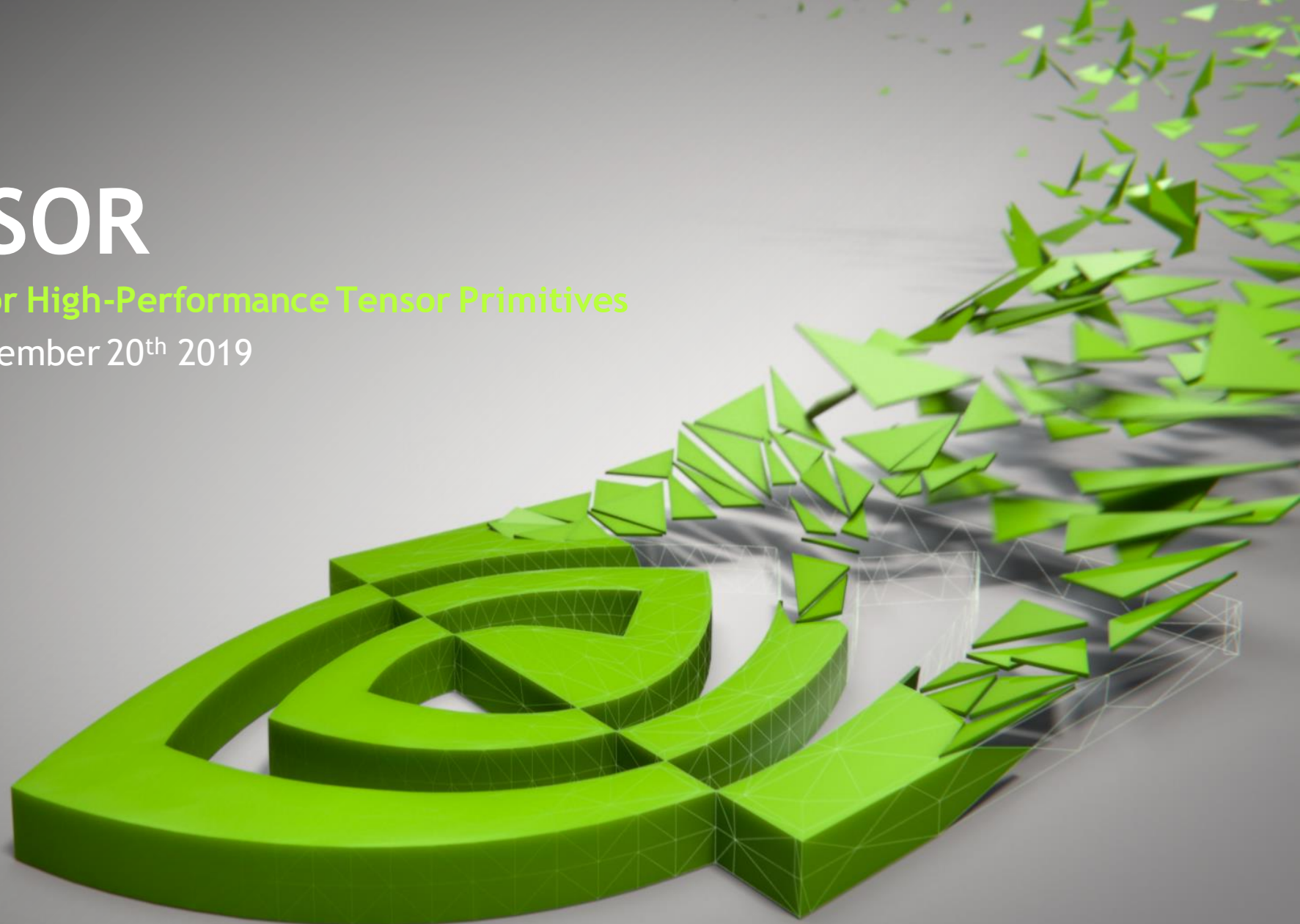


# CUTENSOR

A CUDA Library for High-Performance Tensor Primitives

Paul Springer, November 20<sup>th</sup> 2019



# CONTRIBUTERS

Chenhan Yu

Markus Höhnerbach

Andrew Kerr




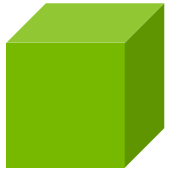
Timothy Costa

Manish Gupta

Miguel Ferrer Avila

Alex Fit-Florea

# WHAT IS A TENSOR?

- mode-0: scalar  $\alpha$  
- mode-1: vector  $A_i$  
- mode-2: matrix  $A_{i,j}$  
- mode-n: general tensor  $A_{i,j,k}$  

