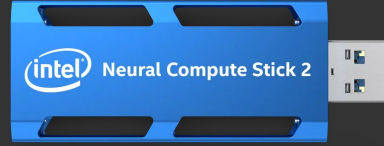
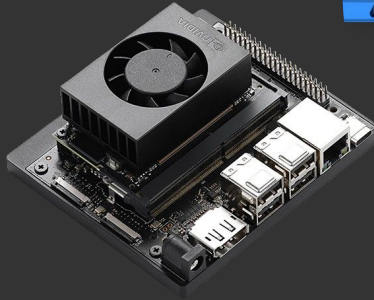
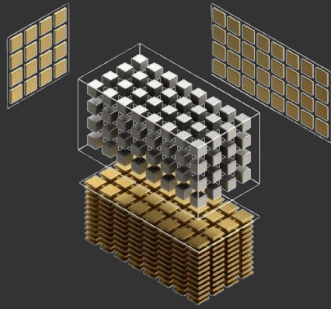


# Deep Learning Tools for Edge Devices



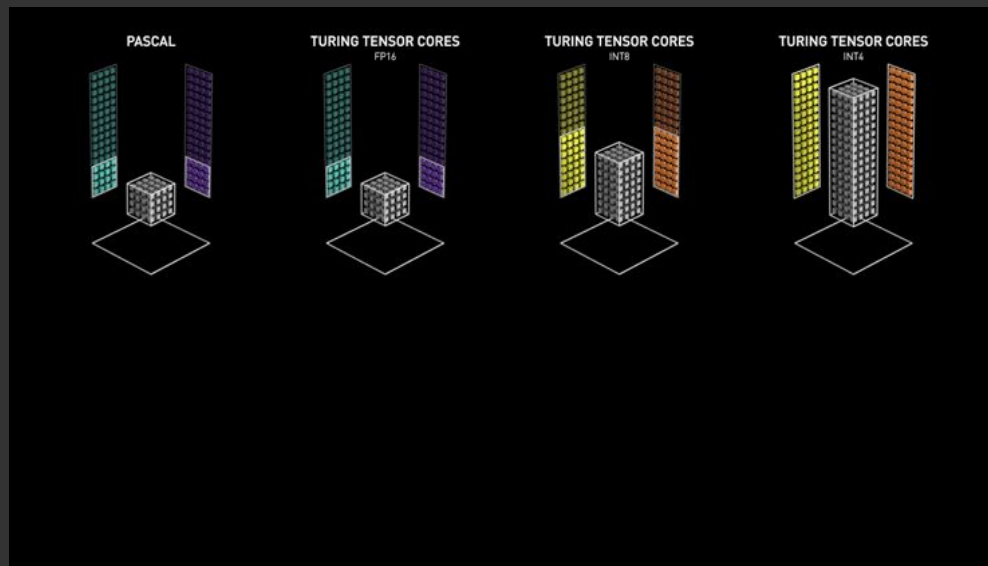
## Lecture 2: Tensor Cores

# What are Tensor Cores?

Dedicated Specialized Hardware Accelerator Cores

Used for Matrix Multiplication and Convolution Operations

Capable of Performing 4x4 matrix multiplication in one go



Source: NVIDIA

# SAXPY

Stands for Single Precision A.X + Y

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32      FP16      FP16      FP16 or FP32

