

Efficient Quantized Sparse Matrix Operations on Tensor Cores

Shigang Li^{*}, Kazuki Osawa⁺, Torsten Hoefler⁺

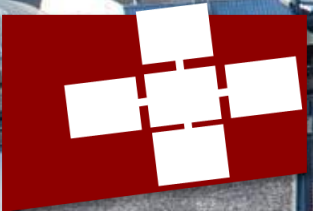
^{*}School of Computer Science, Beijing University of Posts and Telecommunications

⁺Department of Computer Science, ETH Zurich



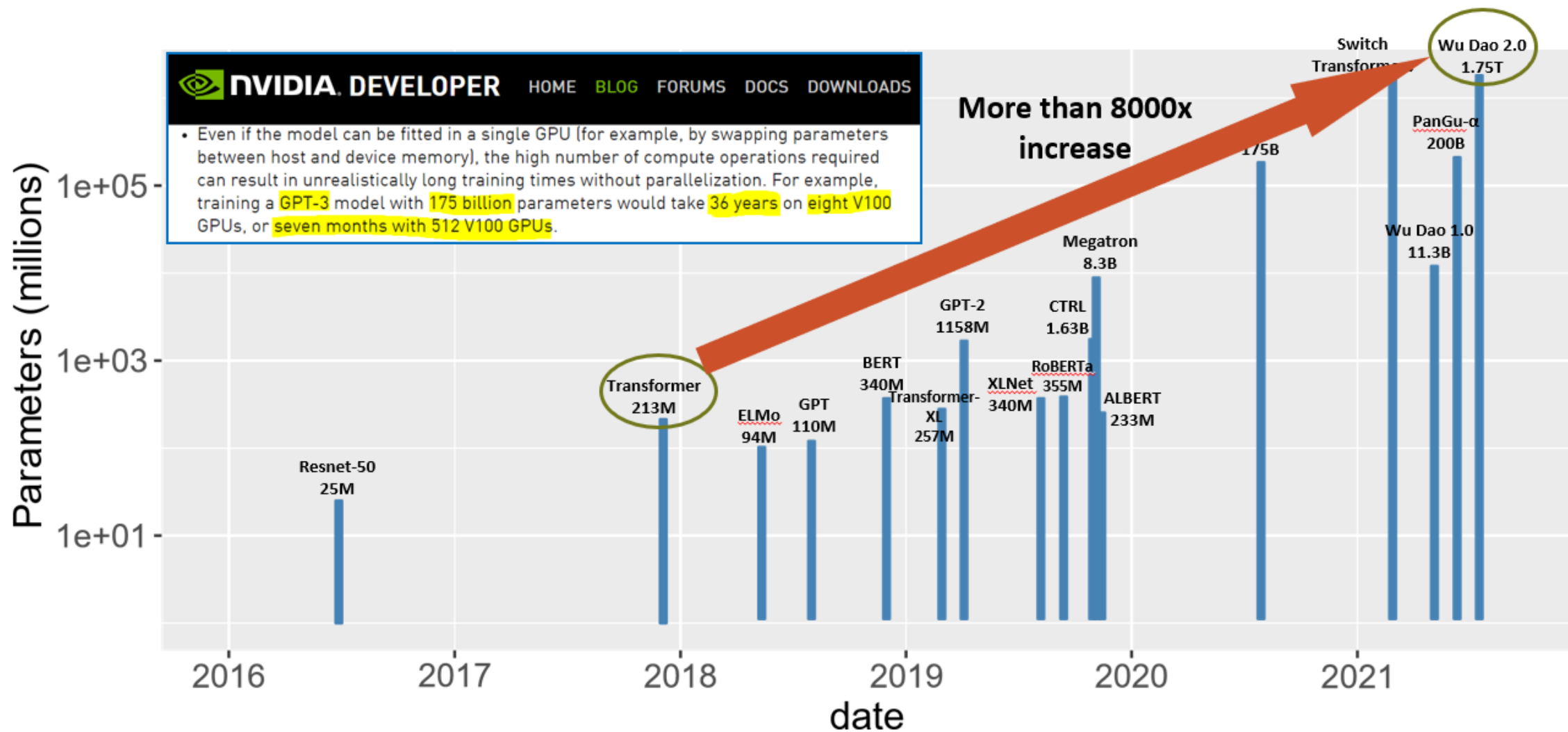
Efficient Quantized Sparse Matrix Operations on Tensor Cores

Shigang Li, Kazuki Osawa, Torsten Hoefler



SC22, Dallas, TX, USA
Nov. 2022

Model size is growing exponentially

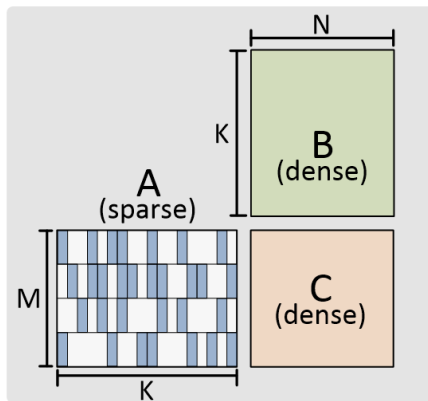


Models are also compressible

Sparsification

SpMM

1. Self-attention in sparse Transformers
2. Forward pass of pruned models
- ...



Sparsity in scientific:

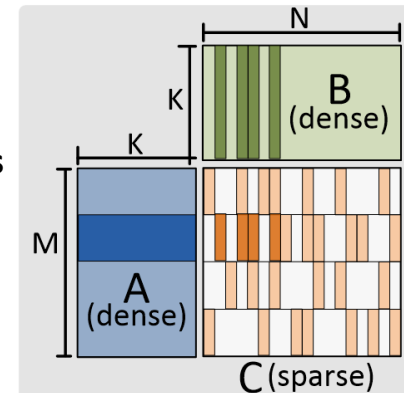
> 99%



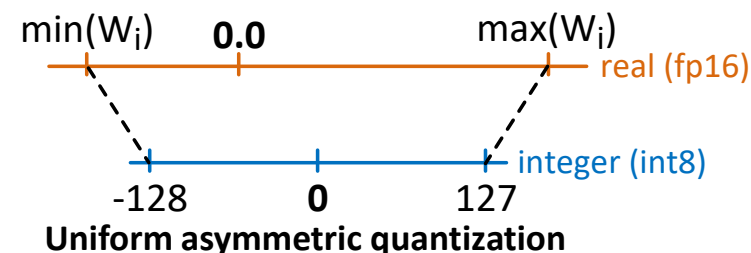
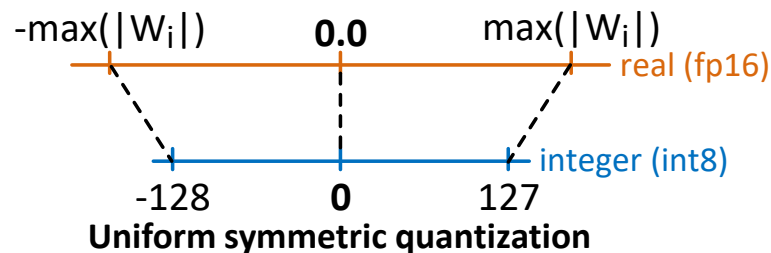
Sparsity in DL:
50% ~ 90%

SDDMM

1. Attention score in sparse Transformers
2. Backward pass of pruned models
- ...



Quantization



Combining sparsification with quantization

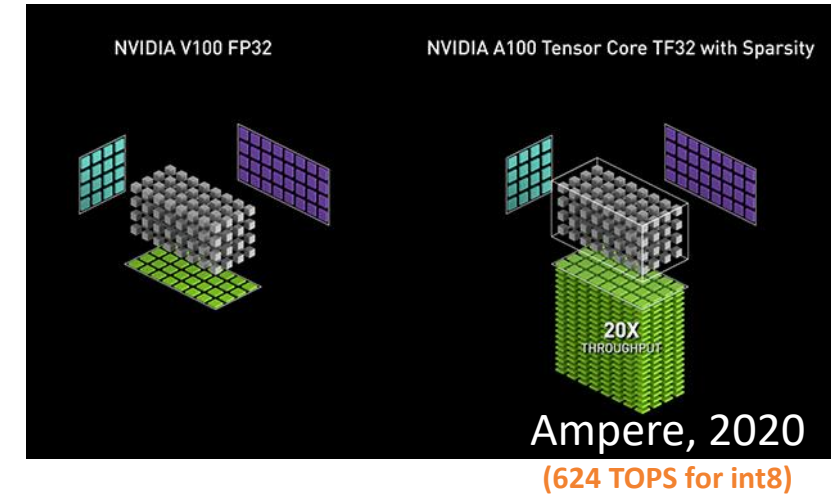
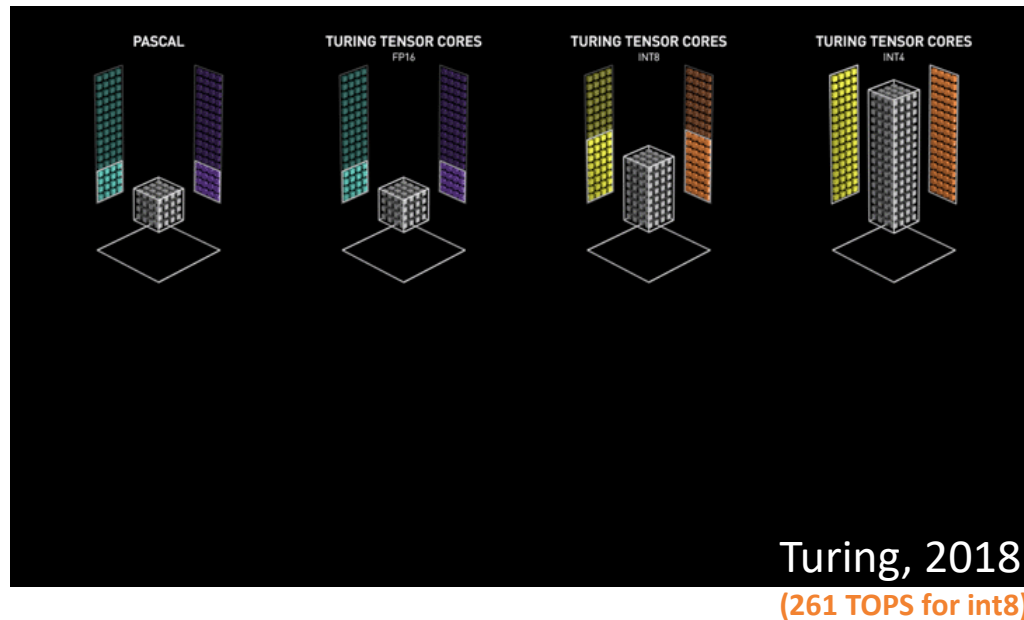
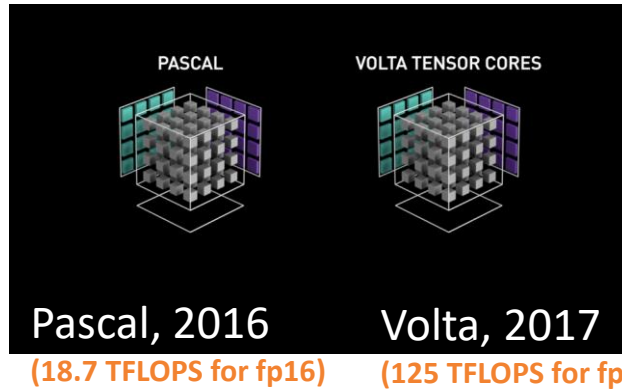
Mart van Baalen et al., Bayesian bits: Unifying quantization and pruning, **NeurIPS 2020**

H. Yang et al., Automatic neural network compression by sparsity-quantization joint learning: A constrained optimization based approach, **CVPR 2020**

S. Han et al., Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, **ICLR 2016**

...

