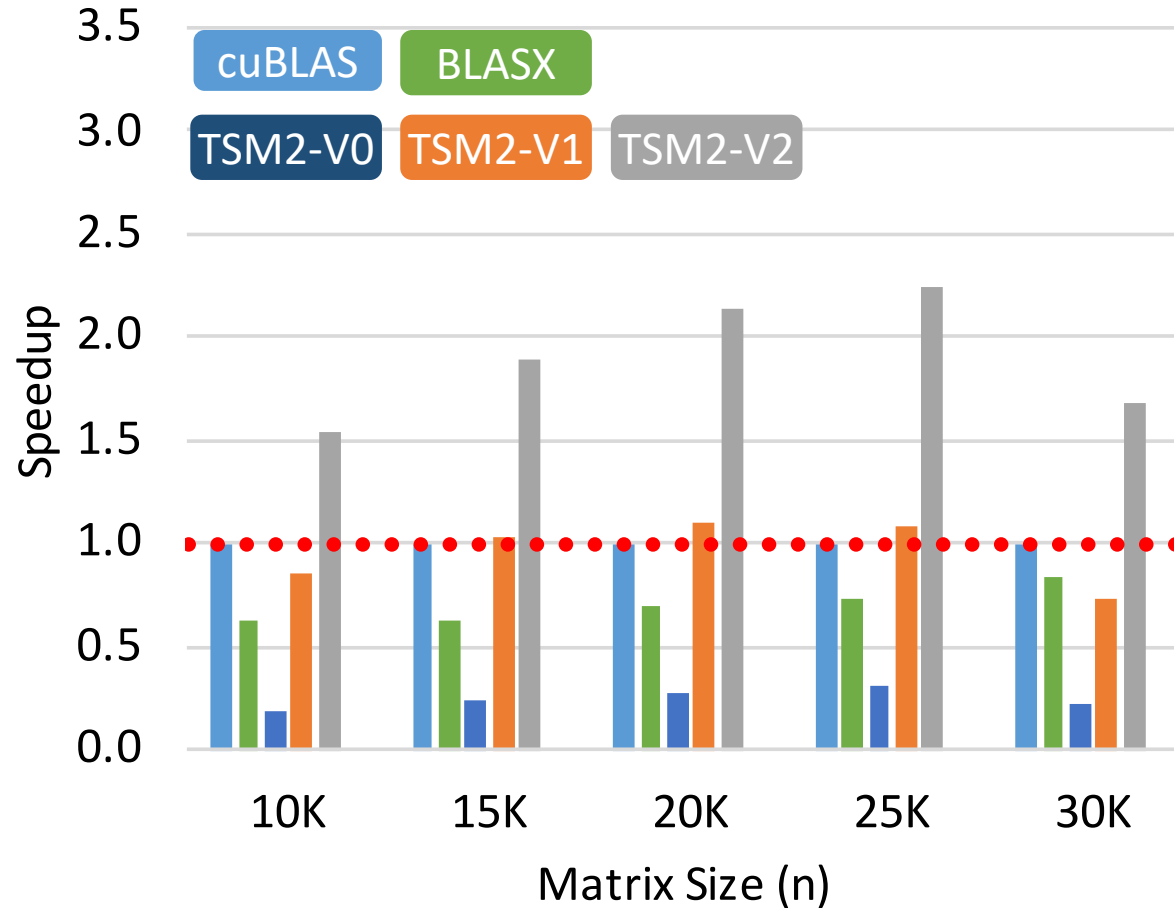
34:  $\ddot{C}[\text{thread\_id}, p+t_2] \leftarrow \ddot{C}_{t_2}$

35: end for

36: end for

**Algorithm 3:** TSM2 with shared memory

- GPU shared memory: sharing data between threads with threadblock
- Benefit: decoupling data load pattern and data use pattern.**



# Improving global memory access efficiency

## Version 2: Outer Product + Shared Mem.

**Require:** input matrix A ( $n \times n$ ) and B ( $n \times k$ )

**Require:** output matrix C ( $n \times k$ )

1:  $t_1 \leftarrow \text{tile\_size\_B}$ ,  $t_2 \leftarrow \text{tile\_size\_C}$ ,  $t_3 \leftarrow \text{tile\_size\_A}$

2: Register:  $A_1, A_2, \dots, A_{t_3}$

3: Register:  $C_1, C_2, \dots, C_{t_2}$

4: Shared Memory: currB with size  $t_1 \times t_2$

5: Threads per thread block  $\leftarrow t_1$

6: Total thread blocks  $\leftarrow n/t_1$

7: **for**  $p = 1$  to  $k$  with step size =  $t_2$  **do**

8:    $C_1 \leftarrow C[\text{thread\_id}, p]$

9:    $C_2 \leftarrow C[\text{thread\_id}, p + 1]$

10:    $\ddot{C}_{t_2} \leftarrow C[\text{thread\_id}, p + t_2 - 1]$

11:   **for**  $j = 0$  to  $n$  with step size =  $t_1$  **do**

*\* Load a tile of B into shared memory \**

    ThreadsSynchronization()

$\text{currB}[\text{thread\_id}, 1] \leftarrow B[j + \text{thread\_id}, p]$

$\text{currB}[\text{thread\_id}, 2] \leftarrow B[j + \text{thread\_id}, p + 1]$

$\dots$

$\text{currB}[\text{thread\_id}, t_2] \leftarrow B[j + \text{thread\_id}, p + t_2 - 1]$

    ThreadsSynchronization()

