

Narrow bit-width number formats for deep learning

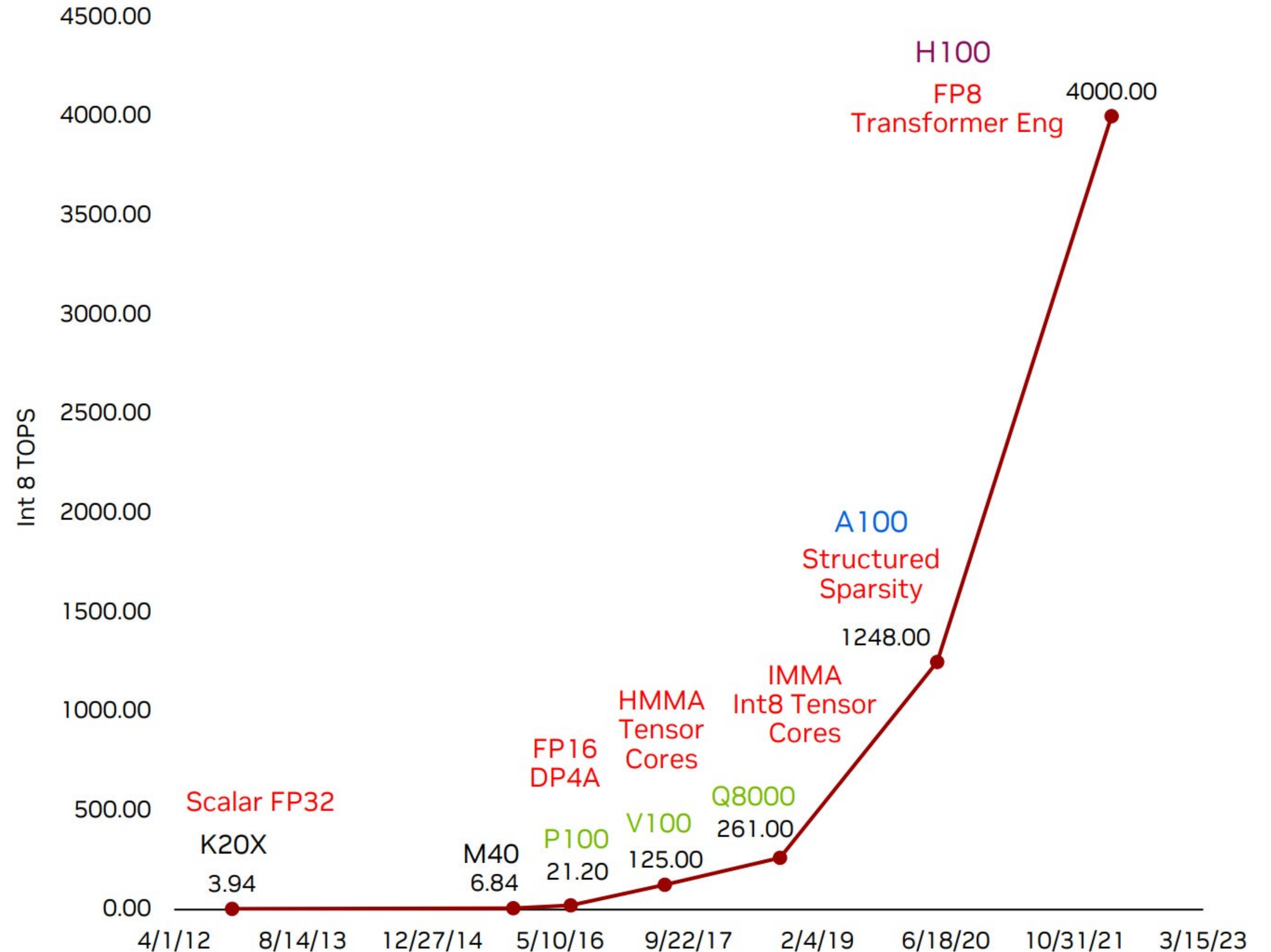
The context: Quantization in Deep Learning

- The basic "atom" of computation in a ML model is matrix multiplication.
- We have a choice on the number of bits used to represent the "weights" in our model (i.e., the learnable parameters in the matrices that we are multiplying the input by).
- Lowering the number of bits used lowers the possible precision of our calculations, and hence the accuracy of our model.
- On the other hand, it speeds up the training and inference of the model while improve its memory and power efficiency.
- In practice, the accuracy loss can be managed and this is an acceptable tradeoff.