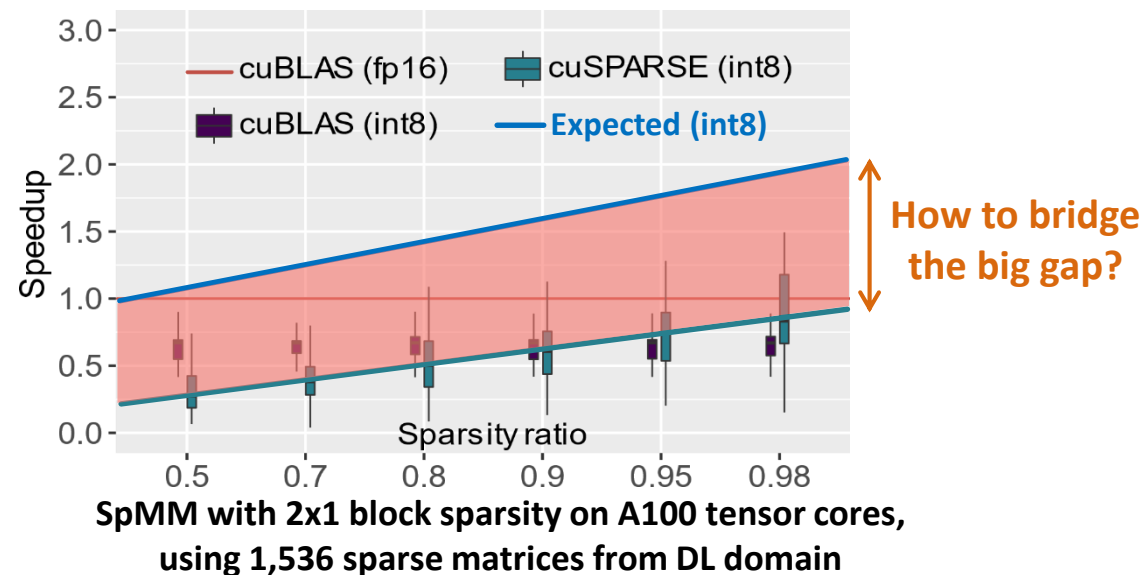
Tensor cores for deep learning acceleration



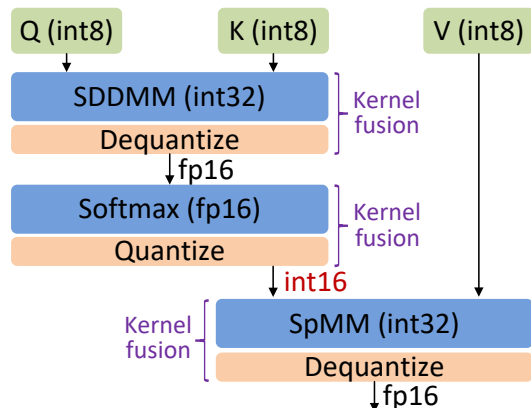
Images and GIFs in this slide are from <https://www.nvidia.com/en-us/data-center/tensor-cores/>

Challenges

(1) How to achieve practical speedup in a large range of sparsity ratio, e.g., **50% ~ 98%**?



(2) How to efficiently support sparse workloads with **mixed precision (two input matrices with different precision)**, e.g., **8-bit weights and 4-bit activation**?



Sparse self-attention with mixed precision

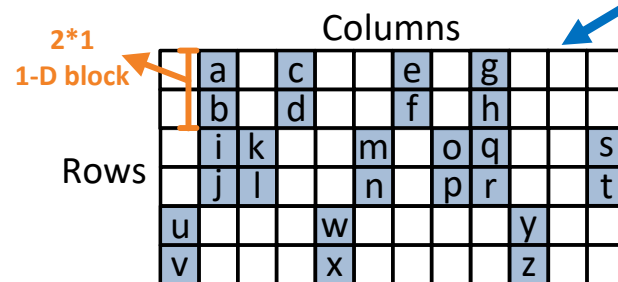
	Hopper	Ampere	Turing	Volta
Supported Tensor Core precisions	FP64, TF32, bfloat16, FP16, FP8, INT8	FP64, TF32, bfloat16, FP16, INT8, INT4, INT1	FP16, INT8, INT4, INT1	FP16

Two input matrices must be the same precision

Libraries of sparse matrix computation

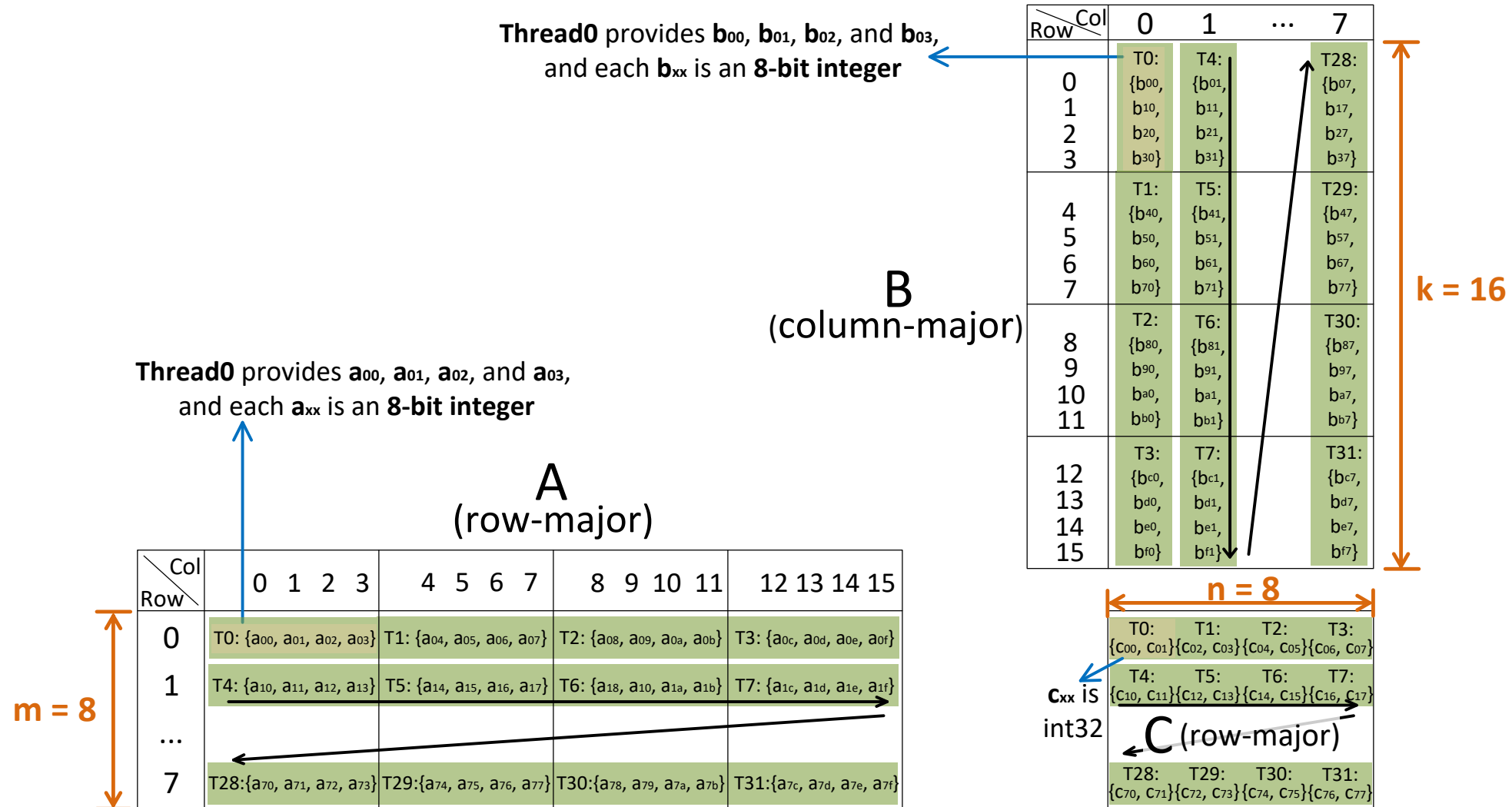
¹ Mixed precision means two input matrices with different precision

Library	Precision				Sparsity		Tensor Core
	fp16	int8	int4	mixed ¹	granularity	DL-friendly?	
cuSPARSE [10]	✓	✓	✗	✗	fine-grained	👎	👎
	✓	✓	✗	✗	block	👍	👍
cuSPARSELt [11]	✓	✓	✓	✗	2:4 structured	👍	👍
Sputnik [13]	✓	✗	✗	✗	fine-grained	👍	👎
vectorSparse [14]	✓	✗	✗	✗	1-D block	👍	👍
Magicube (ours)	✗	✓	✓	✓	<u>1-D block</u>	👍	👍

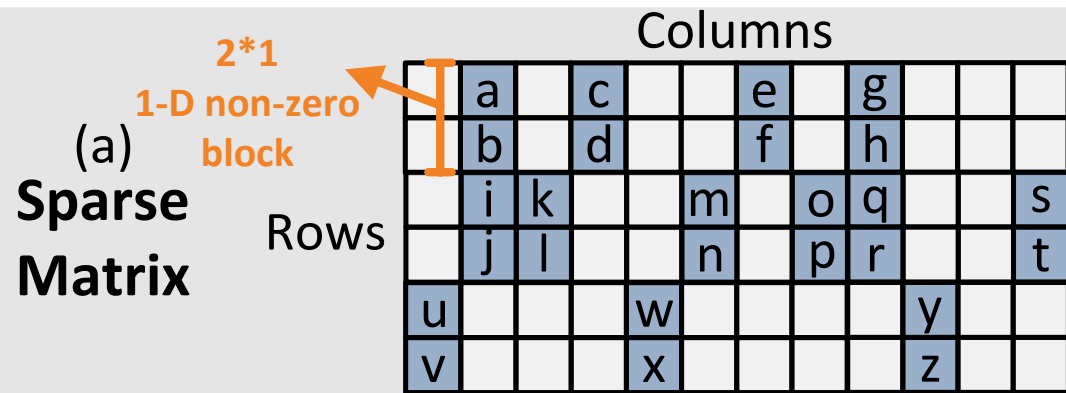


Sparse matrix with 1-D non-zero blocks

Data layout of $m8n8k16$ for $int8$ mma on Tensor Cores



SR-BCRS sparse matrix format



Sparse matrix with 1-D block non-zeros,
the length of the 1-D block = 2, 4, or 8

(b)
BCRS
format

Row pointers = [0, 4, 10, 13]

Column indices =
[1, 3, 6, 8, 1, 2, 5, 7, 8, 11, 0, 4, 9]

Values =

a	b	c	d	e	f	g	h	i	j	k	l	m
n	o	p	q	r	s	t	u	v	w	x	y	z

(c)
SR-BCRS
format
(stride=4)

Row pointers = [0, 4, 4, 10, 12, 15]

Column indices =
[1, 3, 6, 8, 1, 2, 5, 7, 8, 11, *, *, 0, 4, 9, *]