# WHAT IS A TENSOR?

- mode-0: scalar

$\alpha$



- mode-1: vector

$A_i$



- mode-2: matrix

$A_{i,j}$



- mode-n: general tensor

$A_{i,j,k,l}$



# WHAT IS A TENSOR?

- mode-0: scalar

$\alpha$



- mode-1: vector

$A_i$



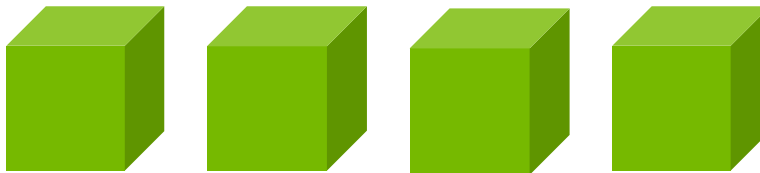
- mode-2: matrix

$A_{i,j}$



- mode-n: general tensor

$A_{i,j,k,l,m}$



# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector

$$\text{green bar} = \alpha \text{ blue bar} + \text{orange bar}$$

# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector

$$\text{Green Vector} = \alpha \text{ Blue Vector} + \text{Orange Vector}$$

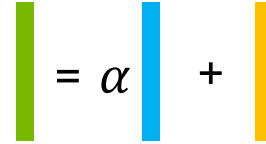
- 1972 - BLAS Level 2: Matrix-Vector

$$\text{Green Vector} = \text{Blue Matrix} * \text{Orange Vector}$$

# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector



A diagram illustrating a vector-vector operation. It consists of a green vertical bar, followed by an equals sign, a scalar symbol  $\alpha$ , a blue vertical bar, a plus sign, and an orange vertical bar.

$$\text{Green Bar} = \alpha \text{ Blue Bar} + \text{Orange Bar}$$

- 1972 - BLAS Level 2: Matrix-Vector



A diagram illustrating a matrix-vector operation. It consists of a green vertical bar, followed by an equals sign, a blue square, an asterisk, and an orange vertical bar.

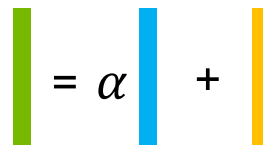
$$\text{Green Bar} = \text{Blue Square} * \text{Orange Bar}$$



# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector



A diagram illustrating a vector-vector operation. It shows a green vertical bar on the left, followed by an equals sign, a scalar symbol  $\alpha$ , a blue vertical bar, a plus sign, and an orange vertical bar.

$$\text{Green Bar} = \alpha \text{ Blue Bar} + \text{Orange Bar}$$


- 1972 - BLAS Level 2: Matrix-Vector



A diagram illustrating a matrix-vector operation. It shows a green vertical bar on the left, followed by an equals sign, a blue square, an asterisk, and an orange vertical bar.

$$\text{Green Bar} = \text{Blue Square} * \text{Orange Bar}$$

- 1980 - BLAS Level 3: Matrix-Matrix



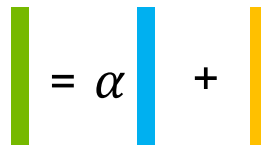
A diagram illustrating a matrix-matrix operation. It shows four green vertical bars on the left, followed by an equals sign, a blue square, an asterisk, and four orange vertical bars.

$$\text{Four Green Bars} = \text{Blue Square} * \text{Four Orange Bars}$$

# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector



A diagram illustrating a vector-vector operation. It shows a green vertical bar representing a result vector, followed by an equals sign, a scalar symbol  $\alpha$ , a blue vertical bar representing a scalar multiplier, a plus sign, and an orange vertical bar representing an input vector.

$$\text{Green Bar} = \alpha \text{ Blue Bar} + \text{Orange Bar}$$


- 1972 - BLAS Level 2: Matrix-Vector



A diagram illustrating a matrix-vector operation. It shows a green vertical bar representing a result vector, followed by an equals sign, a blue square representing a matrix, an asterisk, and an orange vertical bar representing an input vector.

$$\text{Green Bar} = \text{Blue Square} * \text{Orange Bar}$$

- 1980 - BLAS Level 3: Matrix-Matrix



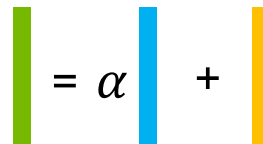
A diagram illustrating a matrix-matrix operation. It shows a green square representing a result matrix, followed by an equals sign, a blue square representing a first matrix, an asterisk, and a yellow square representing a second matrix.

$$\text{Green Square} = \text{Blue Square} * \text{Yellow Square}$$

# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector



A diagram illustrating a vector-vector operation. It shows a green vertical bar representing a result vector, followed by an equals sign, a scalar  $\alpha$ , another green vertical bar representing a first input vector, a plus sign, and a yellow vertical bar representing a second input vector.

$$\text{Green Bar} = \alpha \text{ Green Bar} + \text{Yellow Bar}$$

- 1972 - BLAS Level 2: Matrix-Vector



A diagram illustrating a matrix-vector operation. It shows a green vertical bar representing a result vector, followed by an equals sign, a blue square representing a matrix, an asterisk, and a yellow vertical bar representing an input vector.

$$\text{Green Bar} = \text{Blue Square} * \text{Yellow Bar}$$

- 1980 - BLAS Level 3: Matrix-Matrix



A diagram illustrating a matrix-matrix operation. It shows a green square representing a result matrix, followed by an equals sign, a blue square representing a first input matrix, an asterisk, and a yellow square representing a second input matrix.

$$\text{Green Square} = \text{Blue Square} * \text{Yellow Square}$$

- Now? - BLAS Level 4: Tensor-Tensor



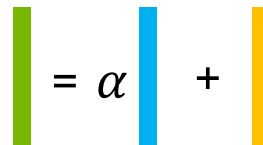
A diagram illustrating a tensor-tensor operation. It shows a stack of three green squares representing a result tensor, followed by an equals sign, a stack of three blue squares representing a first input tensor, an asterisk, and a stack of three yellow squares representing a second input tensor.