Gains from

- Number Representation
 - FP32, FP16, Int8
 - (TF32, BF16)
 - ~16x
- Complex Instructions
 - DP4, HMMA, IMMA
 - ~12.5x
- Process
 - 28nm, 16nm, 7nm, 5nm
 - ~2.5x
- Sparsity
 - ~2x
- Model efficiency has also improved – overall gain > 1000x

Single-Chip Inference Performance - 1000X in 10 years



Cost of Operations

Relative Energy Cost

Operation:	Energy (pJ)
8b Add	0.03
16b Add	0.05
32b Add	0.1
16b FP Add	0.4
32b FP Add	0.9
8b Mult	0.2
32b Mult	3.1
16b FP Mult	1.1
32b FP Mult	3.7
32b SRAM Read (8KB)	5
32b DRAM Read	640

1 10 100 1000 10000

Relative Area Cost

Area (μm^2)
36
67
137
1360
4184
282
3495
1640
7700
N/A
N/A

1 10 100 1000

Energy numbers are from Mark Horowitz "Computing's Energy Problem (and what we can do about it)", ISSCC 2014
 Area numbers are from synthesized result using Design Compiler under TSMC 45nm tech node. FP units used DesignWare Library.

Basics of number formats

- Numbers are stored in binary representation on a computer.
- For a floating point number, you have a sign bit, and then bits allocated to both the *exponent* and *mantissa* (significand).
- For example, the IEEE 754 standard specifies the *single-precision floating point* format FP32 as having 1 sign bit, 8 exponent bits, and 24 mantissa bits.

$$(-1)^{b_{31}} \times 2^{(b_{30}b_{29}\dots b_{23})_2 - 127} \times (1.b_{22}b_{21}\dots b_0)_2$$

- When reducing the number of bits used (always using a power of 2), we have a choice on how to allocate bits to the exponent and mantissa.
- Different choices may work better for training vs. inference (range vs. precision).

Basic Number Formats

Format	Sign	Exponent	Mantissa	Total bits
FP64	1	11	52	64
FP32	1	8	23	32
TF32	1	8	10	19
FP16	1	5	10	16
BF16	1	8	7	16
UINT8	0	0	8	8
INT8	1	0	7	8
FP8 (e5)	1	5	2	8
FP8 (e4)	1	4	3	8
FP8 (e3)	1	3	4	8
UINT4	0	0	4	4
INT4	1	0	3	4

Hopper Tensor Cores by Precision Type

- The H100 GPU has native support for number formats down to 8 bits.

	H100 SXM	H100 PCIe	H100 NVL ¹
FP64	34 teraFLOPS	26 teraFLOPS	68 teraFLOPS
FP64 Tensor Core	67 teraFLOPS	51 teraFLOPS	134 teraFLOPS
FP32	67 teraFLOPS	51 teraFLOPS	134 teraFLOPS
TF32 Tensor Core	989 teraFLOPS ²	756 teraFLOPS ²	1,979 teraFLOPS ²
BFLOAT16 Tensor Core	1,979 teraFLOPS ²	1,513 teraFLOPS ²	3,958 teraFLOPS ²
FP16 Tensor Core	1,979 teraFLOPS ²	1,513 teraFLOPS ²	3,958 teraFLOPS ²
FP8 Tensor Core	3,958 teraFLOPS ²	3,026 teraFLOPS ²	7,916 teraFLOPS ²
INT8 Tensor Core	3,958 TOPS ²	3,026 TOPS ²	7,916 TOPS ²

Targeting the Hopper Tensor Cores in CUDA

- On Hopper, we have WGMMMA (Warpgroup Matrix Multiply Accumulate) for doing accelerated matrix multiplications $C = A * B$.
- WGMMMA corresponds to the instruction *wgmma.mma_async* in the PTX ISA.
- In *wgmma.mma_async*, one warpgroup (=128 threads) cooperatively performs a matrix-multiplication on tiles of the matrices A and B (of certain allowed dimensions, like $\{M,N,K\} = \{64,64,16\}$).
- *wgmma.mma_async* makes assumptions on the layouts of the fragments of the operand matrices held by the different threads.