Values =

a	c	e	g	i	k	m	o	q	s	u	w	y
b	d	f	h	j	l	n	p	r	t	v	x	z

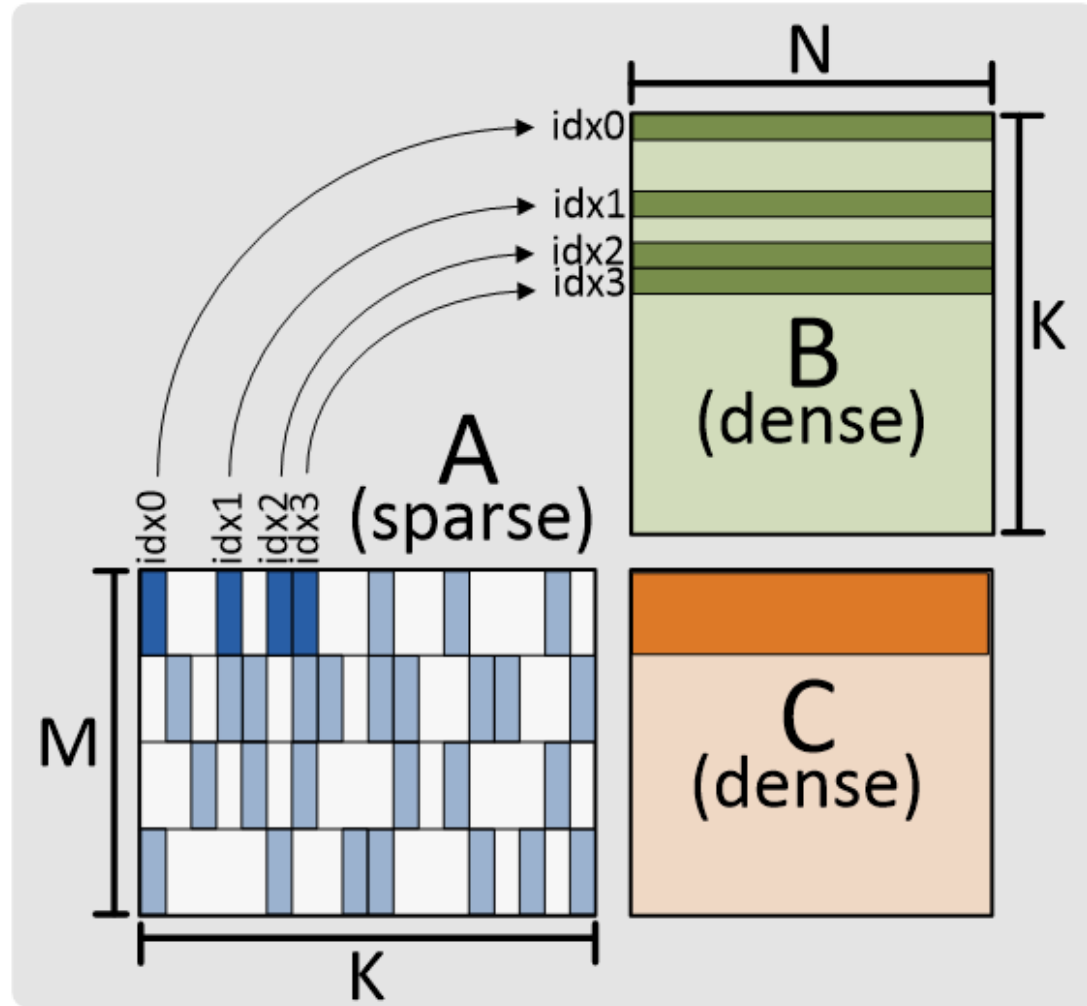
SR-BCRS (ours) is more friendly to Tensor Cores

Matrix A for *mma*

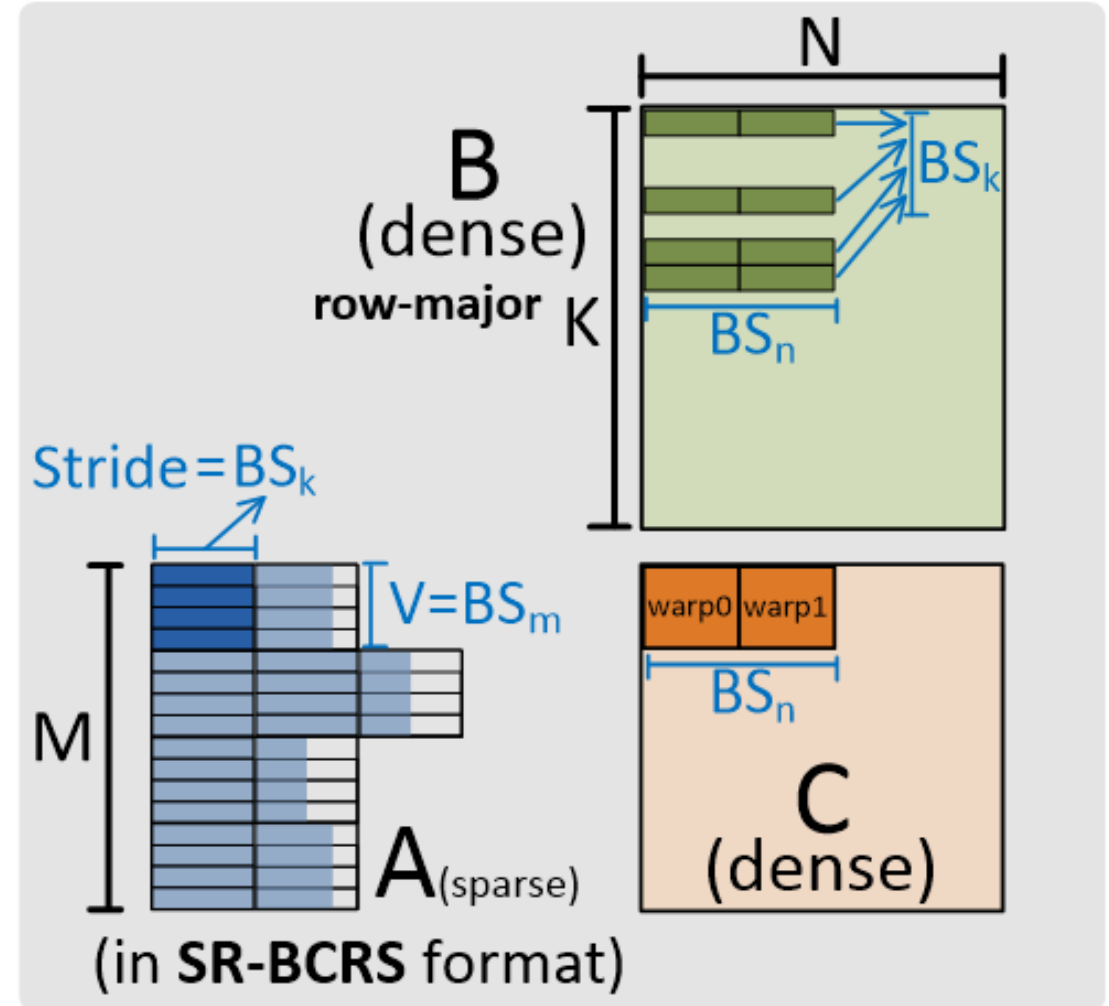
A (row-major)

Col \ Row	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	T0: {a00, a01, a02, a03}				T1: {a04, a05, a06, a07}				T2: {a08, a09, a0a, a0b}				T3: {a0c, a0d, a0e, a0f}			
1	T4: {a10, a11, a12, a13}				T5: {a14, a15, a16, a17}				T6: {a18, a19, a1a, a1b}				T7: {a1c, a1d, a1e, a1f}			
...																
7	T28: {a70, a71, a72, a73}				T29: {a74, a75, a76, a77}				T30: {a78, a79, a7a, a7b}				T31: {a7c, a7d, a7e, a7f}			