$$\text{Stack of 3 Green Squares} = \text{Stack of 3 Blue Squares} * \text{Stack of 3 Yellow Squares}$$

# BASIC LINEAR SUBPROGRAMS

## A Success Story

- 1969 - BLAS Level 1: Vector-Vector



A diagram illustrating a vector-vector operation. It shows a green vertical bar representing a result vector, followed by an equals sign, a scalar  $\alpha$ , another green vertical bar representing a first input vector, a plus sign, and a yellow vertical bar representing a second input vector.

$$\text{Green Bar} = \alpha \text{ Green Bar} + \text{Yellow Bar}$$

- 1972 - BLAS Level 2: Matrix-Vector



A diagram illustrating a matrix-vector operation. It shows a green vertical bar representing a result vector, followed by an equals sign, a blue square representing a matrix, an asterisk, and a yellow vertical bar representing an input vector.

$$\text{Green Bar} = \text{Blue Square} * \text{Yellow Bar}$$

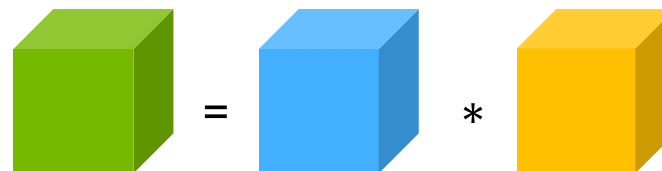
- 1980 - BLAS Level 3: Matrix-Matrix



A diagram illustrating a matrix-matrix operation. It shows a green square representing a result matrix, followed by an equals sign, a blue square representing a first input matrix, an asterisk, and a yellow square representing a second input matrix.

$$\text{Green Square} = \text{Blue Square} * \text{Yellow Square}$$

- Now? - BLAS Level 4: Tensor-Tensor



A diagram illustrating a tensor-tensor operation. It shows a green 3D cube representing a result tensor, followed by an equals sign, a blue 3D cube representing a first input tensor, an asterisk, and a yellow 3D cube representing a second input tensor.

$$\text{Green Cube} = \text{Blue Cube} * \text{Yellow Cube}$$

# CUTENSOR

## A High-Performance CUDA Library for Tensor Primitives

- Tensor contractions (generalization of matrix-matrix multiplication)

$$\text{D} = \sum ( \text{A} * \text{B} ) + \text{C}$$

- Tensor reductions

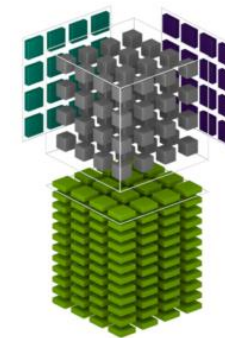
$$\text{D} = \sum ( \text{A} ) + \text{C}$$

- Element-wise operations (e.g., permutations, additions)

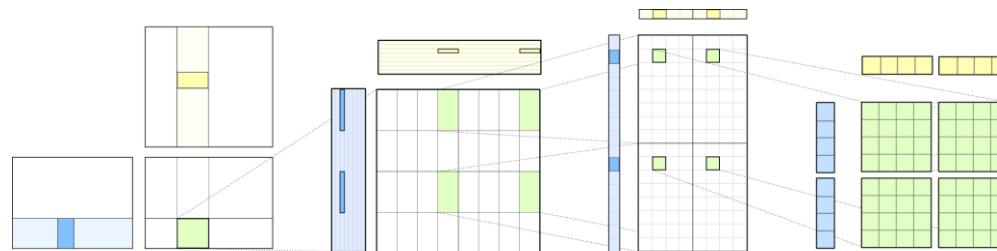
$$\text{D} = \text{A} + \text{B} + \text{C}$$

# CUTENSOR

## Key Features



- Transpose-free contractions
- Arbitrary data layouts
- Extensive mixed-precision support
  - Complex-times-Real
  - FP64 data + FP32 compute
  - FP32 data + FP16 compute
- Tensor-Core support
  - Built on CUTLASS 2.0
- Flexible and modern interface
  - No mallocs inside the lib
  - Tensor Cores active by default



# TENSORS ARE UBIQUITOUS

Potential Use Cases: HPC & AI

 PyTorch

 Pyro

 TensorFlow

Scikit-CUDA

 **NWCHEM**  
HIGH-PERFORMANCE COMPUTATIONAL  
CHEMISTRY SOFTWARE

TAL-SH

 **ITENSOR**

 **julia**

 CuPy

# Tensor Contractions



# TENSOR CONTRACTIONS

## Examples



A diagram illustrating a tensor contraction. It shows a green cube labeled 'D' on the left, followed by an equals sign, then a summation symbol (Σ) in parentheses. Inside the parentheses is a blue cube labeled 'A' followed by an asterisk (\*) and a yellow cube labeled 'B'. To the right of the parentheses is a plus sign (+) and an orange cube labeled 'C'.