**for**  $l = j$  to  $j + t_1$  with step size =  $t_3$  **do**

*\* Load a tile of A into registers \**

$A_1 \leftarrow A[\text{thread\_id}, l]$

$A_2 \leftarrow A[\text{thread\_id}, l + 2]$

$\dots$

$A_{t_3} \leftarrow A[\text{thread\_id}, l + t_3 - 1]$

$C_1 + = A_{[1 \dots t_3]} \times \text{currB}[[l \dots l + t_3], 1]$

$C_2 + = A_{[1 \dots t_3]} \times \text{currB}[[l \dots l + t_3], 2]$

$\dots$

$\ddot{C}_{t_2} + = A_{[1 \dots t_3]} \times \text{currB}[[l \dots l + t_3], t_2]$

**end for**

**end for**

26:    $C[\text{thread\_id}, p] \leftarrow C_1$

27:    $C[\text{thread\_id}, p + 1] \leftarrow C_2$

28:    $\ddot{C}[\text{thread\_id}, p + t_2] \leftarrow \ddot{C}_{t_2}$

29: **end for**

**Algorithm 3:** TSM2 with shared memory

Load

Use

Data dependency between data load and data use instructions

- Even with efficient global memory loading pattern, it still brings high GPU underutilization
  - Main cause: long memory access latency can be hard to hide.

# Data prefetch: Improving GPU utilization

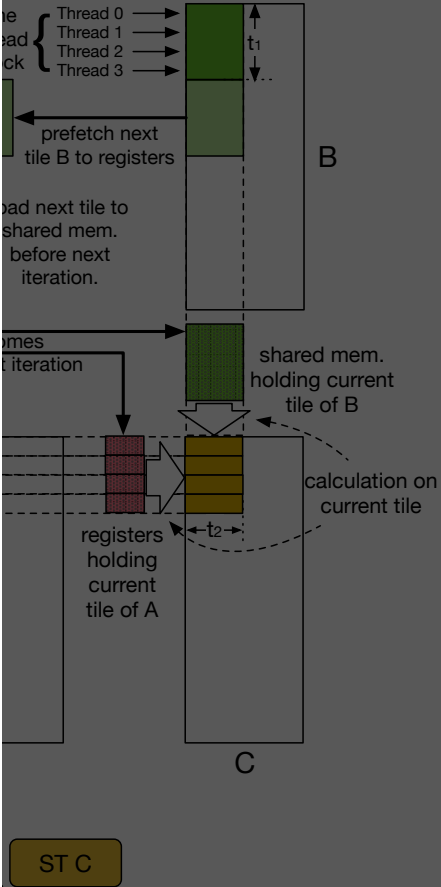
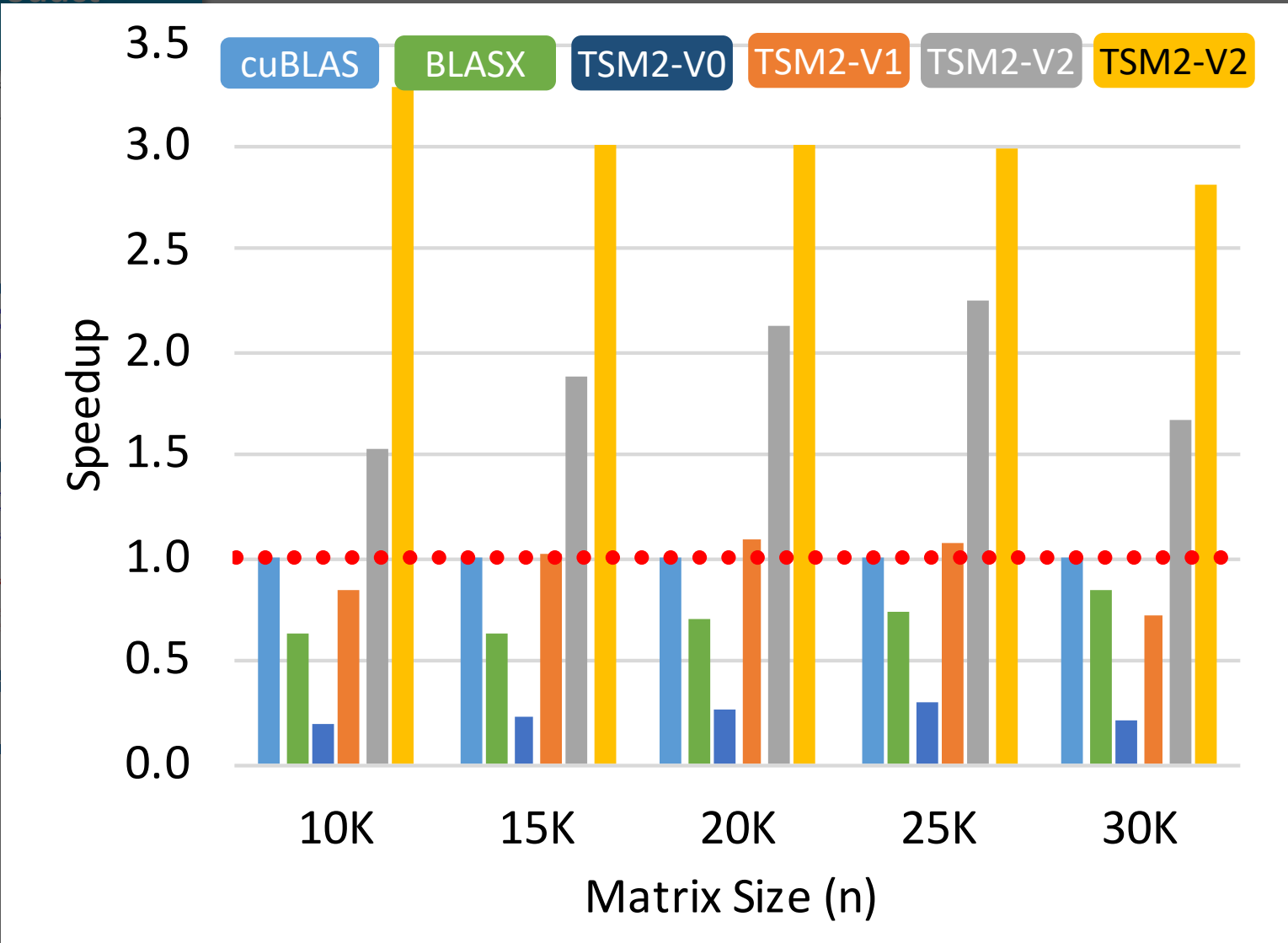
Version 3: Outer Product + Shared Mem. + Data Prefetch