R\C	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	T0:{a0, a1}		T1:{a0, a1}		T2:{a0, a1}		T3:{a0, a1}		T0:{a4, a5}		T1:{a4, a5}		T2:{a4, a5}		T3:{a4, a5}	
1	T4:{a0, a1}		T5:{a0, a1}		T6:{a0, a1}		T7:{a0, a1}		T4:{a4, a5}		T5:{a4, a5}		T6:{a4, a5}		T7:{a4, a5}	
...	→							←								
7	T28:{a0, a1}		T29:{a0, a1}		T30:{a0, a1}		T31:{a0, a1}		T28:{a4, a5}		T29:{a4, a5}		T30:{a4, a5}		T31:{a4, a5}	
8	T0:{a2, a3}		T1:{a2, a3}		T2:{a2, a3}		T3:{a2, a3}		T0:{a6, a7}		T1:{a6, a7}		T2:{a6, a7}		T3:{a6, a7}	
.	→							←								
15	T0:{a2, a3}		T1:{a2, a3}		T2:{a2, a3}		T3:{a2, a3}		T28:{a6, a7}		T29:{a6, a7}		T30:{a6, a7}		T31:{a6, a7}	
16	T32:{a0, a1}		T33:{a0, a1}		T34:{a0, a1}		T35:{a0, a1}		T32:{a4, a5}		T33:{a4, a5}		T34:{a4, a5}		T35:{a4, a5}	
..	→							←								
31	T60:{a2, a3}		T61:{a2, a3}		T62:{a2, a3}		T63:{a2, a3}		T60:{a6, a7}		T61:{a6, a7}		T62:{a6, a7}		T63:{a6, a7}	
32	T64:{a0, a1}		T65:{a0, a1}		T66:{a0, a1}		T67:{a0, a1}		T64:{a4, a5}		T65:{a4, a5}		T66:{a4, a5}		T67:{a4, a5}	
..	→							←								
47	T92:{a2, a3}		T93:{a2, a3}		T94:{a2, a3}		T95:{a2, a3}		T92:{a6, a7}		T93:{a6, a7}		T94:{a6, a7}		T95:{a6, a7}	
48	T96:{a0, a1}		T97:{a0, a1}		T98:{a0, a1}		T99:{a0, a1}		T96:{a4, a5}		T97:{a4, a5}		T98:{a4, a5}		T99:{a4, a5}	
..	→							←								
63	T124:{a2, a3}		T125:{a2, a3}		T126:{a2, a3}		T127:{a2, a3}		T124:{a6, a7}		T125:{a6, a7}		T126:{a6, a7}		T127:{a6, a7}	