SpMM in Magicube

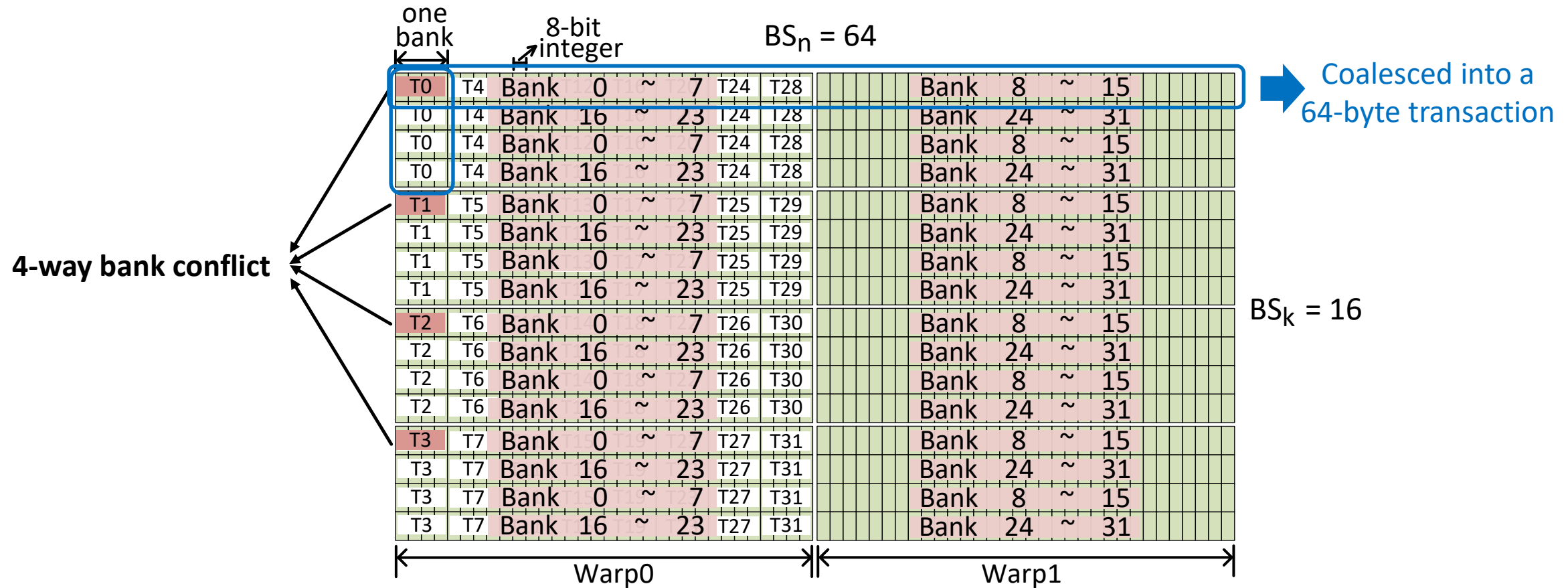


(a) SpMM

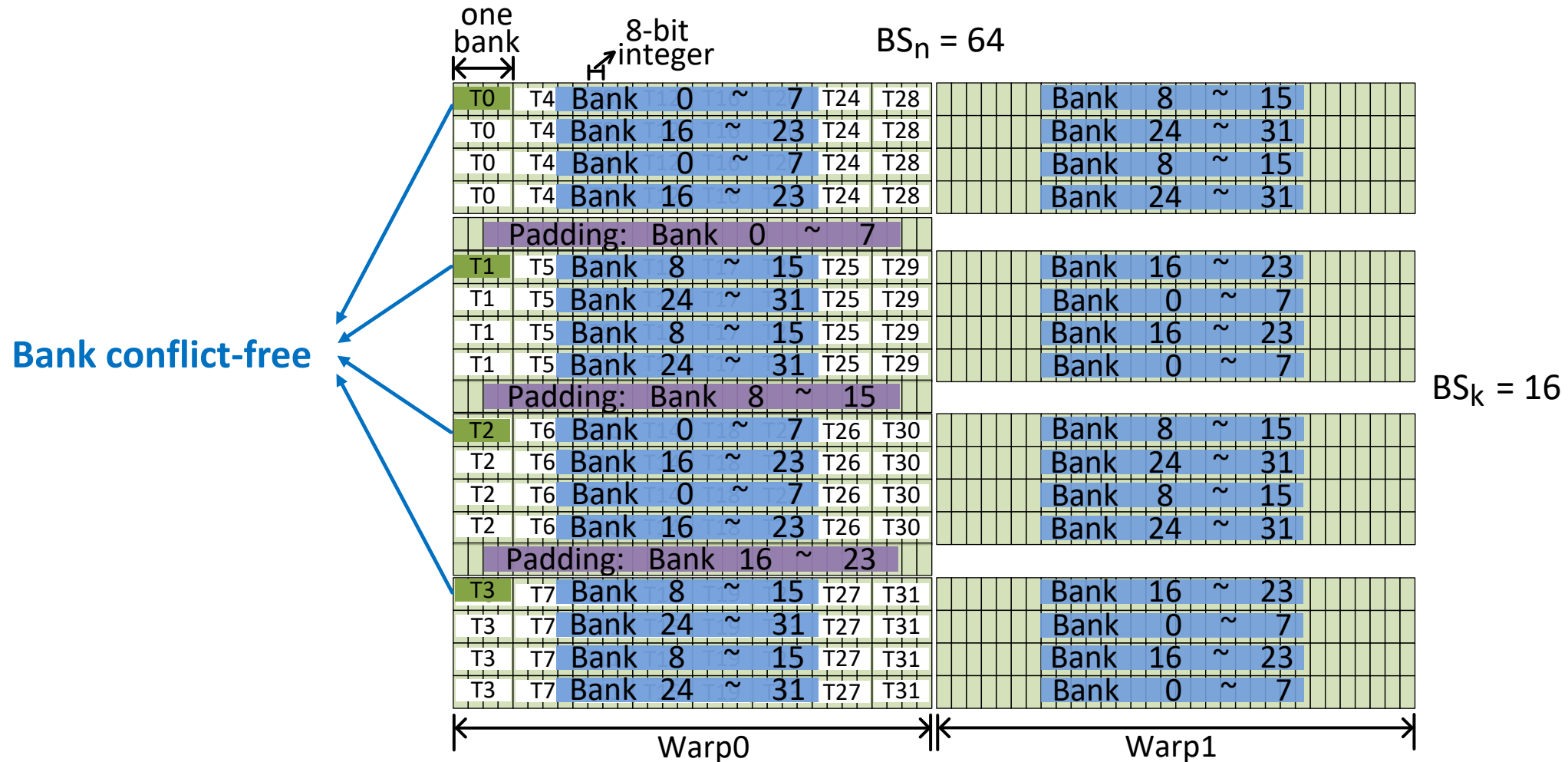


(b) SpMM in Magicube at thread-block level

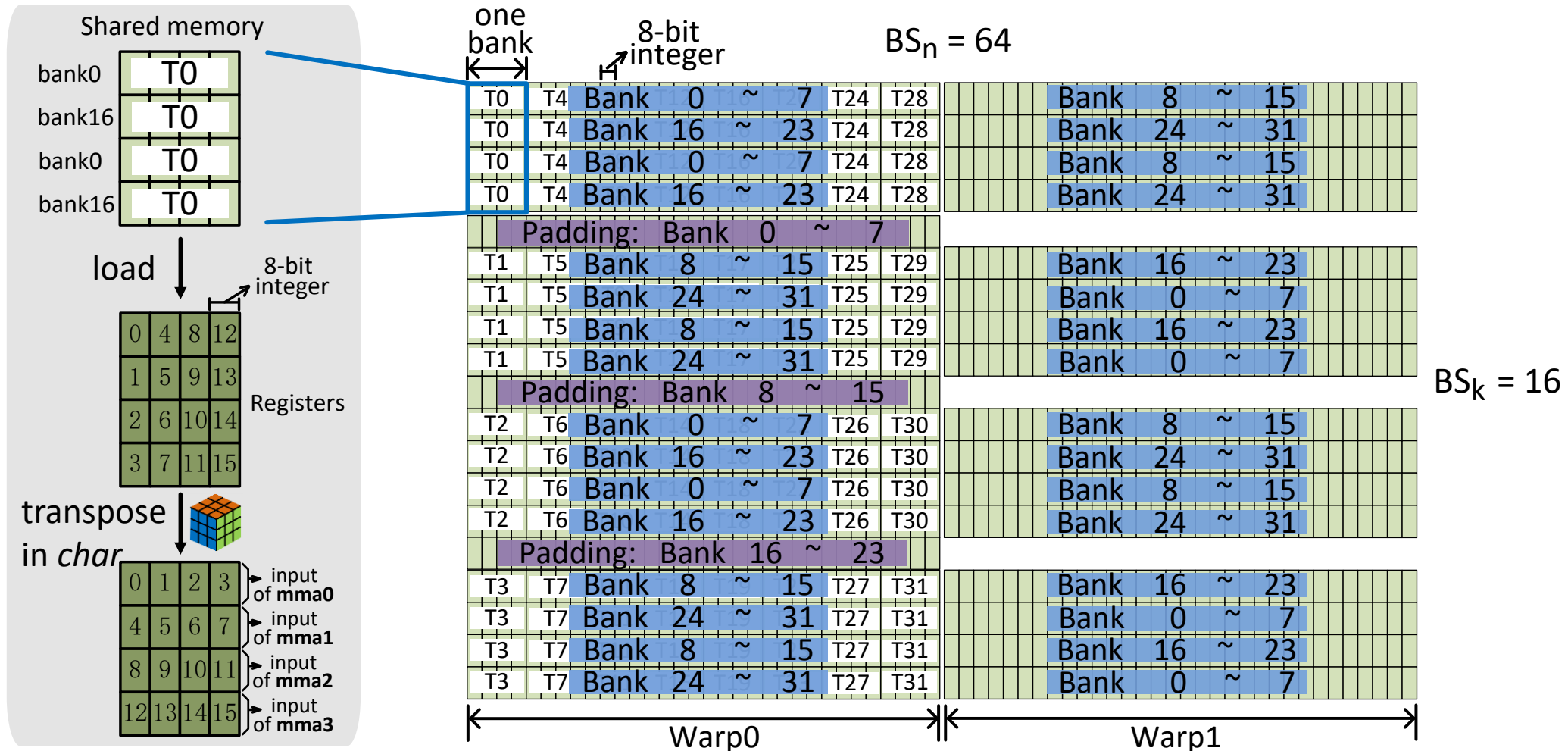
Load rows of matrix B to shared memory for *int8*



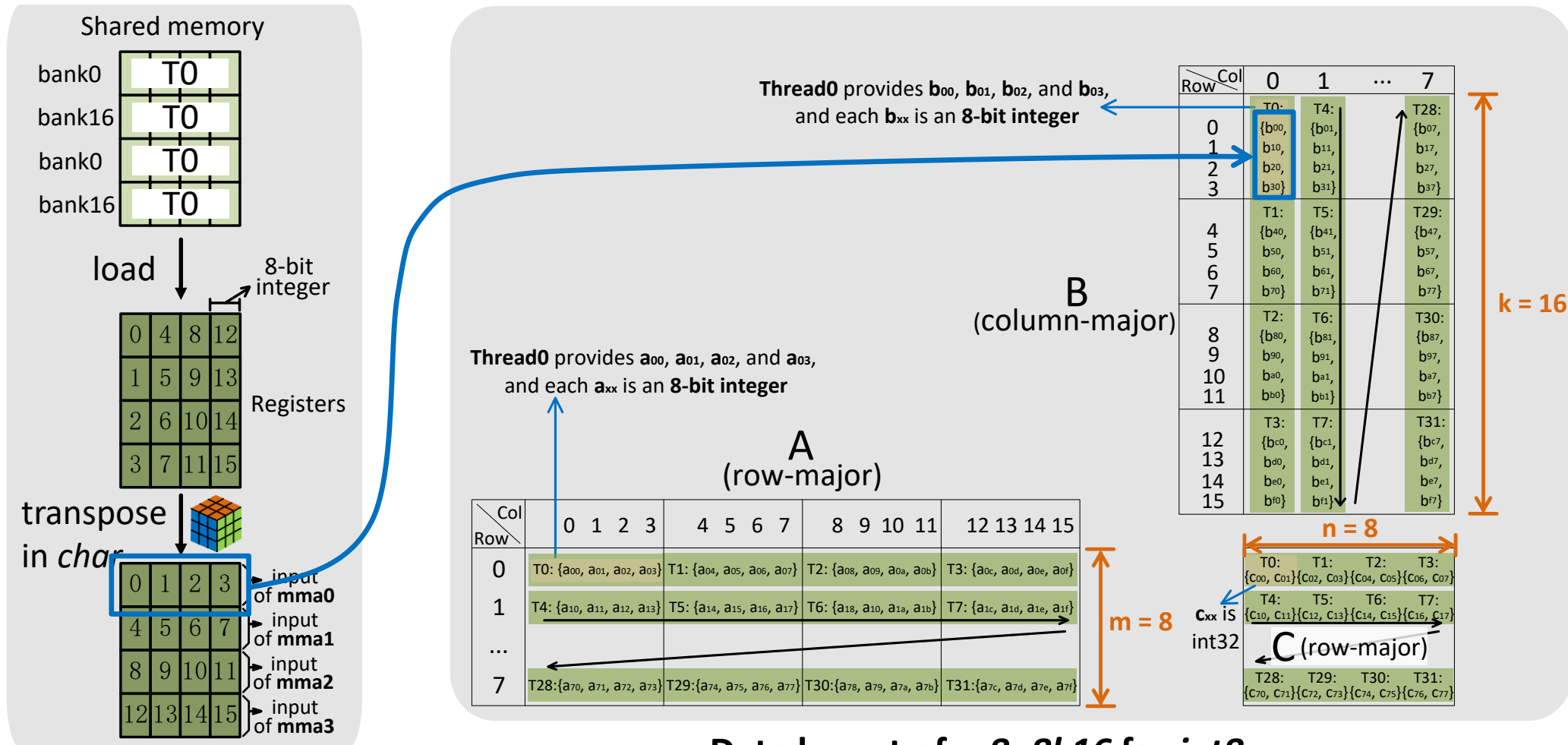
Load blocks of matrix B to shared memory for *int8*