Require: input matrix A ( $n \times n$ ) and B ( $n \times k$ )

Require: output matrix C ( $n \times k$ )

```
1:  $t_1 \leftarrow \text{tile\_size\_B}, t_2 \leftarrow \text{tile\_size\_C}, t_3 \leftarrow \text{tile\_size\_A}$ 
2: Register:  $\text{currA}_1, \text{currA}_2, \dots, \text{currA}_{t_3}$ 
3: Register:  $\text{nextA}_1, \text{nextA}_2, \dots, \text{nextA}_{t_3}$ 
4: Register:  $\text{nextB}_1, \text{nextB}_2, \dots, \text{nextB}_{t_2}$ 
5: Register:  $C_1, C_2, \dots, C_{t_2}$ 
6: Shared Memory:  $\text{currB}$  with size  $t_1 \times t_2$ 
7: Threads per thread block  $\leftarrow t_1$ 
8: Total thread blocks  $\leftarrow n/t_1$ 
9: for  $p = 1$  to  $k$  with step size  $= t_2$  do
10:    $C_1 \leftarrow C[\text{thread}, p]$ 
11:    $C_2 \leftarrow C[\text{thread}, p+1]$ 
12:    $C_{t_2} \leftarrow C[\text{thread}, p+t_2-1]$ 
13:    $\text{currB}[\text{thread\_id}, 1] \leftarrow B[\text{thread\_id}, 1]$ 
14:    $\text{currB}[\text{thread\_id}, 2] \leftarrow B[\text{thread\_id}, 2]$ 
15:   ...
16:    $\text{currB}[\text{thread\_id}, t_2] \leftarrow B[\text{thread\_id}, t_2]$ 
17:    $\text{currA}_1 \leftarrow A[\text{thread\_id}, 1]$ 
18:    $\text{currA}_2 \leftarrow A[\text{thread\_id}, 2]$ 
19:   ...
20:    $\text{currA}_{t_3} \leftarrow A[\text{thread\_id}, t_3]$ 
21:   for  $j = 0$  to  $n$  with step size  $= t_1$  do
22:     ThreadsSynchronization()
23:     /* prefetch the next tile of B into registers */
24:     if  $j+t_1 < n$  then
25:        $\text{nextB}_1 \leftarrow B[j+t_1+\text{thread\_id}, 1]$ 
26:        $\text{nextB}_2 \leftarrow B[j+t_1+\text{thread\_id}, 2]$ 
27:       ...
28:        $\text{nextB}_{t_2} \leftarrow B[j+t_1+\text{thread\_id}, t_2]$ 
29:     end if
30:     for  $l = j$  to  $j+t_1$  with step size  $= t_3$  do
31:       /* prefetch the next tile of A into registers */
32:       if  $l+t_3 < n$  then
33:          $\text{nextA}_1 \leftarrow A[\text{thread\_id}, l+t_3+1]$ 
34:          $\text{nextA}_2 \leftarrow A[\text{thread\_id}, l+t_3+2]$ 
35:         ...
36:          $\text{nextA}_{t_3} \leftarrow A[\text{thread\_id}, l+t_3+t_3]$ 
37:       end if
38:        $C_{t_2} += \text{currA}_{t_3} \times \text{currB}[\text{thread\_id}, t_2]$ 
39:       ...
40:        $C_2 += \text{currA}_{t_3} \times \text{currB}[\text{thread\_id}, 2]$ 
41:        $C_1 += \text{currA}_{t_3} \times \text{currB}[\text{thread\_id}, 1]$ 
42:     end for
43:     ThreadsSynchronization()
44:      $\text{currB}[\text{thread\_id}, 1] \leftarrow \text{nextB}_1$ 
45:      $\text{currB}[\text{thread\_id}, 2] \leftarrow \text{nextB}_2$ 
46:     ...
47:      $\text{currB}[\text{thread\_id}, t_2] \leftarrow \text{nextB}_{t_2}$ 
48:   end for
49:    $C[\text{thread}, p] \leftarrow C_1$ 
50:    $C[\text{thread}, p+1] \leftarrow C_2$ 
51:   ...
52:    $C[\text{thread}, p+t_2] \leftarrow C_{t_2}$ 
53: end for
```

Algorithm 4: TSM2 with shared memory and data prefetching



Tall-and-skinny GEMM with K=8 on Nvidia Tesla K40c

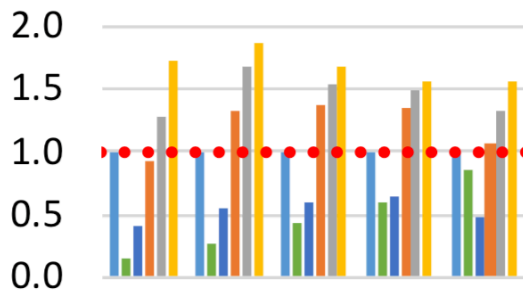
## Experimental evaluation:

---

GPU Model	Micro-architectures	Memory	Peak performance	Peak memory bandwidth
Tesla K40c	Kepler	12 GB	1430 GFLOPS	288 GB/s
Tesla M40	Maxwell	24 GB	213 GFLOPS	288 GB/s
Tesla P100	Pascal	16 GB	4600 GFLOPS	720 GB/s