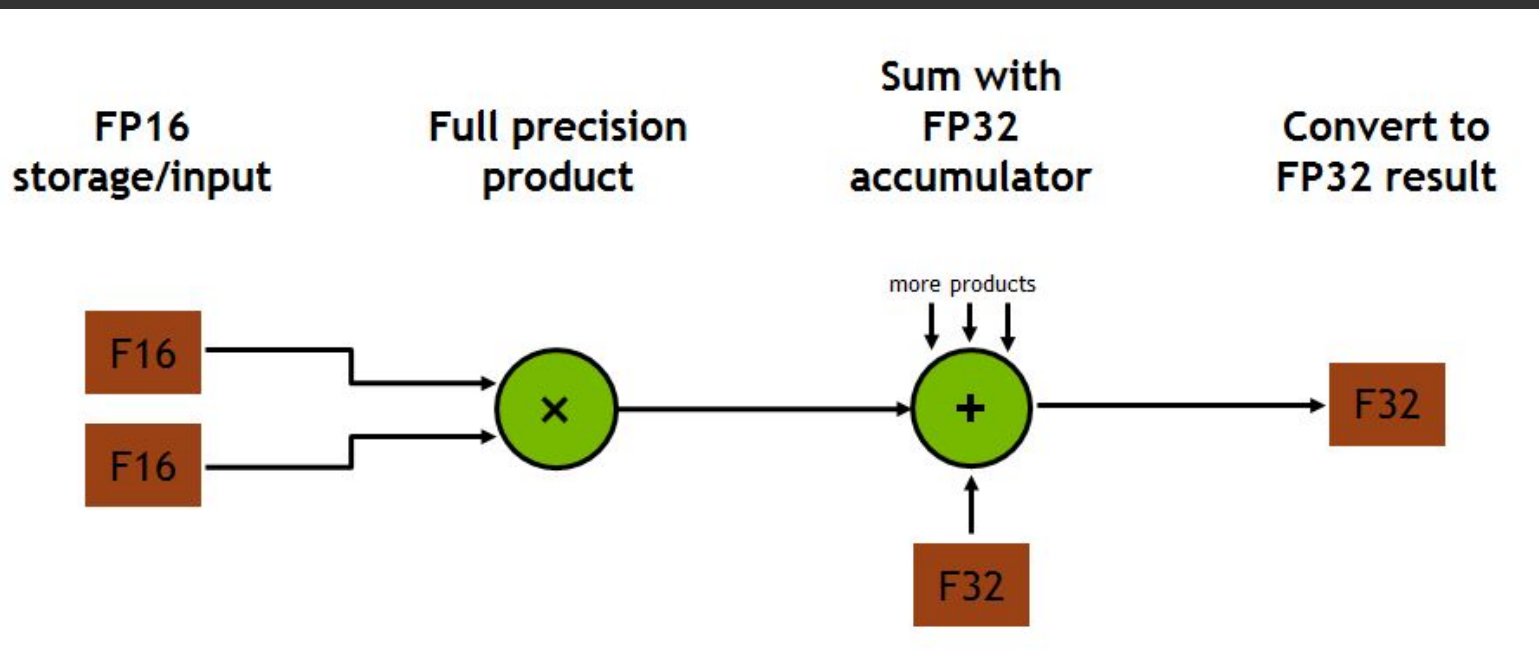
Source: NVIDIA

A combination of scalar multiplication and vector/matrix addition

4x4 Matrix Multiplication -> 64 Scalar multiplications

Requires 64 Cuda core vs 1 Tensor core

# Automatic Mixed Precision



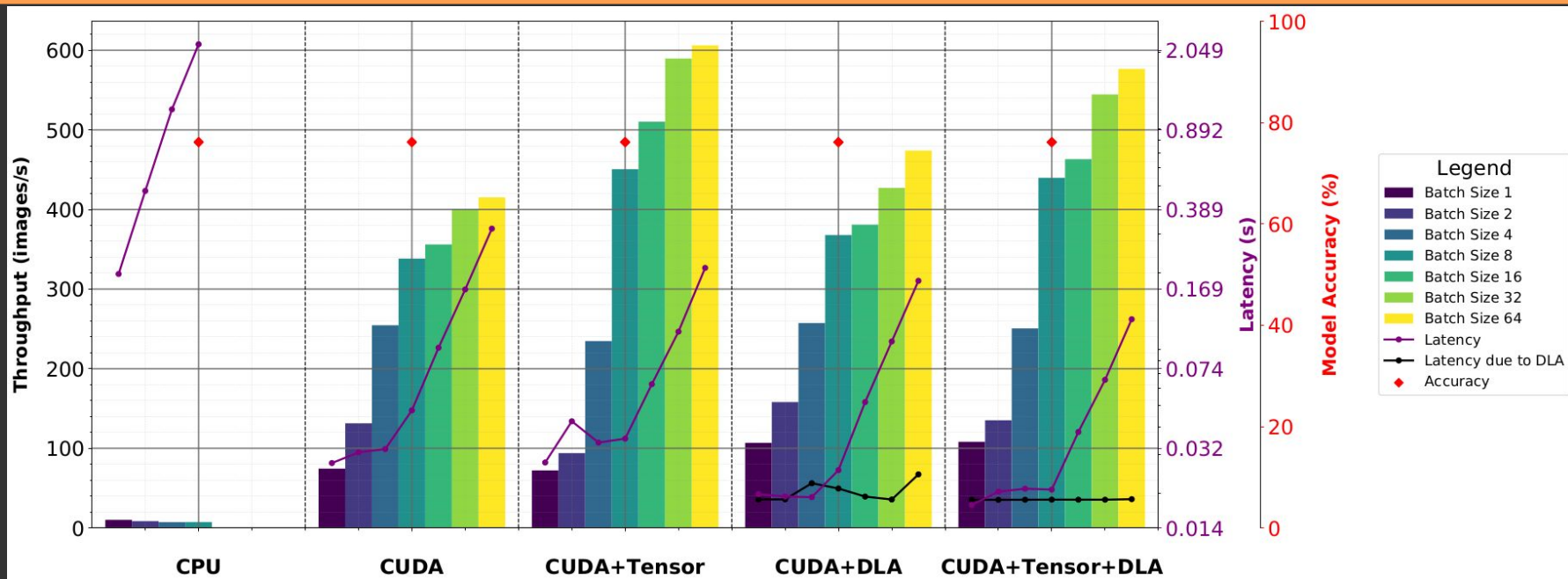
# Example

---

Now, let us consider an example,

Model:	ResNet-50
Framework:	Pytorch
Precision:	FP16, FP32, AMP
Batch Sizes:	1, 2, 4, 8, 16, 32, 64
Concurrency:	2

# Example Results



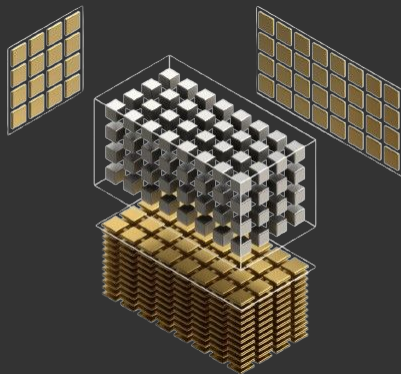
*Evaluating Multi-Instance DNN Inferencing on Multiple Accelerators of an Edge Device*, by M. Tayal, Y. Simmhan  
(IEEE International Conference on High Performance Computing, Data, and Analytics)

*Any References?*

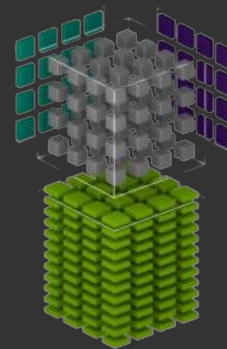
# Useful Articles and Videos



YT Channel: 0Mean1Sigma [\[Video\]](#)



Article: Programming Tensor Cores in CUDA 9 [\[Article\]](#)



Article: Tips for Optimizing GPU Performance Using Tensor Cores [\[Article\]](#)



Thank You

---