(tid % 128) : fragments

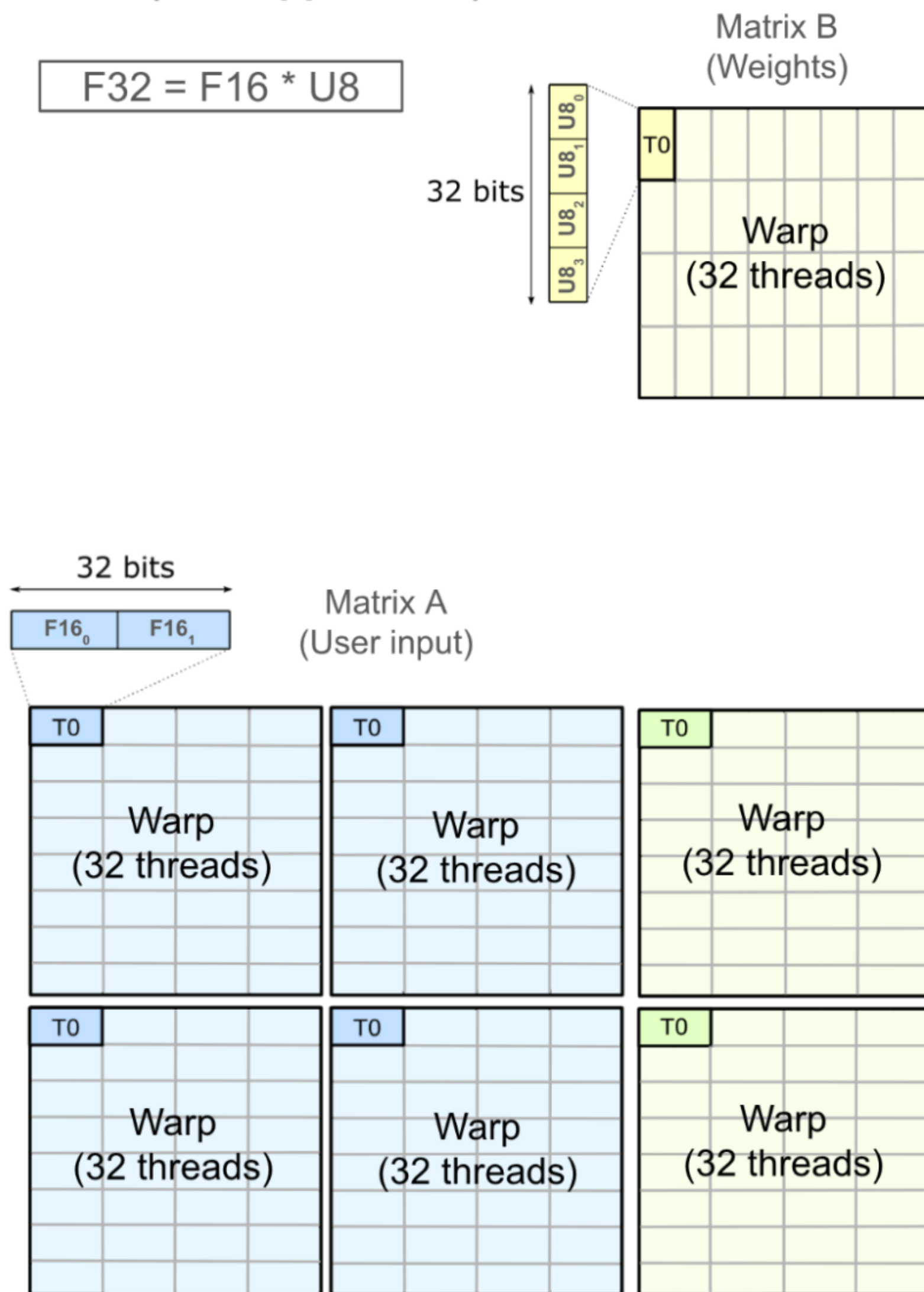
R\C	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31							
0	T0:{a0, a1, a2, a3}	T1:{a0, a1, a2, a3}	T2:{a0, a1, a2, a3}	T3:{a0, a1, a2, a3}	T0:{a8, a9, a10, a11}	T1:{a8, a9, a10, a11}	T2:{a8, a9, a10, a11}	T3:{a8, a9, a10, a11}																															
1	T4:{a0, a1, a2, a3}	T5:{a0, a1, a2, a3}	T6:{a0, a1, a2, a3}	T7:{a0, a1, a2, a3}	T4:{a8, a9, a10, a11}	T5:{a8, a9, a10, a11}	T6:{a8, a9, a10, a11}	T7:{a8, a9, a10, a11}																															
...																																							
7	T28:{a0, a1, a2, a3}	T29:{a0, a1, a2, a3}	T30:{a0, a1, a2, a3}	T31:{a0, a1, a2, a3}	T28:{a8, a9, a10, a11}	T29:{a8, a9, a10, a11}	T30:{a8, a9, a10, a11}	T31:{a8, a9, a10, a11}																															
8	T0:{a4, a5, a6, a7}	T1:{a4, a5, a6, a7}	T2:{a4, a5, a6, a7}	T3:{a4, a5, a6, a7}	T0:{a12, a13, a14, a15}	T1:{a12, a13, a14, a15}	T2:{a12, a13, a14, a15}	T3:{a12, a13, a14, a15}																															
.																																							
15	T28:{a4, a5, a6, a7}	T29:{a4, a5, a6, a7}	T30:{a4, a5, a6, a7}	T31:{a4, a5, a6, a7}	T28:{a12, a13, a14, a15}	T29:{a12, a13, a14, a15}	T30:{a12, a13, a14, a15}	T31:{a12, a13, a14, a15}																															