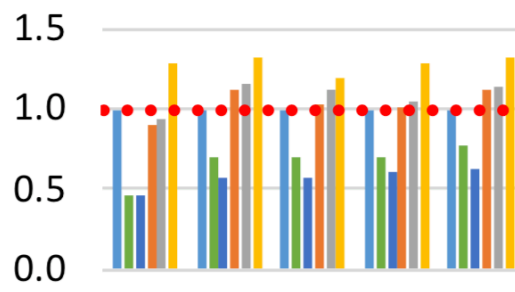
# Experimental evaluation: Speedup (on Nvidia Tesla K40c)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

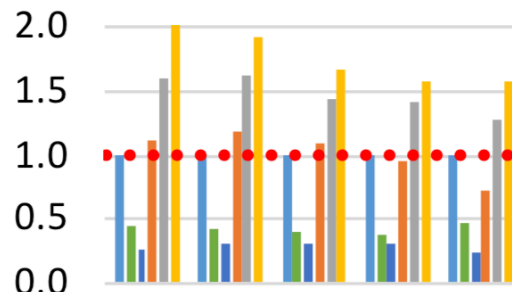
(a) Single precision (k=2)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

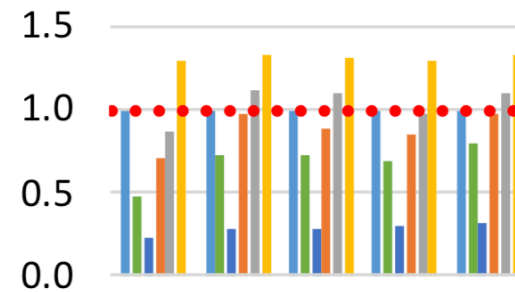
(b) Double precision (k=2)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

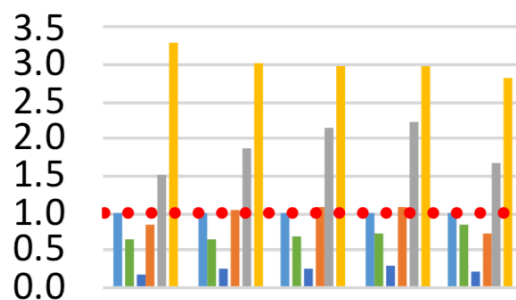
(c) Single precision (k=4)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

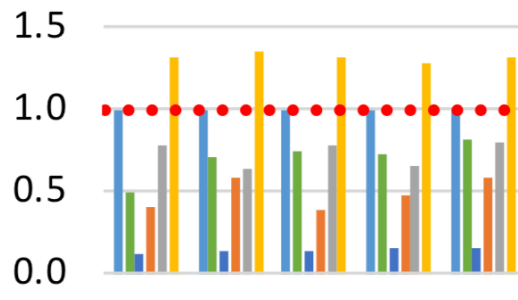
(d) Double precision (k=4)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

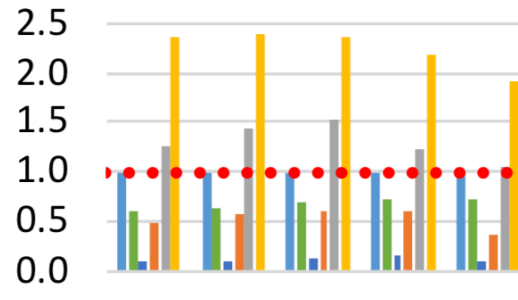
(e) Single precision (k=8)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

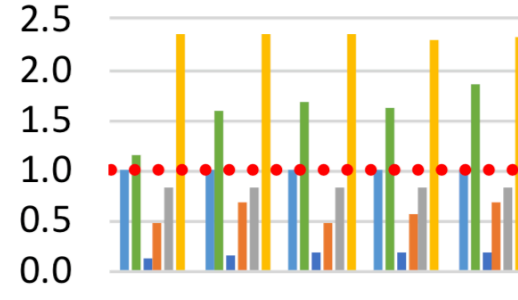
(f) Double precision (k=8)



10K 15K 20K 25K 30K  
Matrix Size (n)

cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

(g) Single precision (k=16)