- Einstein notation (einsum)
  - Modes that appear in A and B are contracted


- Examples

- $D_{m,n} = \alpha \sum_k A_{m,k} * B_{k,n}$

// GEMM

# TENSOR CONTRACTIONS

## Examples


$$D = \sum (A * B) + C$$

- Einstein notation (einsum)
  - Modes that appear in A and B are contracted

- Examples

- $D_{m,n} = \alpha A_{m,k} * B_{k,n}$

// GEMM

- $D_{m_1,n,m_2} = \alpha A_{m_1,k,m_2} * B_{k,n}$

// Tensor Contraction

- $D_{m_1,n_1,n_2,m_2} = \alpha A_{m_1,k,m_2} * B_{k,n_2,n_1}$


// Tensor Contraction

- $D_{m_1,n_1,n_2,m_2} = \alpha A_{m_1,k_1,m_2,k_2} * B_{k_1,k_2,n_2,n_1}$

// Multi-mode Tensor Contraction

# TENSOR CONTRACTIONS

Examples (cont.)

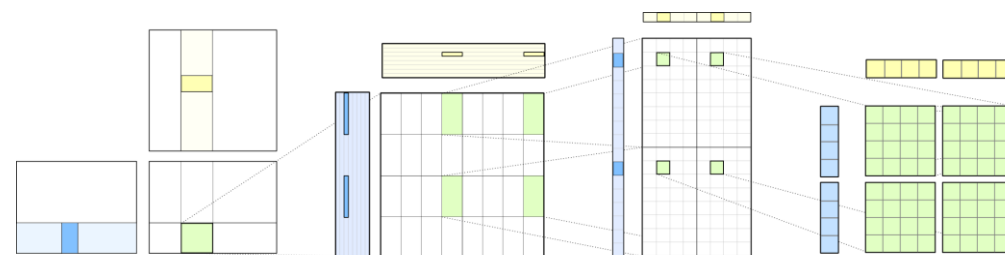
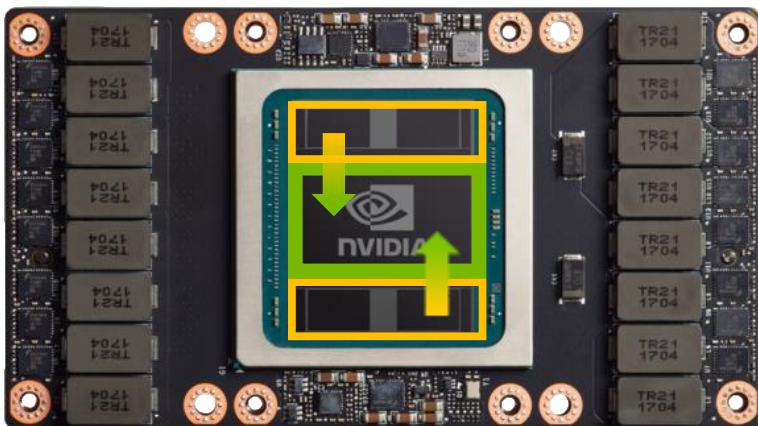

$$D = \sum (A * B) + C$$

- Examples

- $D_{m,n} = \alpha A_m * B_n$  // outer product
- $D_{m_1,n,m_2} = \alpha A_{m_1,m_2} * B_n$  // outer product
- $D_{m_1,m_2} = \alpha A_{m_1,k_1,m_2,k_2} * B_{k_1,k_2}$  // GEMV-like tensor contraction
- $D_{m_1,n_1,l_1} = \alpha A_{m_1,k,l_1} * B_{k,n_1,l_1}$  // batched GEMM
- $D_{m_1,n_1,l_1,n_2,m_2} = \alpha A_{m_1,k,l_1,m_2} * B_{k,n_2,n_1,l_1}$  // single-mode batched tensor contraction
- $D_{m_1,n_1,l_1,n_2,m_2,l_2} = \alpha A_{m_1,k,l_2,l_1,m_2} * B_{k,n_2,n_1,l_1,l_2}$  // multi-mode batched tensor contraction

# TENSOR CONTRACTIONS

## Key Challenges



CUTLASS 2.0

- Keep the **fast FPUs** busy
  - Reuse data in **shared memory & registers** as much as possible
  - Coalesced accesses to/from **global memory**