16	T32:{a0, a1, a2, a3}	T33:{a0, a1, a2, a3}	T34:{a0, a1, a2, a3}	T35:{a0, a1, a2, a3}	T32:{a8, a9, a10, a11}	T33:{a8, a9, a10, a11}	T34:{a8, a9, a10, a11}	T35:{a8, a9, a10, a11}																															
..																																							
31	T60:{a4, a5, a6, a7}	T61:{a4, a5, a6, a7}	T62:{a4, a5, a6, a7}	T63:{a4, a5, a6, a7}	T60:{a12, a13, a14, a15}	T61:{a12, a13, a14, a15}	T62:{a12, a13, a14, a15}	T63:{a12, a13, a14, a15}																															
32	T64:{a0, a1, a2, a3}	T65:{a0, a1, a2, a3}	T66:{a0, a1, a2, a3}	T67:{a0, a1, a2, a3}	T64:{a8, a9, a10, a11}	T65:{a8, a9, a10, a11}	T66:{a8, a9, a10, a11}	T67:{a8, a9, a10, a11}																															
..																																							
47	T92:{a4, a5, a6, a7}	T93:{a4, a5, a6, a7}	T94:{a4, a5, a6, a7}	T95:{a4, a5, a6, a7}	T92:{a12, a13, a14, a15}	T93:{a12, a13, a14, a15}	T94:{a12, a13, a14, a15}	T95:{a12, a13, a14, a15}																															
48	T96:{a0, a1, a2, a3}	T97:{a0, a1, a2, a3}	T98:{a0, a1, a2, a3}	T99:{a0, a1, a2, a3}	T96:{a8, a9, a10, a11}	T97:{a8, a9, a10, a11}	T98:{a8, a9, a10, a11}	T99:{a8, a9, a10, a11}																															
..																																							
63	T124:{a4, a5, a6, a7}	T125:{a4, a5, a6, a7}	T126:{a4, a5, a6, a7}	T127:{a4, a5, a6, a7}	T124:{a12, a13, a14, a15}	T125:{a12, a13, a14, a15}	T126:{a12, a13, a14, a15}	T127:{a12, a13, a14, a15}																															