Local transpose on registers for *int8*



Local transpose on registers for *int8*



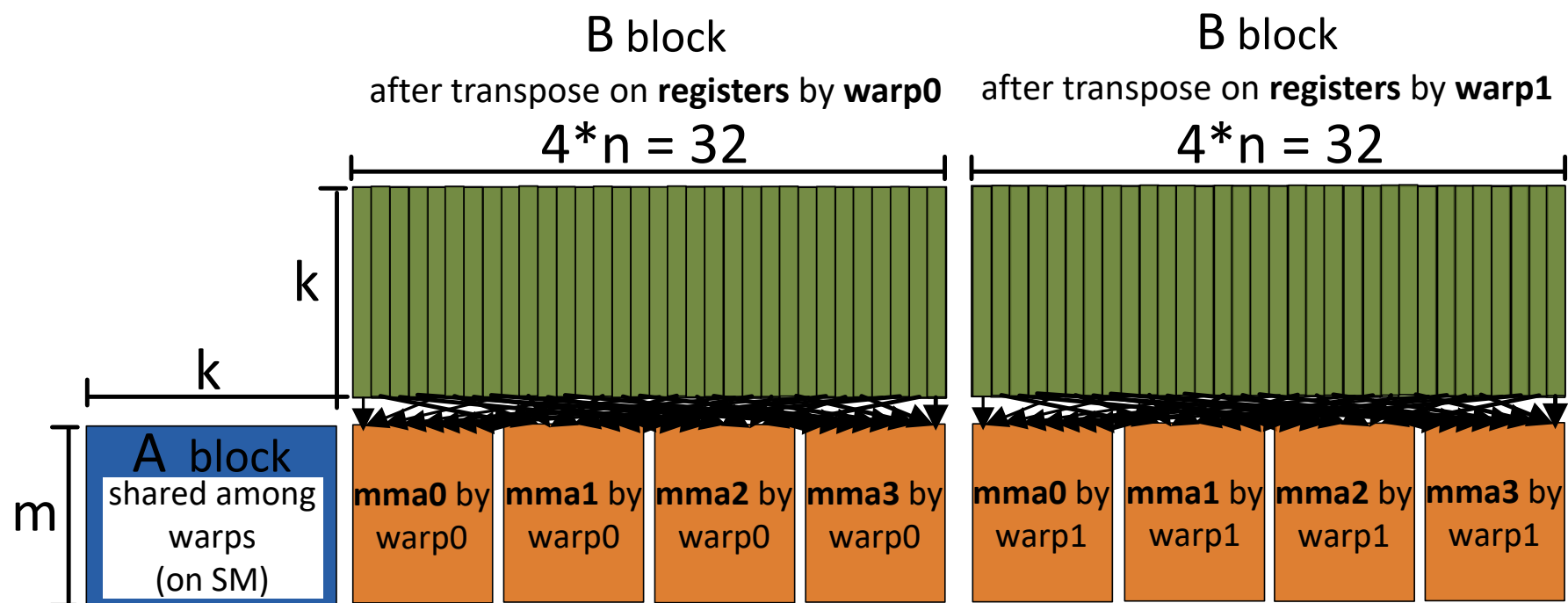
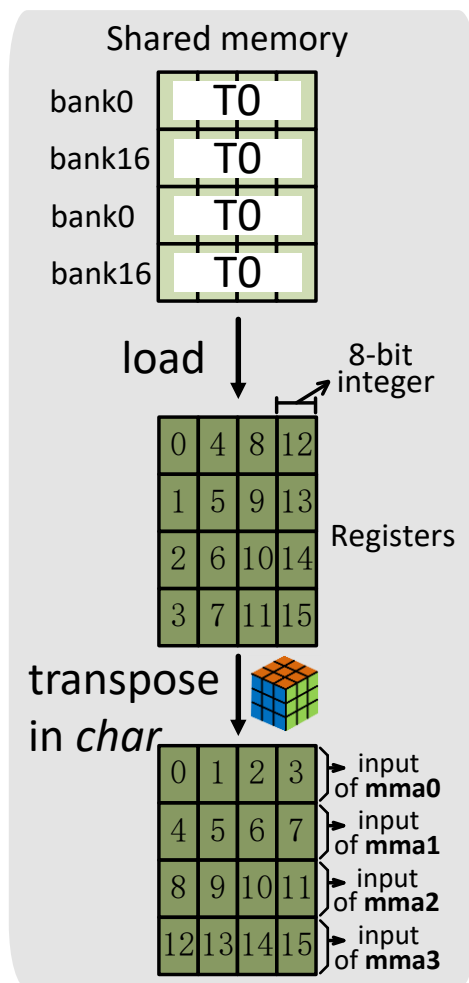
Efficiency is guaranteed by:

(1) coalesced global memory access

(2) conflict-free shared memory access

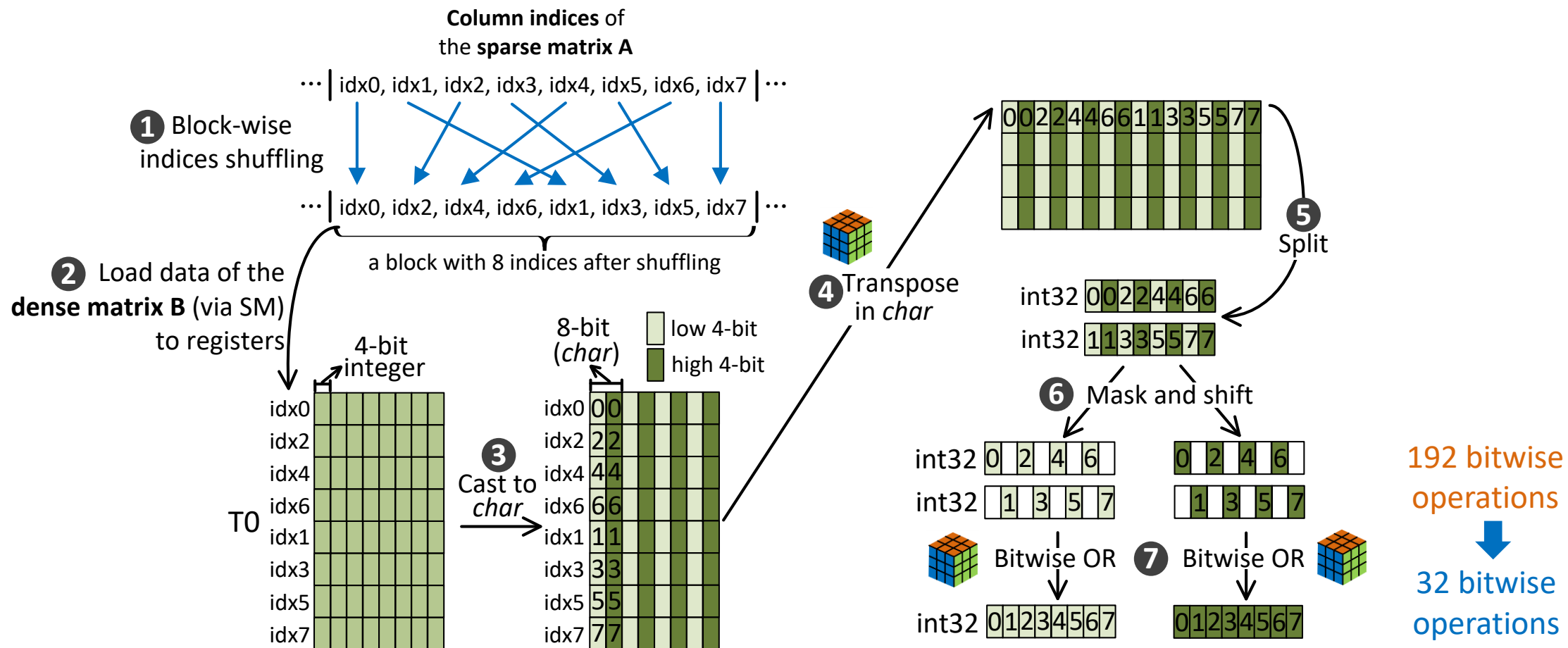
(3) fast transpose on registers

MMAAs in SpMM with *int8*

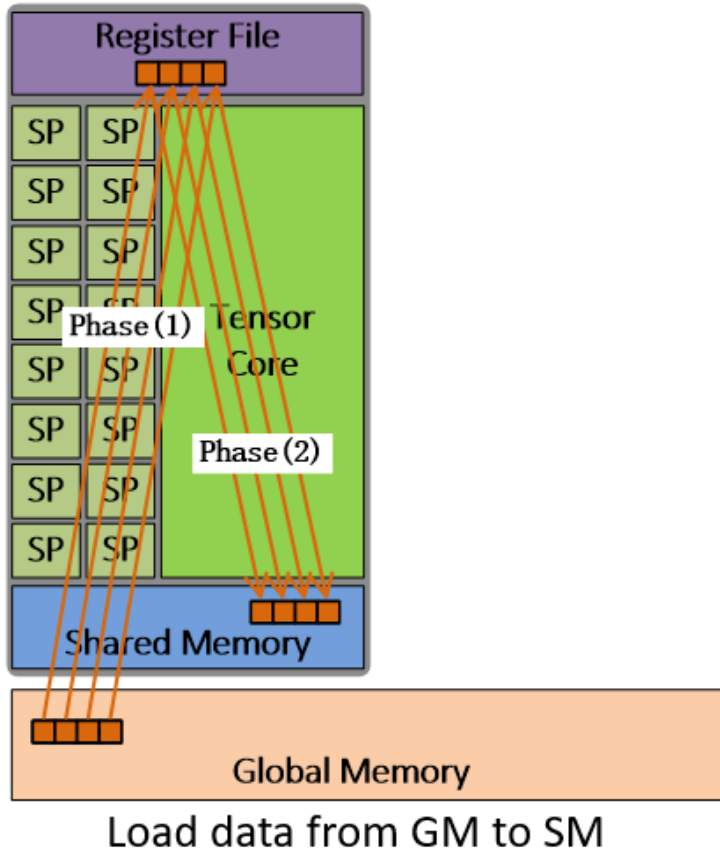


The warp-level view of MMAAs in SpMM with *int8*

Efficient local transpose for *int4* with indices shuffling



Prefetch data blocks of matrix B of SpMM



Algorithm 1 Prefetch the data block of dense matrix B

```
steps = nnz /  $BS_k$ ;
```

```
Load_A_values_and_indices_to_shared(0);  
__syncthreads();  
Prefetch_B_values_to_registers(0);
```

➔ Cold start

```
for i=1; i < steps; i++ do
```

```
Store_B_values_on_regs_to_shared(i-1);  
Load_A_values_and_indices_to_shared(i);  
__syncthreads();
```

➔ Load data and indices to SM

```
Prefetch_B_values_to_registers(i);  
MMA_compute_tiles(i-1);  
__syncthreads();
```