# TENSOR CONTRACTIONS

## Key Challenges

- Loading a scalar

$\alpha$



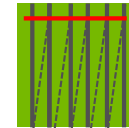
- Loading a vector

$A_i$



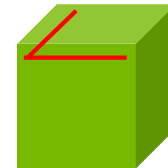
- Loading a matrix

$A_{i,j}$



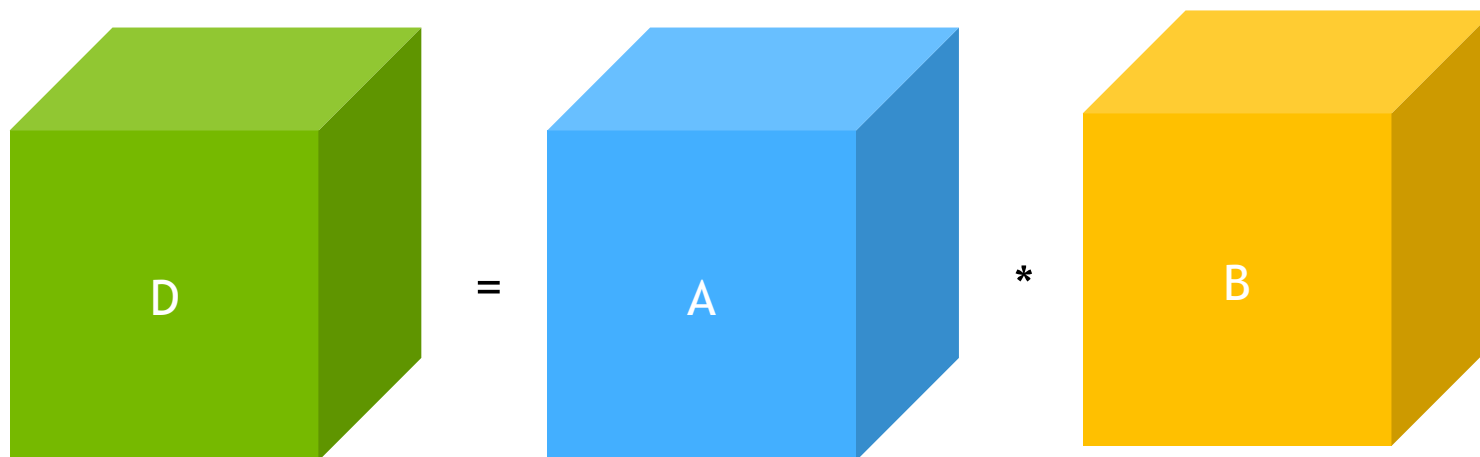
- Loading a general tensor

$A_{i,j,k}$



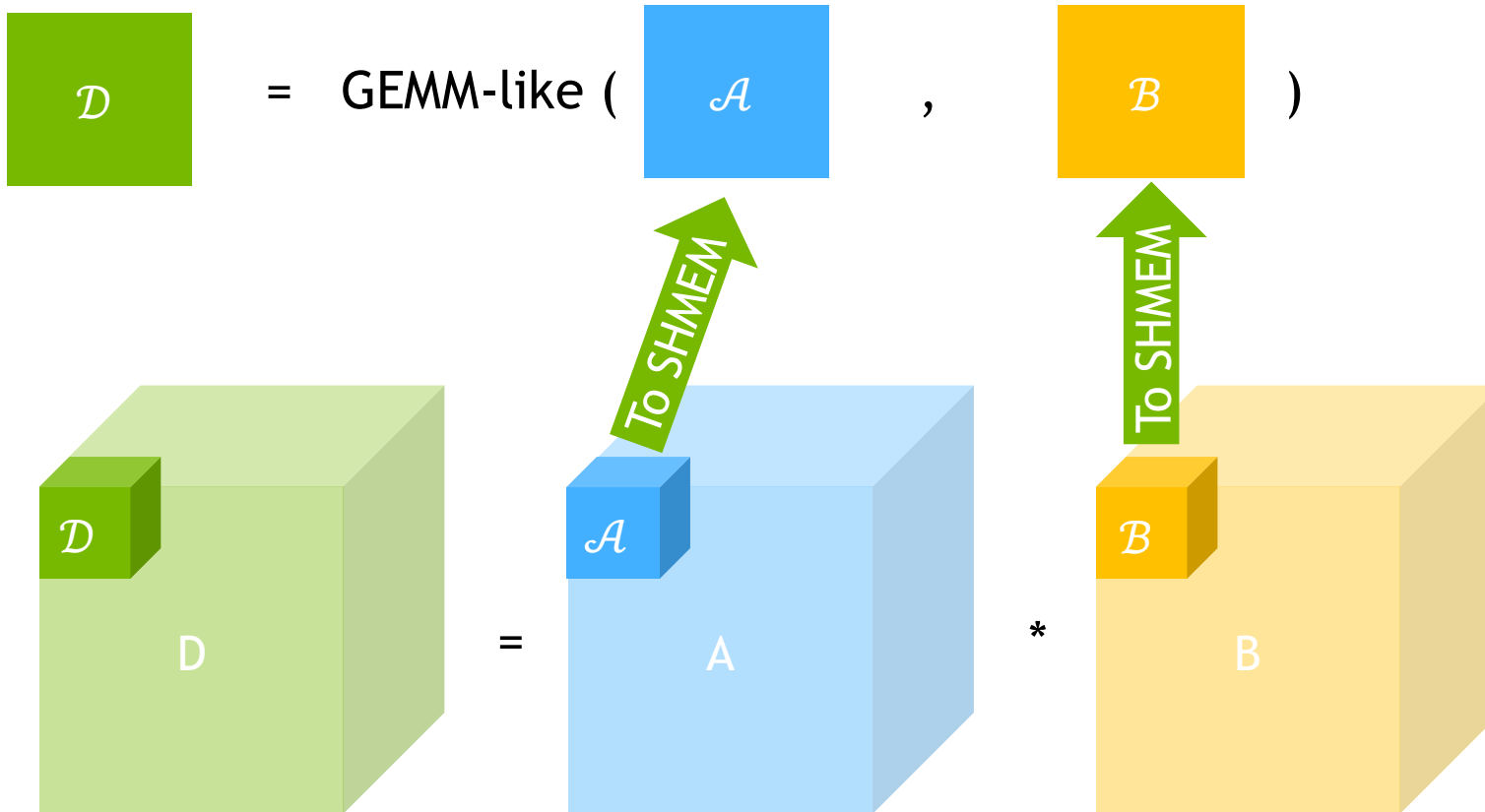