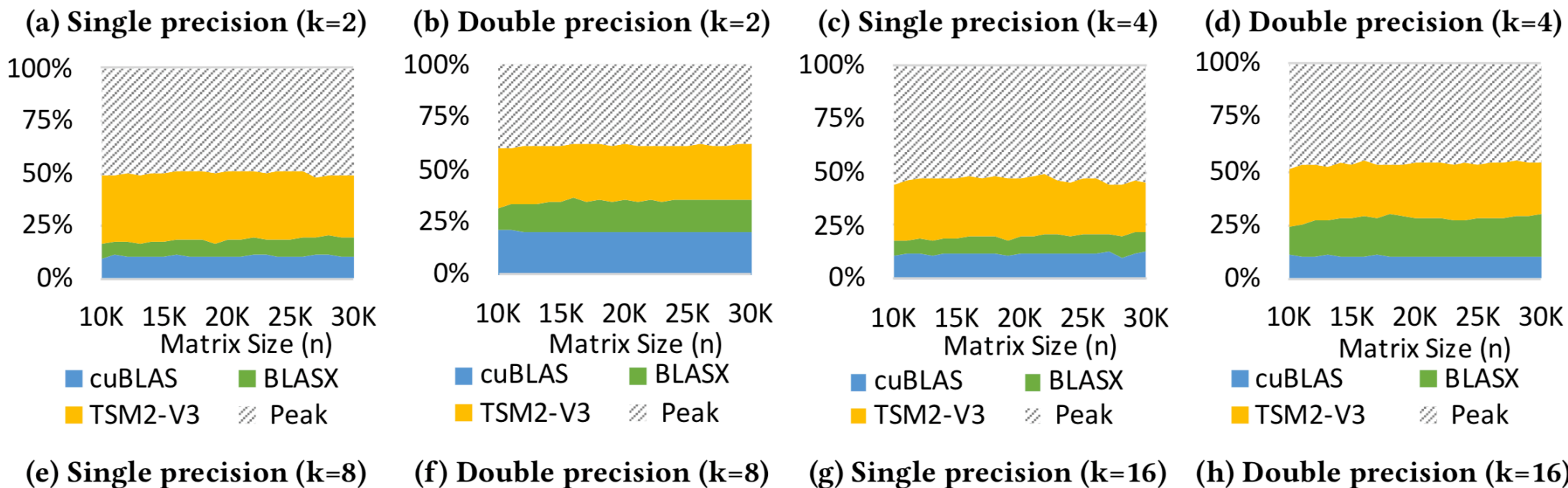
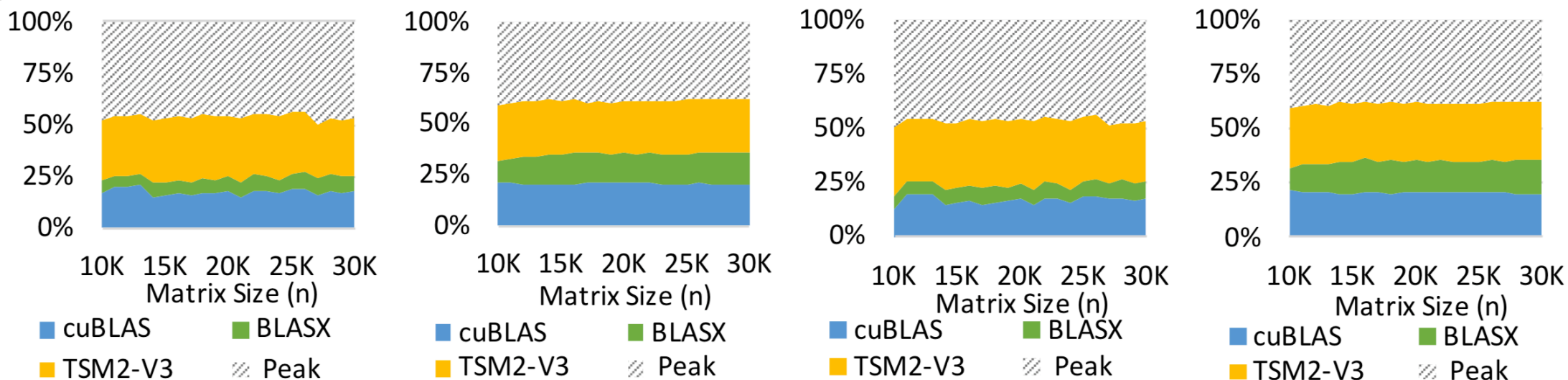


10K 15K 20K 25K 30K  
Matrix Size (n)

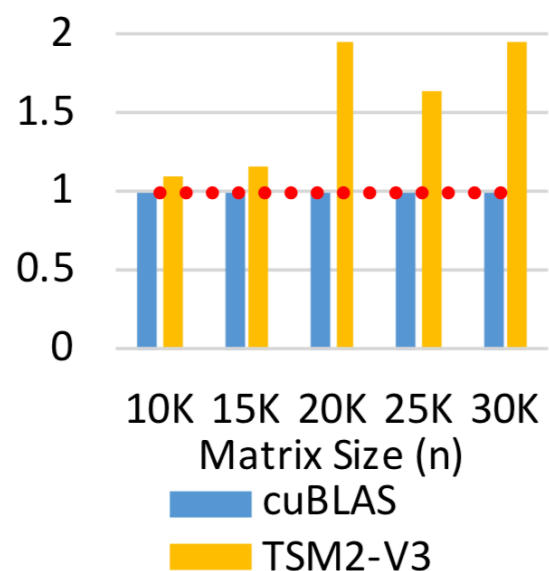
cuBLAS BLASX  
TSM2-V0 TSM2-V1  
TSM2-V2 TSM2-V3

(h) Double precision (k=16)

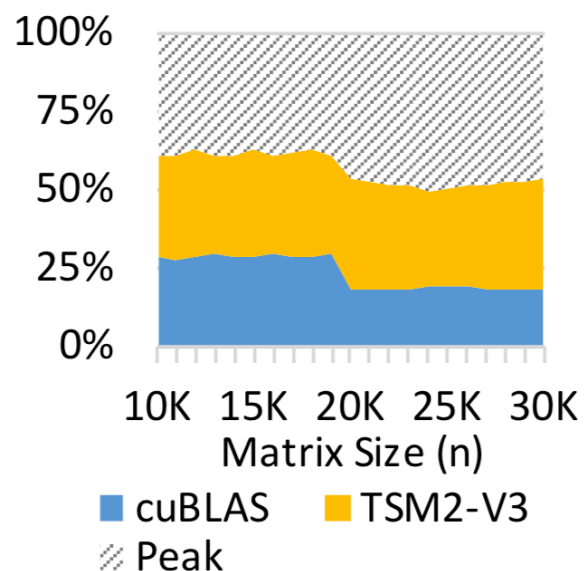
# Experimental evaluation: Memory bandwidth (on Nvidia Tesla K40c)