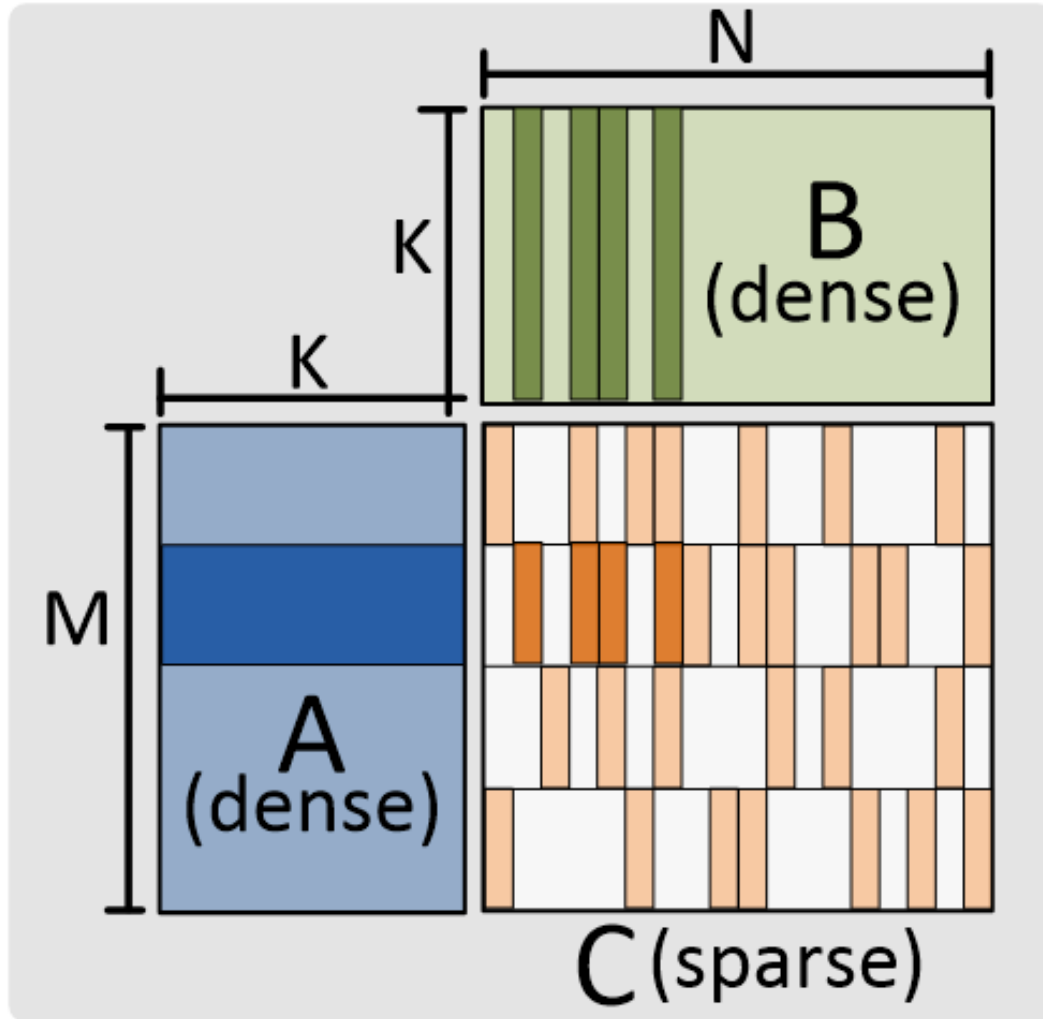
➔ Overlap prefetch with MMA

```
end for
```

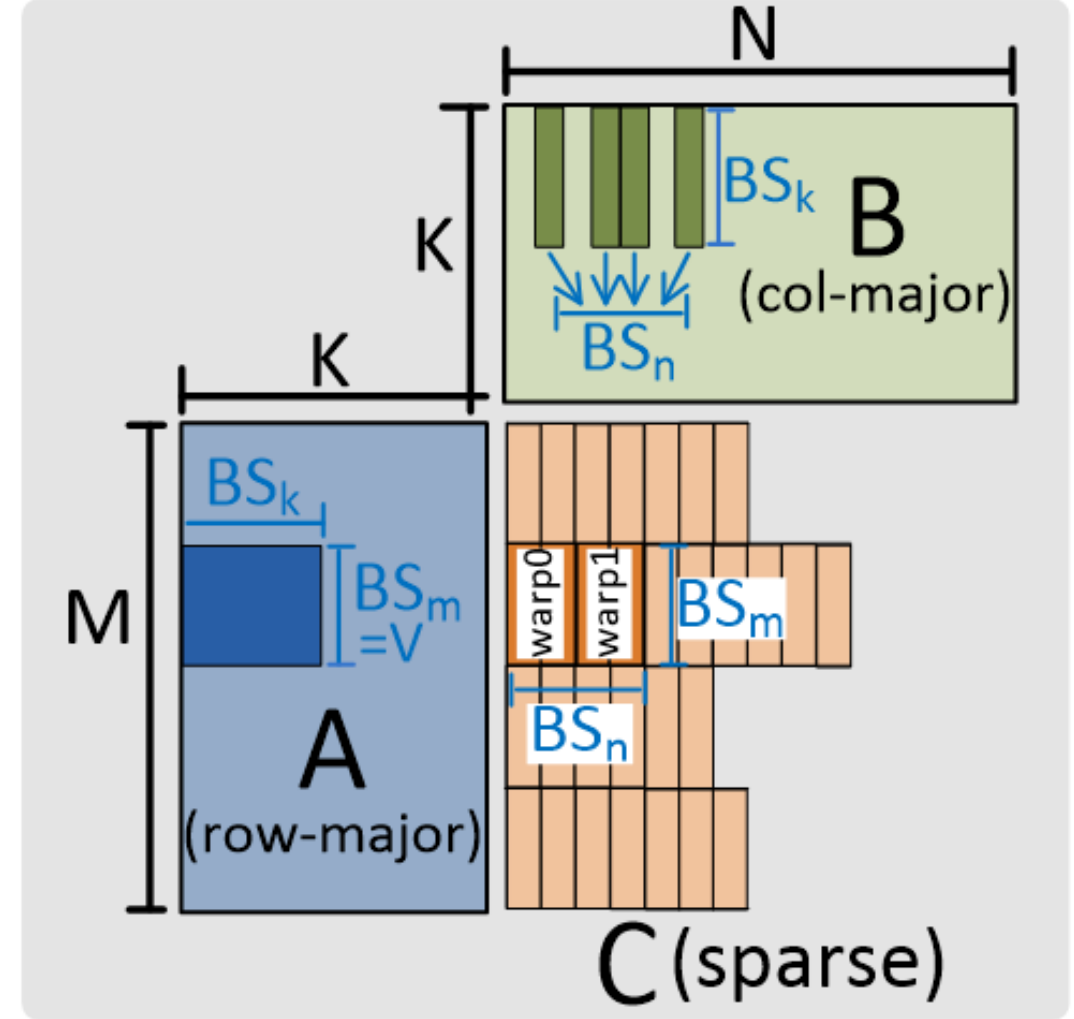
```
Store_B_values_on_regs_to_shared(i-1);  
__syncthreads();  
MMA_compute_tiles(i-1);
```

➔ The tail of pipeline

SDDMM in Magicube

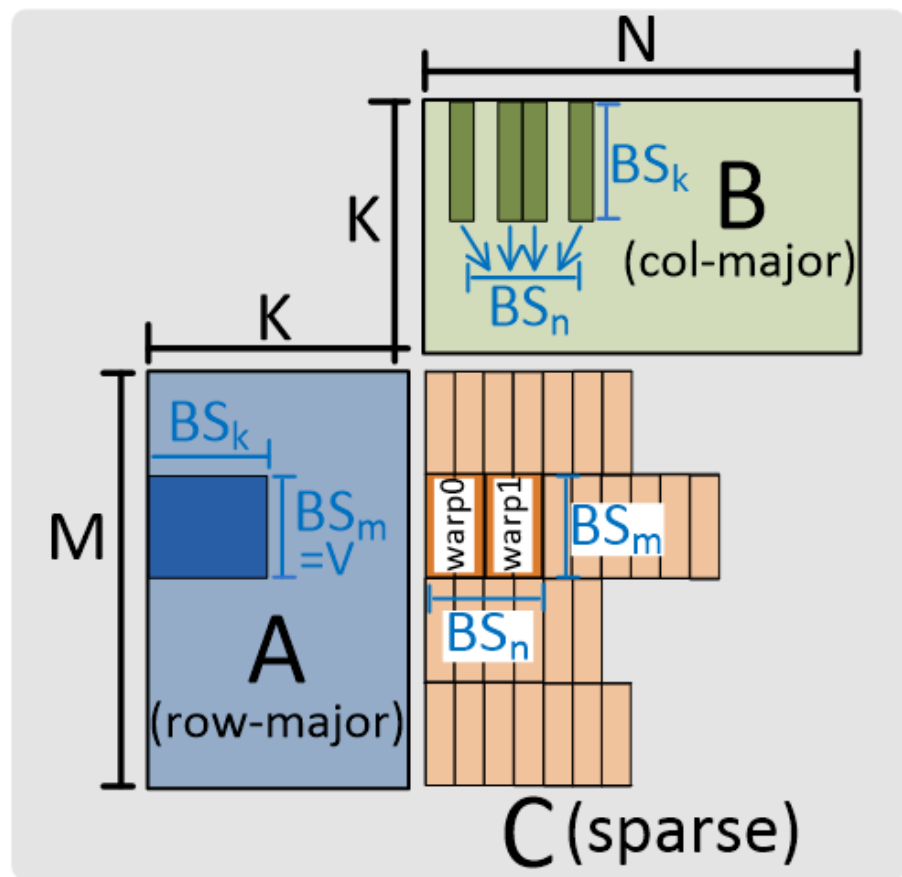


(a) SDDMM

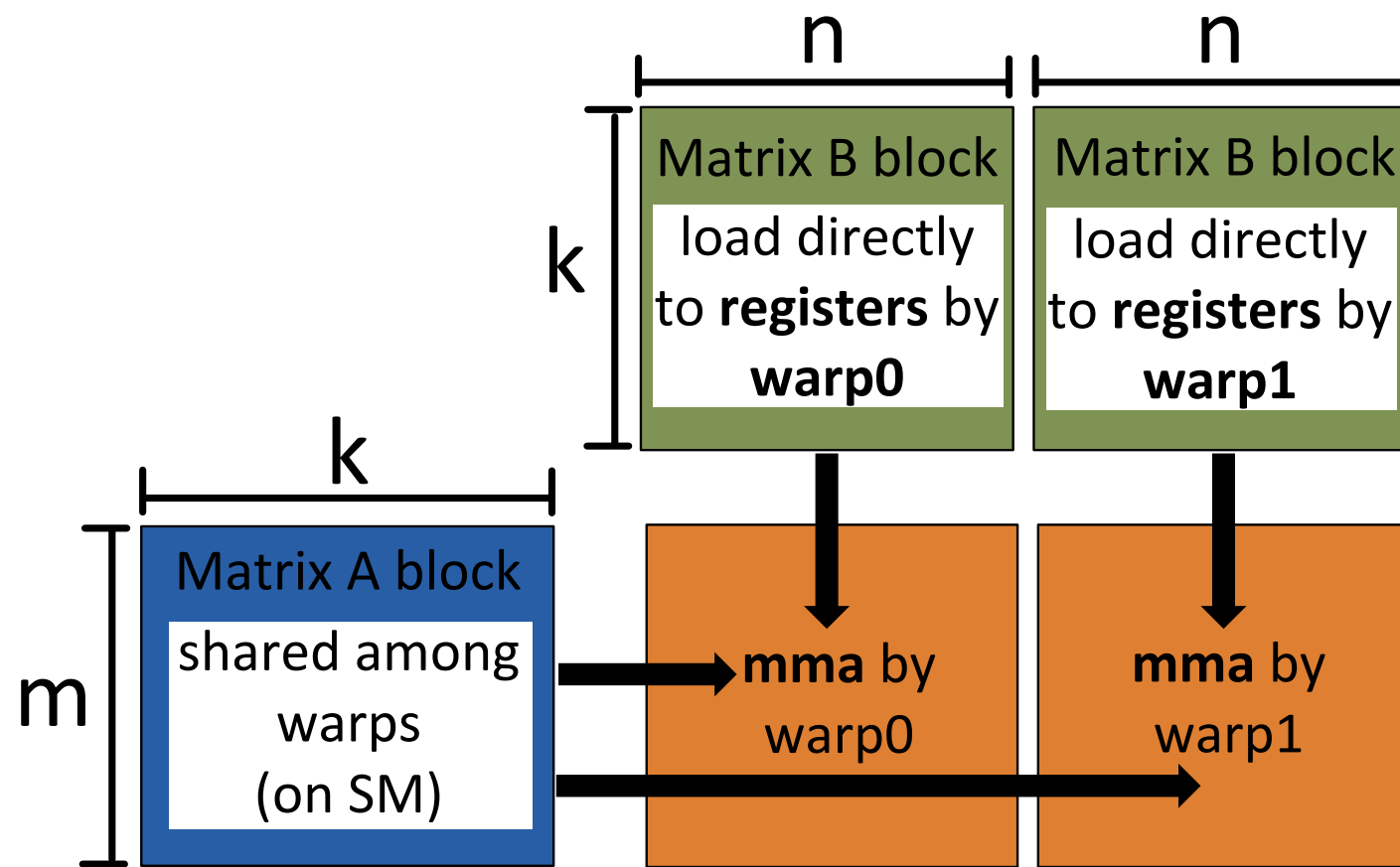


(b) SDDMM in Magicube at thread-block level

MMAAs in SDDMM



The **thread-block level view** of SDDMM



The **warp-level view** of MMAAs in SDDMM

Mixed precision

- a is an 8-bit **unsigned** integer, b is unsigned 4-bit

$a = 11101101$ (237 in decimal)