# TENSOR CONTRACTIONS

Technical insight



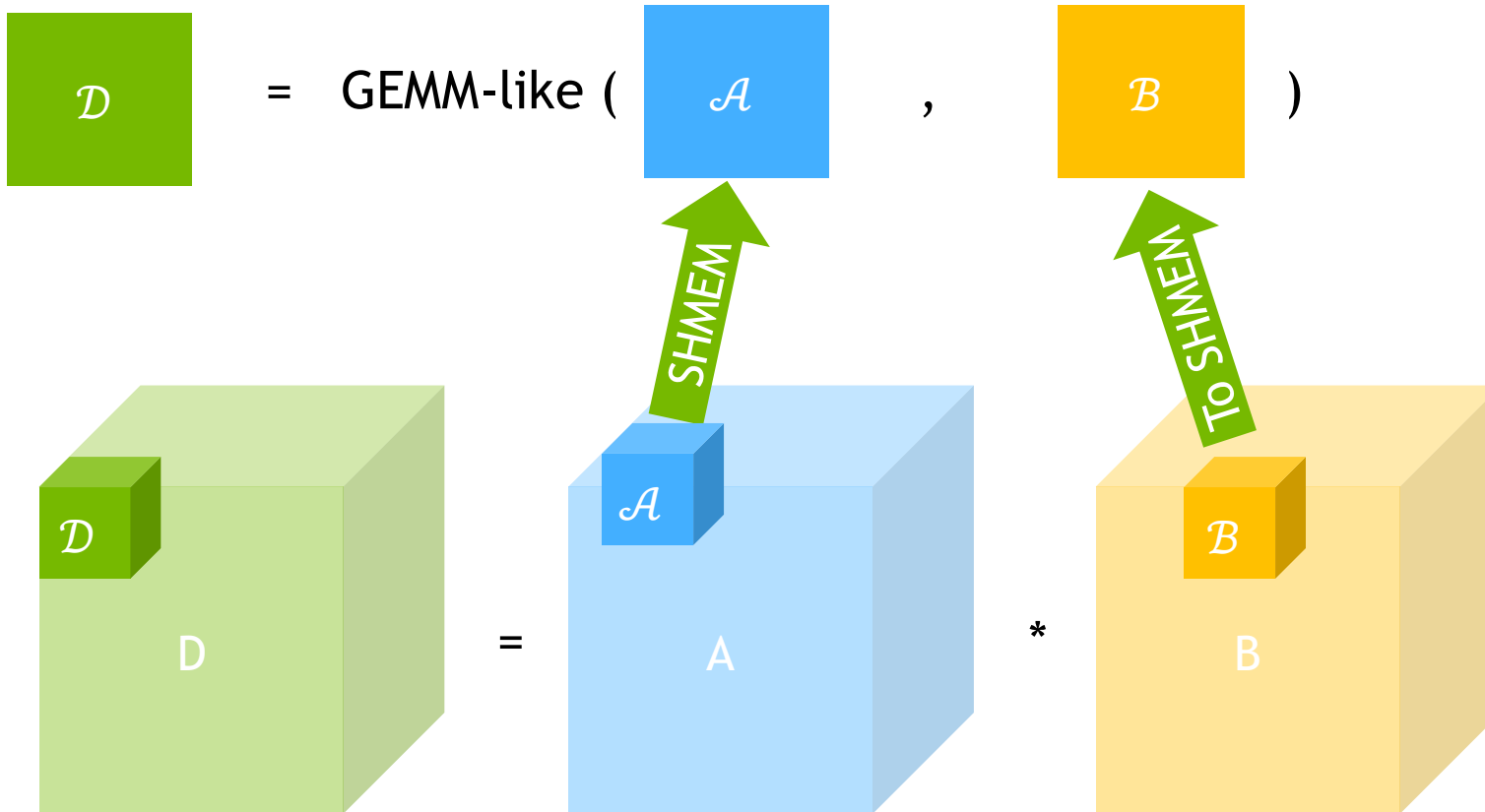
# TENSOR CONTRACTIONS

Technical insight



# TENSOR CONTRACTIONS

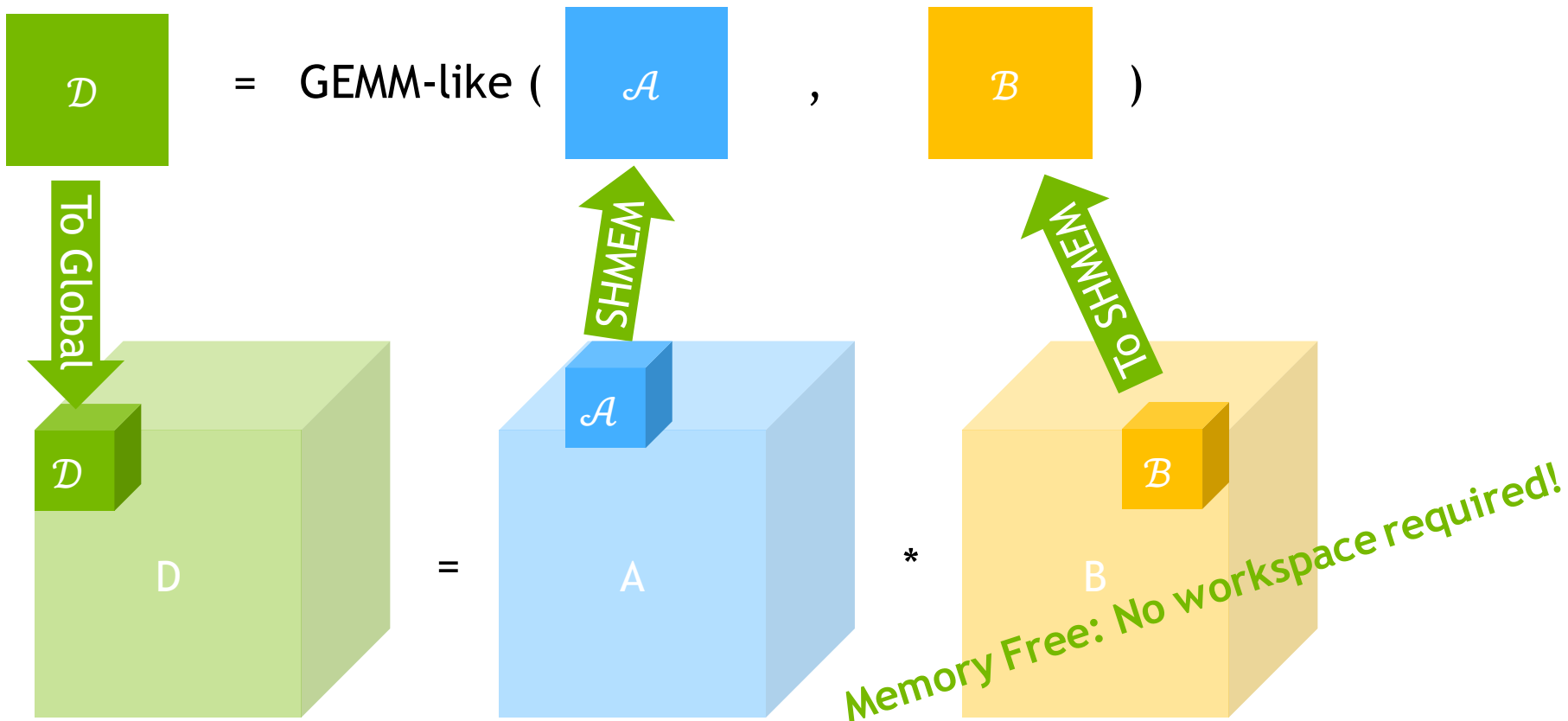
Technical insight





# TENSOR CONTRACTIONS

Technical insight



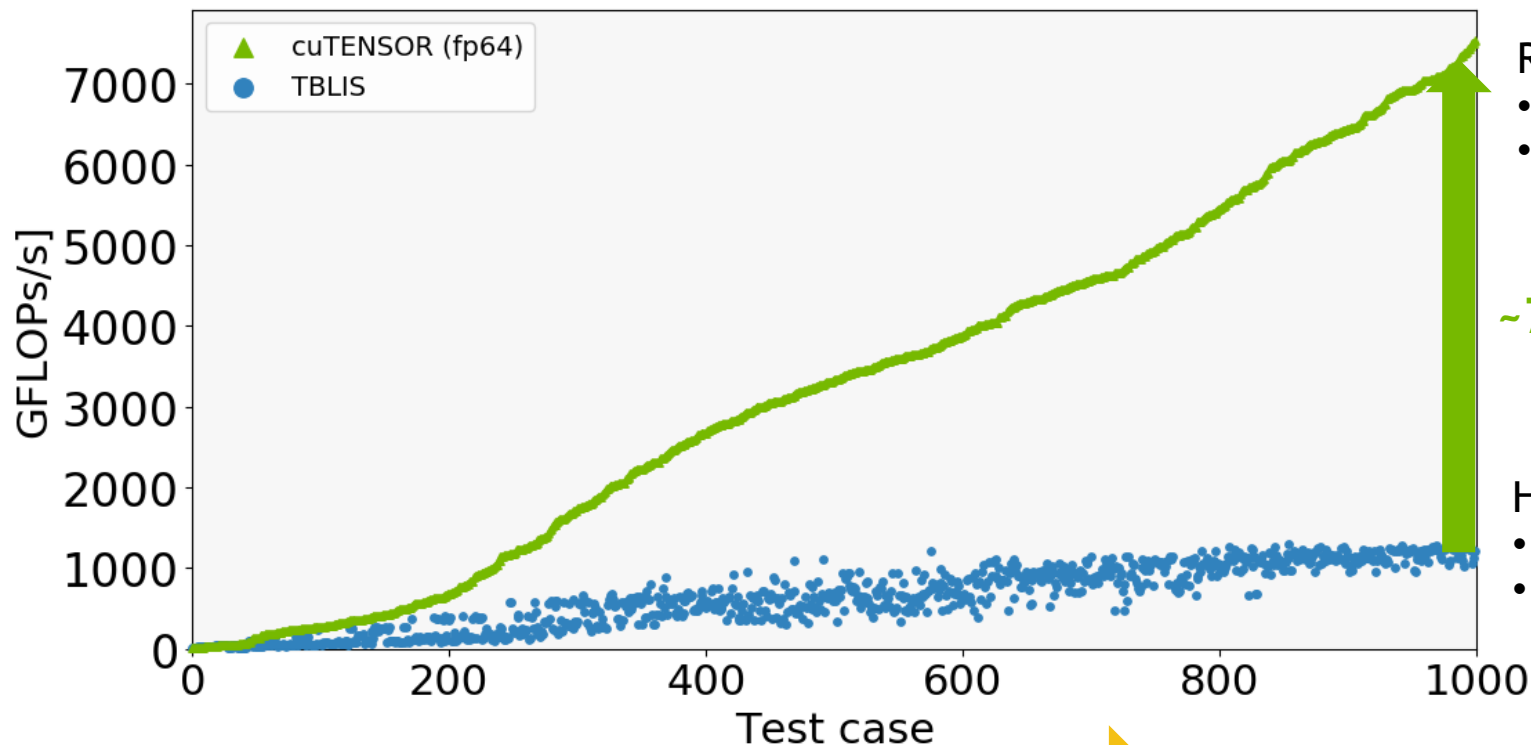