# Experimental evaluation on Nvidia Tesla M40 and P100

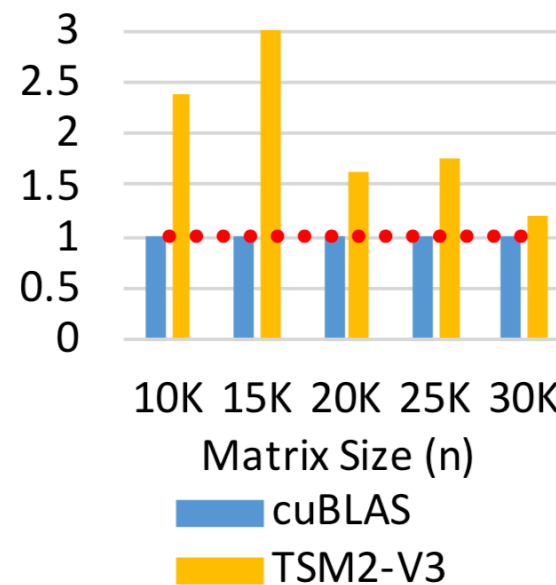


(a) Speedup

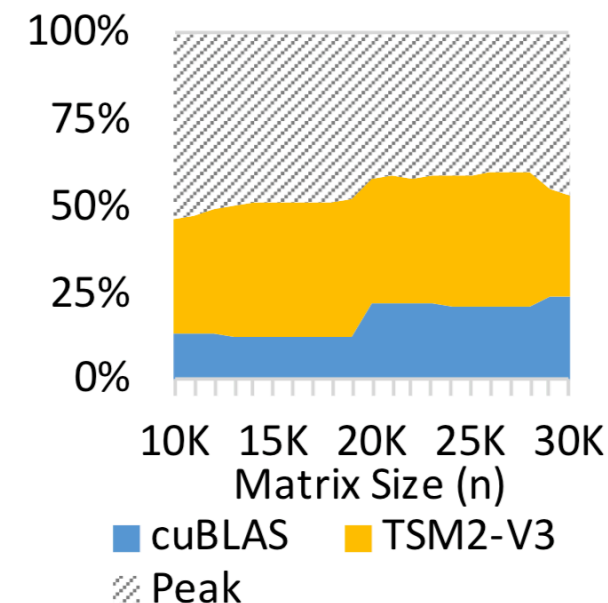


(b) Computational power util.

Tesla M40



(a) Speedup



(b) Memory throughput util.

Tesla P100

# Showcase 1: K-means

Core computation of Lloyd's K-means: **distance calculation**.  
Common choice: **Euclidean Distance**

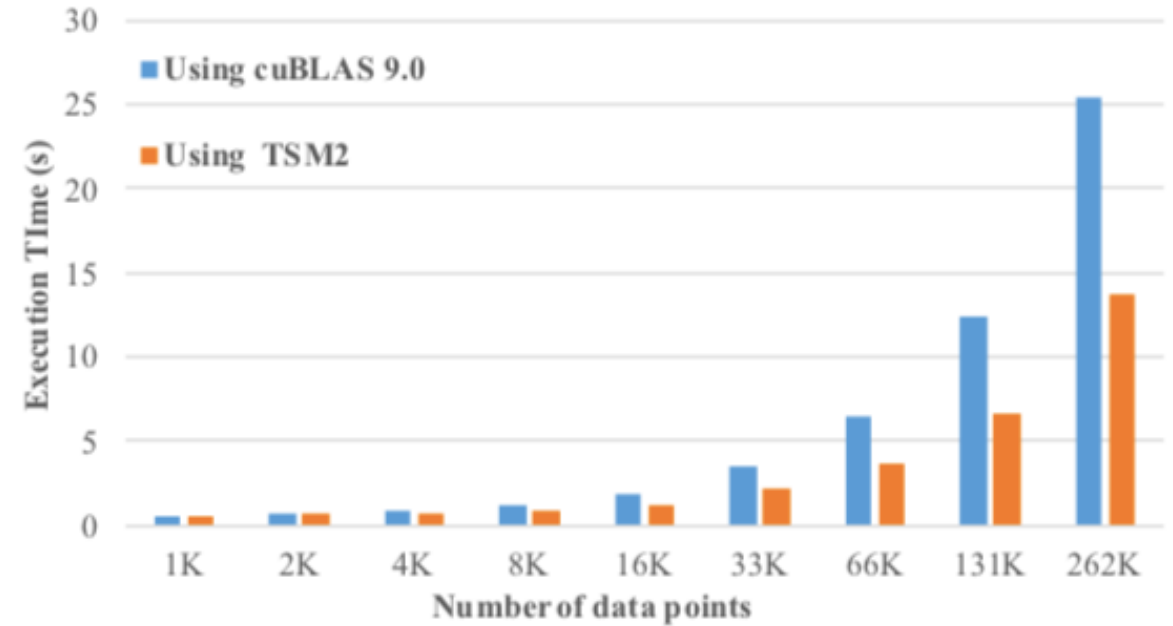
$$||x - y||^2 = ||x||^2 + ||y||^2 - 2xy$$

When we have multiple x and y:

Group x → matrix X } calculating xy → XY  
Group y → matrix Y } (matrix matrix multiplication)

Calculating distance between:

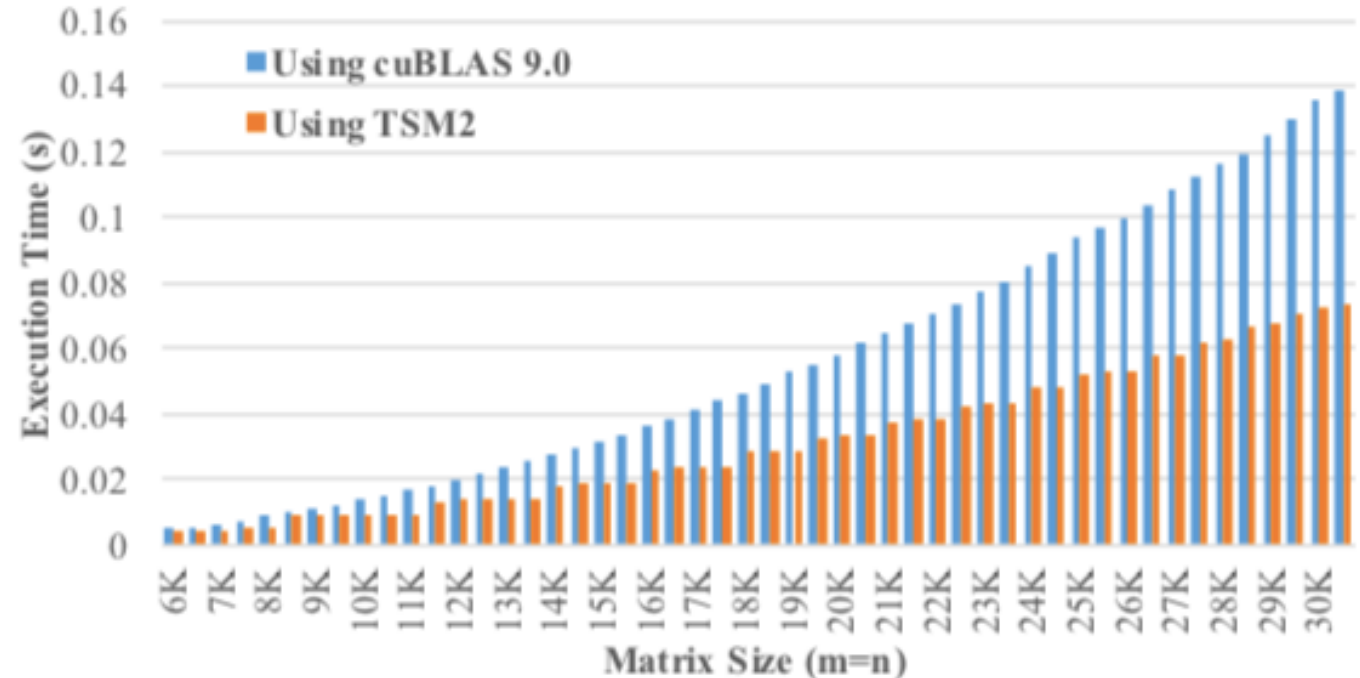
- Data points X (n points with d dimensions);
- Centroids C (k centroids with d dimensions);
- → matrix-matrix multiplication: (n\*d) times (d\*k).
- Usually  $k \ll n, d \rightarrow$  tall-and-skinny



- Execution time of the first 100 iterations of Lloyd's K-means algorithm on K40c (d = 4096, k = 16).
- Using our TSM2, we speedup K-means by 1.06x - 1.89x (**avg. 1.53x**).
- GPU version K-means originally developed by NVIDIA:  
<https://github.com/NVIDIA/kmeans>

## Showcase 2: ABFT Matrix Checksum Encoding

- Core computation of ABFT: **calculating checksum** (*encode redundant info*)
- E.g., calculate the checksum of matrix A with checksum weight vector v:  
$$\text{checksum}(A) = Av$$
- Usually use multiple different checksum weight vectors.
- If we use c different checksum weight vectors  $\rightarrow$  (m-by-n) times (n-by-c)
- **Common choice:  $c = 2 \ll m, n \rightarrow$  tall-and-skinny**



We compare the checksum encoding performance by using cuBLAS and TSM2 on K40c. As we can see, our TSM2 significantly improve the checksum encoding calculation with 1.10x to 1.90x speedup (**avg. 1.67x**).

## Conclusion:

---

- We first analyzed the performance of current GEMM in the latest cuBLAS library.
- We discovered the potential challenges of optimizing tall-and-skinny GEMM since its workload is memory bound.
- We redesigned an optimized tall-and-skinny GEMM with several optimization techniques focusing on GPU resource utilization.
- Experiment results show that our optimized implementation can achieve better performance on three modern GPU micro-architectures.