↓ Split

$a_{7\sim4}$	$a_{3\sim0}$
1110	1101
unsigned	unsigned

↓ Recover

$$a = 2^4 * a_{7\sim4} + a_{3\sim0}$$

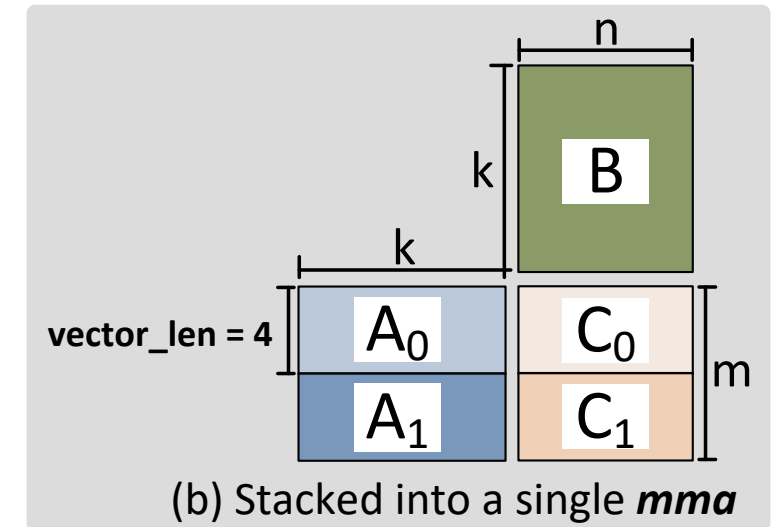
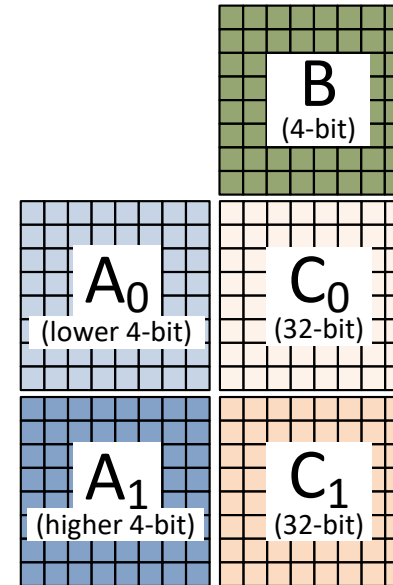
$$a * b = 2^4 * a_{7\sim4} * b + a_{3\sim0} * b$$

- a is an 8-bit **signed** integer, b is signed 4-bit

$a = 11101101$ (-19 in decimal)

↓ Split

$a_{7\sim4}$	$a_{3\sim0}$
1110	1101
signed	unsigned



$$C = 2^0 * C_0 + 2^4 * C_1$$

(a) Emulation of A (8-bit) * B (4-bit) using 4-bit mma

Evaluation

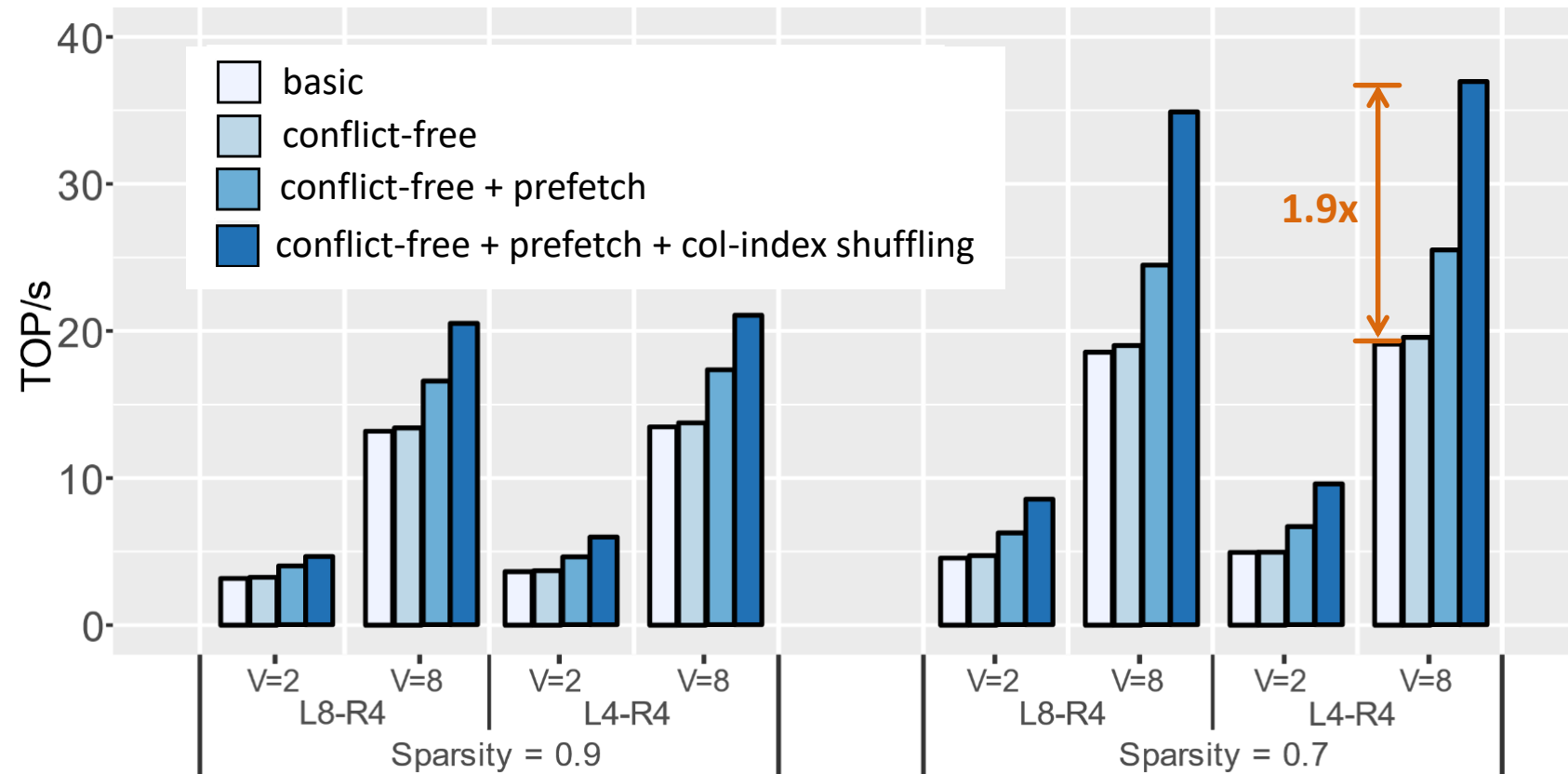
- NVIDIA A100-SXM4-40GB GPU
 - total 108 SMs
 - each SM has 192KB configurable L1 cache and shared memory, and 256KB registers
 - supported datatypes on Tensor Core: int8, int4, int1, fp16, bf16, tf32, fp64
- Compare the performance of **Magicube** with sparse libraries (**vectorSparse**, **cuSPARSE**) and dense libraries (**cuBLAS**, **cuDNN**)
- **Micro-benchmarks**: 1,536 sparse matrices from Deep Learning Matrix Collection (DLMC) with sparsity 50%~98%, dilating each scalar with 1-D blocks (length $V = 2, 4, 8$)
- **Case study**: end-to-end sparse Transformer inference



One streaming multiprocessor (SM) of GA100

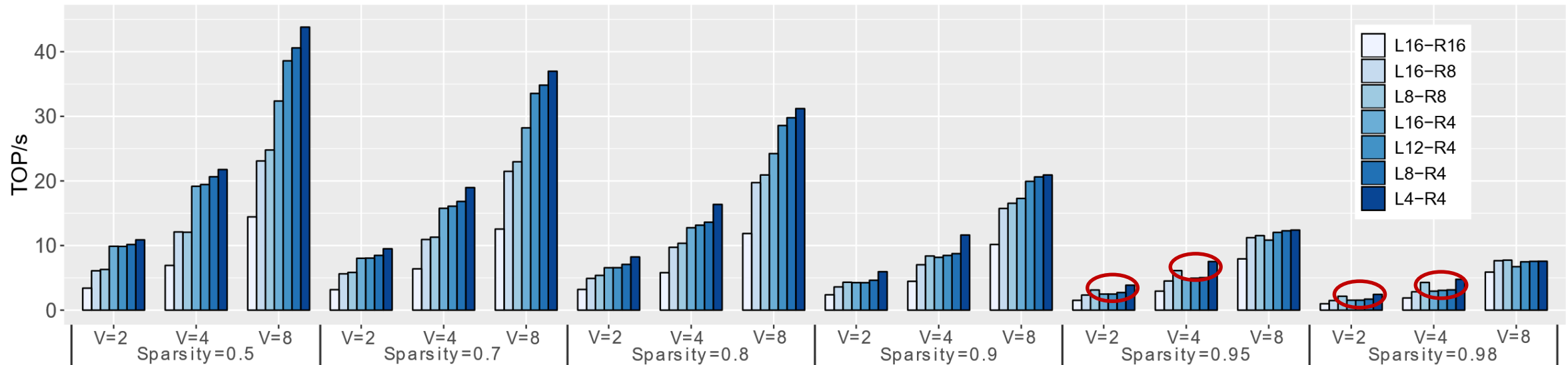
This image is from: R. Krashinsky, et al. <https://developer.nvidia.com/blog/nvidia-ampere-architecture-in-depth/> May, 2020