(tid % 128) : fragments

WGMMMA continued

- WGMMMA operations are *mixed-precision*: operand matrices are lower bit-width than the accumulator matrix (e.g., FP16 or FP8 operand matrices with FP32 accumulator).
- Mixed-precision should be contrasted with *mixed-input* where the A operand is of a different bit-width than the B operand, e.g. as in weight-only quantization.
- The ISA doesn't currently support mixed-input wgmma.
- However, the CUDA programmer can support mixed-input wgmma on the software side through upcasting the narrower bit-width operand.
- Then, *layout conformance* becomes a challenge with using wgmma.

(a) **Mixed-input
(LLM application)**



Software solution

Datatype conversion
Layout conformance

(b) **Mixed-precision
(Tensor Cores in Hardware)**

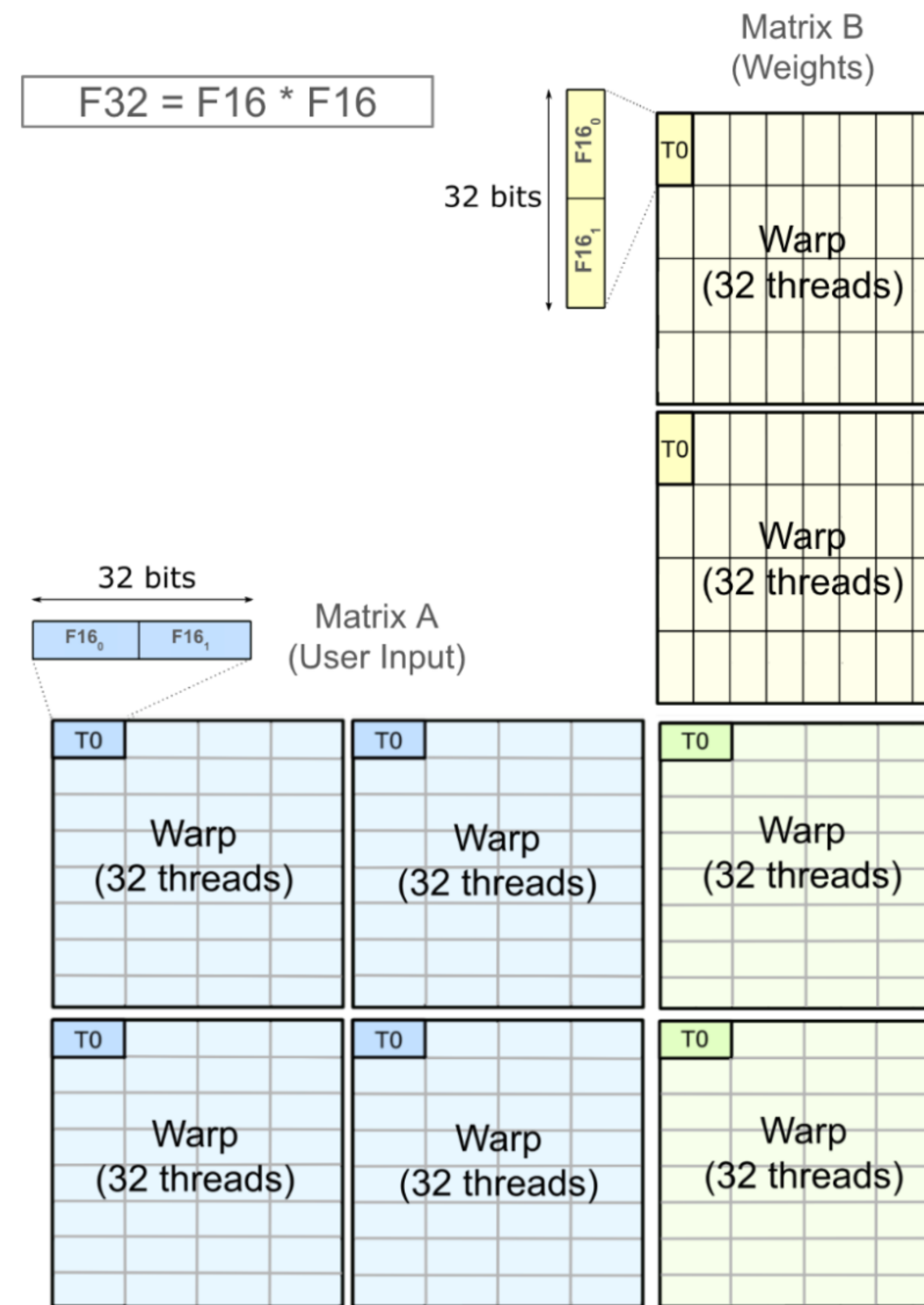


Illustration of upcasting 8 to 16 bits in mixed-input GEMM. Source: "Developing CUDA kernels to push Tensor Cores to the Absolute Limit on NVIDIA A100."

References

- "Neural Network Quantization & Number Formats From First Principles", Dylan Patel. <https://www.semianalysis.com/p/neural-network-quantization-and-number>
- "Mixed-input matrix multiplication performance optimizations", Manish Gupta. <https://blog.research.google/2024/01/mixed-input-matrix-multiplication.html>