# We have an optimized design, but when do we use it?

**Algorithm 4** Tall-and-skinny matrix-matrix multiplication with data prefetching

**Require:** input matrix A ( $n \times n$ ) and B ( $n \times k$ )

**Require:** output matrix C ( $n \times k$ )

```

1:  $t_1 \leftarrow \text{tile\_size\_B}$ ,  $t_2 \leftarrow \text{tile\_size\_C}$ ,  $t_3 \leftarrow \text{tile\_size\_A}$ 
2: Register:  $\text{currA}_1, \text{currA}_2, \dots, \text{currA}_{t_1}$ 
3: Register:  $\text{nextA}_1, \text{nextA}_2, \dots, \text{nextA}_{t_1}$ 
4: Register:  $\text{nextB}_1, \text{nextB}_2, \dots, \text{nextB}_{t_2}$ 
5: Register:  $C_1, C_2, \dots, C_{t_2}$ 
6: Shared Memory:  $\text{currB}$  with size  $t_1 \times t_2$ 
7: Threads per thread block  $\leftarrow t_1$ 
8: Total thread blocks  $\leftarrow n/t_1$ 
9: for  $p = 1$  to  $k$  with step size  $= t_2$  do
10:    $C_1 \leftarrow C[\text{thread}, p]$ 
11:    $C_2 \leftarrow C[\text{thread}, p+1]$ 
12:   ...
13:    $\text{currB}[\text{thread\_id}, 1] \leftarrow B[\text{thread\_id}, p]$ 
14:    $\text{currB}[\text{thread\_id}, 2] \leftarrow B[\text{thread\_id}, p+1]$ 
15:   ...
16:    $\text{currB}[\text{thread\_id}, t_2] \leftarrow B[\text{thread\_id}, p+t_2-1]$ 
17:    $\text{currA}_1 \leftarrow A[\text{thread\_id}, 1]$ 
18:    $\text{currA}_2 \leftarrow A[\text{thread\_id}, 2]$ 
19:   ...
20:    $\text{currA}_{t_1} \leftarrow A[\text{thread\_id}, t_1]$ 
21:   for  $j = 0$  to  $n$  with step size  $= t_1$  do
22:     ThreadsSynchronization()
23:     /* prefetch the next tile of B into registers */
24:     if  $j+t_1 < n$  then
25:        $\text{nextB}_1 \leftarrow B[j+t_1+\text{thread\_id}, p]$ 
26:        $\text{nextB}_2 \leftarrow B[j+t_1+\text{thread\_id}, p+1]$ 
27:       ...
28:        $\text{nextB}_{t_2} \leftarrow B[j+t_1+\text{thread\_id}, p+t_2-1]$ 
29:     end if
30:     for  $l = j$  to  $j+t_1$  with step size  $= t_3$  do
31:       /* prefetch the next tile of A into registers */
32:       if  $l+t_3 < n$  then
33:          $\text{nextA}_1 \leftarrow A[\text{thread\_id}, l+t_3]$ 
34:          $\text{nextA}_2 \leftarrow A[\text{thread\_id}, l+t_3+1]$ 
35:         ...
36:          $\text{nextA}_{t_1} \leftarrow A[\text{thread\_id}, l+t_3+t_1-1]$ 
37:       end if
38:        $C_1 += \text{currA}_{[1..t_1]} \times \text{currB}[[l..l+t_3], 1]$ 
39:        $C_2 += \text{currA}_{[1..t_1]} \times \text{currB}[[l..l+t_3], 2]$ 
40:       ...
41:        $C_{t_2} += \text{currA}_{[1..t_1]} \times \text{currB}[[l..l+t_3], t_2]$ 
42:        $\text{currA}_1 \leftarrow \text{nextA}_1$ 
43:        $\text{currA}_2 \leftarrow \text{nextA}_2$ 
44:       ...
45:        $\text{currA}_{t_1} \leftarrow \text{nextA}_{t_1}$ 
46:     end for
47:     ThreadsSynchronization()
48:      $\text{currB}[\text{thread\_id}, 1] \leftarrow \text{nextB}_1$ 
49:      $\text{currB}[\text{thread\_id}, 2] \leftarrow \text{nextB}_2$ 
50:     ...
51:      $\text{currB}[\text{thread\_id}, t_2] \leftarrow \text{nextB}_{t_2}$ 
52:   end for
53:    $C[\text{thread}, p] \leftarrow C_1$ 
54:    $C[\text{thread}, p+1] \leftarrow C_2$ 
55:   ...
56:    $C[\text{thread}, p+t_2] \leftarrow C_{t_2}$ 
57: end for

```

How to determine when the computation is memory bound and when it is not?

GPU Peak Perf.

GPU Peak Mem. Band.

Hardware parameters

**Algorithm 5** Parameter Optimization for Tall-and-skinny matrix-matrix multiplication

```

1: if  $k \leq t_2^{\text{threshold}}$  then
2:    $\text{Total\_memory} \approx n \times n \times \frac{k}{t_2} \times \text{bytes\_per\_elem.}$ 
3:    $\text{Bandwidth} = \text{PeakBand.} \times \text{Util}_{\text{mem}}$ 
4:   Use Gradient Descent to Optimize ( $t_2$  and  $t_3$ ):  $\text{Time} = \frac{\text{Total\_memory}}{\text{Bandwidth}}$  with  $1 \leq t_2 \leq k$  and  $1 \leq t_3$ 
5:   Output:  $t_2$  and  $t_3$ 
6: else
7:    $\text{Total\_flops} = n \times n \times k \times 2$ 
8:    $\text{Compute\_power} = \text{PeakPerf.} \times \text{Util}_{\text{comp}}$ 
9:   Use Gradient Descent to Optimize ( $t_2$  and  $t_3$ ):  $\text{Time}_1 = \frac{\text{Total\_flops}}{\text{Compute\_power}}$  with  $t_2^{\text{threshold}} \leq t_2 \leq k$  and  $1 \leq t_3$ 
10:   $t_{2(\text{time}_1)} \leftarrow t_2$ 
11:   $t_{3(\text{time}_1)} \leftarrow t_3$ 
12:   $\text{Total\_memory} \approx n \times n \times \frac{k}{t_2} \times \text{bytes\_per\_elem.}$ 
13:   $\text{Bandwidth} = \text{PeakBand.} \times \text{Util}_{\text{mem}}$ 
14:  Use Gradient Descent to Optimize ( $t_2$  and  $t_3$ ) in  $\text{Time}_2 = \frac{\text{Total\_memory}}{\text{Bandwidth}}$  with  $1 \leq t_2 \leq t_{2(\text{time}_1)}$  and  $1 \leq t_3$ 
15:   $t_{2(\text{time}_2)} \leftarrow t_2$ 
16:   $t_{3(\text{time}_2)} \leftarrow t_3$ 
17:  if  $\text{Time}_1 < \text{Time}_2$  then
18:    Output:  $t_{2(\text{time}_1)}$  and  $t_{3(\text{time}_1)}$ 
19:  else
20:    Output:  $t_{2(\text{time}_2)}$  and  $t_{3(\text{time}_2)}$ 
21:  end if
22: end if

```

Tuning parameters

NVIDIA Tesla K40:

K=40

Computation bound

memory bound

NVIDIA Tesla M40:

K=6

Computation bound

memory bound

NVIDIA Tesla P100:

K=50

Computation bound

memory bound