Ablation study for SpMM in Magicube



Ablation study for optimizations of SpMM

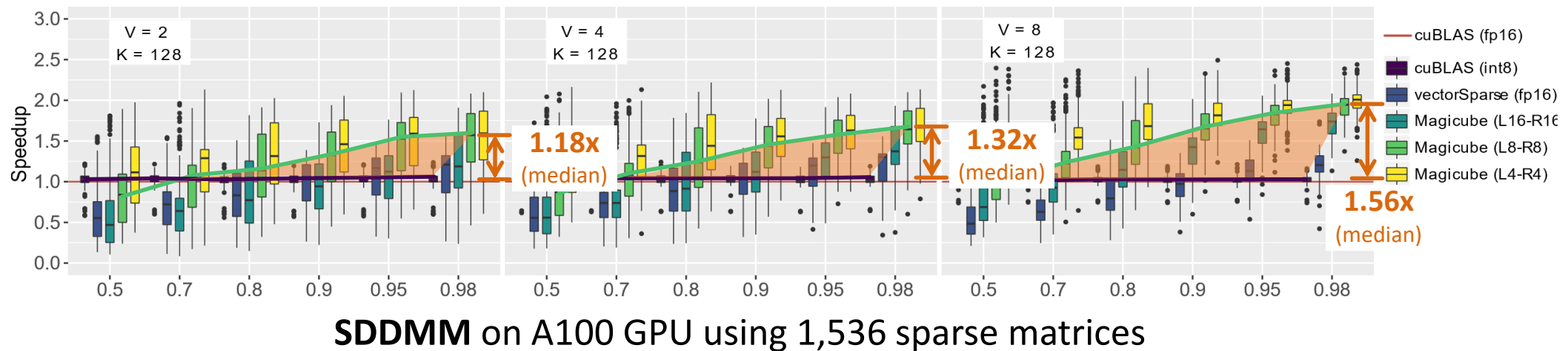
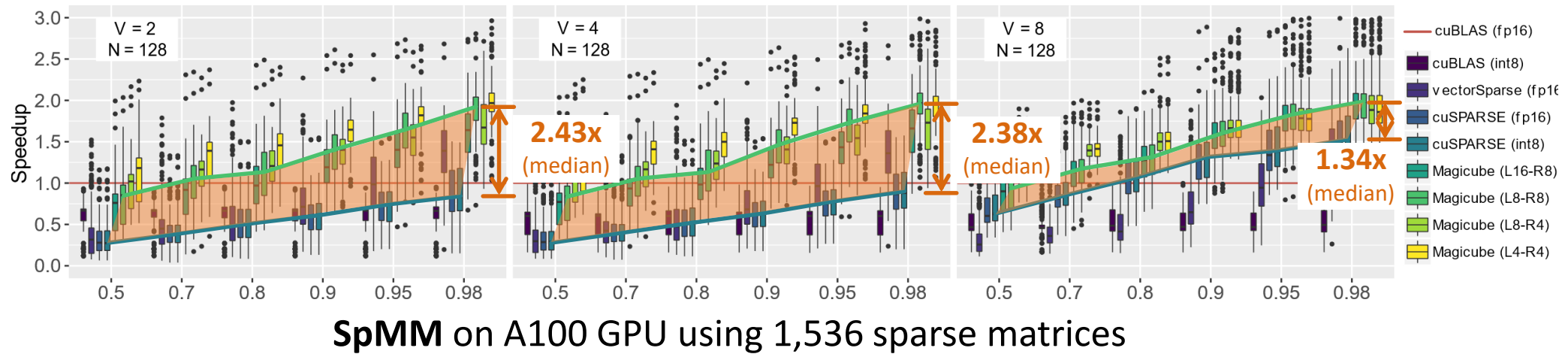
SpMM with mixed precision in Magicube



SpMM with mixed precision

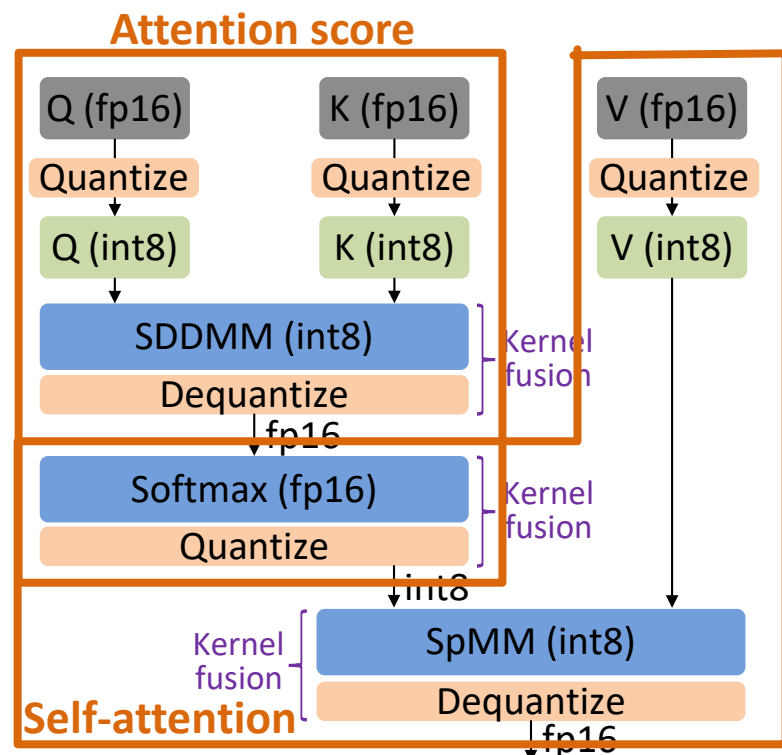
Lx-Ry means x-bit A matrix multiplied by y-bit B matrix

Benchmarking SpMM and SDDMM

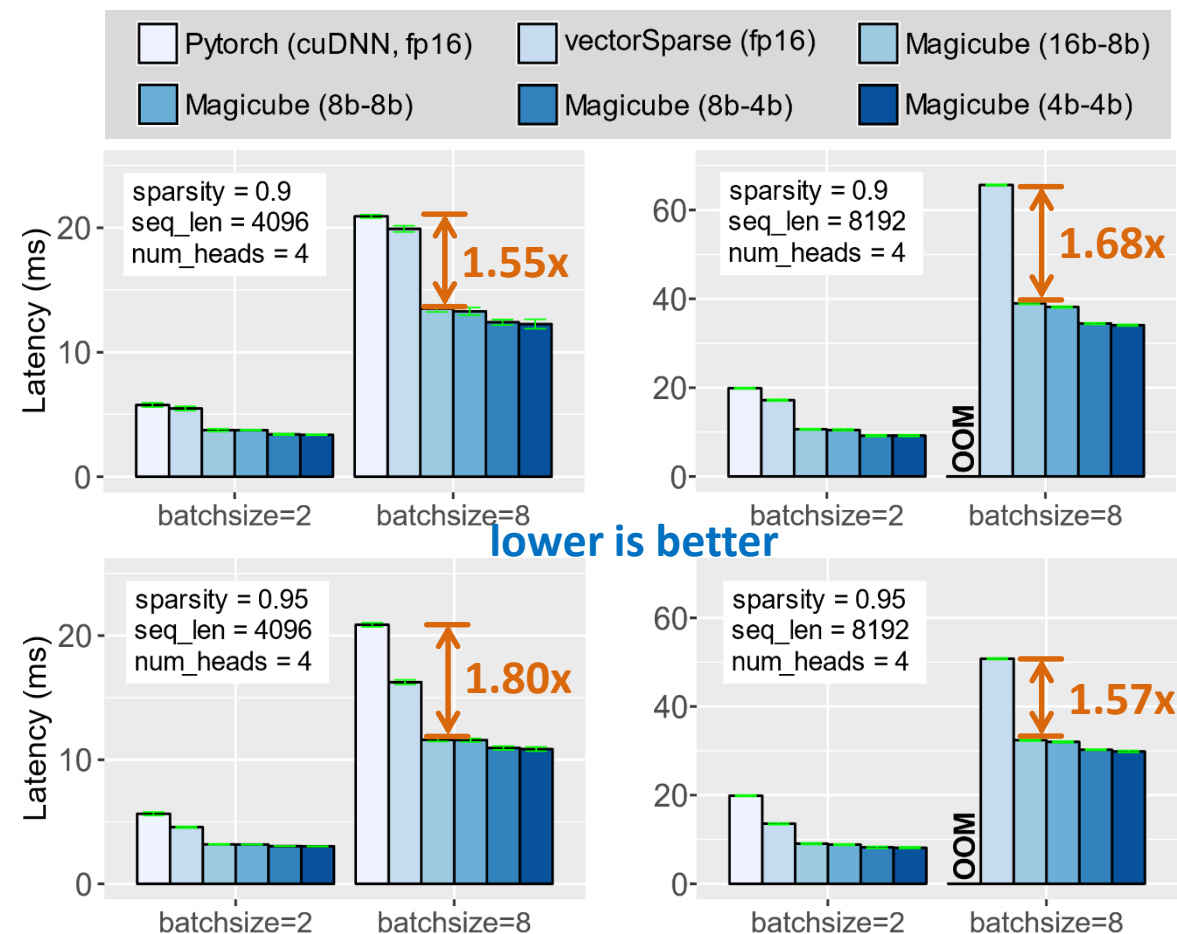


End-to-end sparse Transformer inference

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T \odot M}{\sqrt{d_k}}\right) V$$



Quantized self-attention with sparse attention mask



Latency of end-to-end inference of sparse Transformer

End-to-end sparse Transformer inference

dense		sparsity=0.9			
PyTorch (cuDNN, fp32)	PyTorch (cuDNN, fp16)	vectorSparse (fp16)	Magicube (16b-8b)	Magicube (8b-8b)	Magicube (8b-4b)
57.36%	57.50%	57.14%	57.32% ←	57.11%	56.79%

Test accuracy of text classification using sparse Transformer
 with num_heads=4 and seq_len=4,096

Conclusion

1. Challenges

Challenges

(1) How to achieve practical speedup in a large range of sparsity ratio, e.g., 50% ~ 98%?

How to bridge the big gap?

(2) How to efficiently support sparse workloads with mixed precision (two input matrices with different precision), e.g., 8-bit weights and 4-bit activation?

	Hopper	Ampere	Turing	Volta
Supported Tensor Core precisions	FP64, TF32, bfloat16, FP16, FP8, INT8	FP64, TF32, bfloat16, FP16, INT8, INT4, INT1	FP16, INT8, INT4, INT1	FP16

Two input matrices must be the same precision

2. SR-BCRS format

SR-BCRS sparse matrix format

(a) Sparse Matrix

(b) BCRS format

Row pointers = [0, 4, 7, 13]
Column indices = [1, 3, 6, 8, 0, 4, 9, 1, 2, 5, 7, 8, 11]
Values = [a b c d e f g h i j k l m n o p q r s t u v w x y z]

(c) SR-BCRS format (stride=4)

Row pointers = [0, 4, 4, 7, 8, 14]
Column indices = [1, 3, 6, 8, 0, 4, 9, x, 1, 2, 5, 7, 8, 11, x, x]
Values = [a c e g i k m o q s u w y i j l n p r t v x z]

SR-BCRS is more friendly to Tensor Cores

3. SpMM in Magicube

SpMM in Magicube

(a) SpMM

(b) SpMM in Magicube at thread-block level

4. SDDMM in Magicube

SDDMM in Magicube

(a) SDDMM

(b) SDDMM in Magicube at thread-block level

5. Mixed precision

Mixed precision

- a is an 8-bit unsigned integer, b is unsigned 4-bit

a = 11101101 (237 in decimal)

Split: a₇₋₄ = 1110, a₃₋₀ = 1101

Recover: a = 2⁴ * a₇₋₄ + a₃₋₀

a * b = 2⁴ * a₇₋₄ * b + a₃₋₀ * b

- a is an 8-bit signed integer, b is signed 4-bit

a = 11101101 (-19 in decimal)

Split: a₇₋₄ = 1110, a₃₋₀ = 1101

Recover: a = 2⁴ * a₇₋₄ + a₃₋₀

a * b = 2⁴ * a₇₋₄ * b + a₃₋₀ * b

(a) Emulation of A (8-bit) * B (4-bit) using 4-bit mma

(b) Stacked into a single mma

6. Evaluation

End-to-end sparse Transformer inference

Attention score: $\text{softmax}\left(\frac{QK^T \odot M}{\sqrt{d_k}}\right) V$

Quantized self-attention with sparse attention mask

Latency of end-to-end inference of sparse Transformer