# TENSOR CONTRACTIONS

## Performance Guidelines

- Arrange modes similarly in all tensors
  - E.g.,  $C_{a,b,c} = A_{a,k,c} B_{k,b}$  is preferable to  $C_{a,b,c} = A_{c,k,a} B_{k,b}$
- Keep the extent of the fastest-varying mode as large as possible
  - E.g.,  $C_{a,b,c} \in R^{1000 \times 100 \times 10}$  is preferable to  $C_{c,b,a} \in R^{10 \times 100 \times 1000}$
- Keep batched modes as the slowest-varying modes (i.e., with the largest strides)
  - E.g.,  $C_{a,b,c,l} = A_{a,k,c,l} B_{k,b,l}$  is preferable to  $C_{a,b,c,l} = A_{a,k,l,c} B_{l,k,b}$

# PERFORMANCE

## Tensor Contractions



Random tensor contractions:

- 3D to 6D tensors
- FP64

~7x Speedup

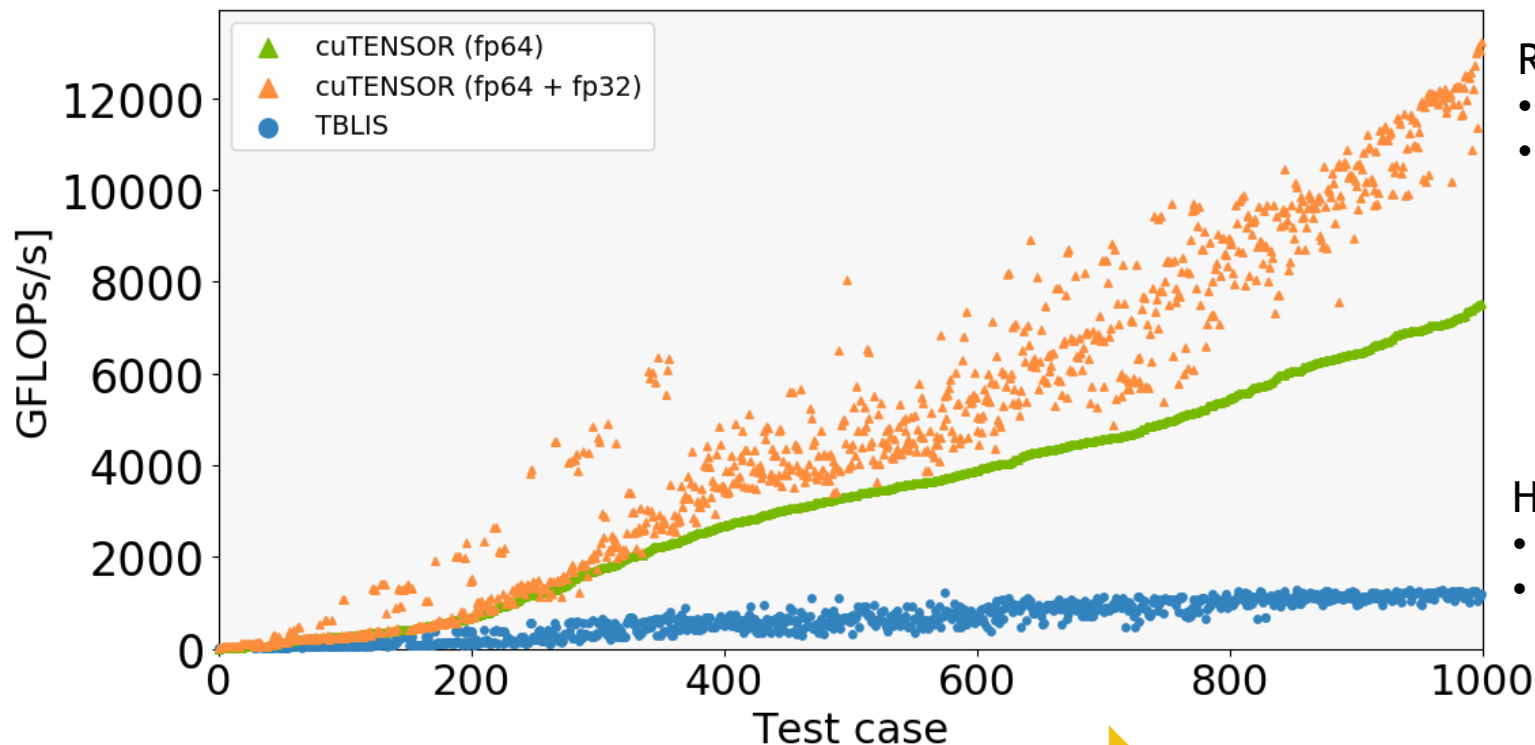
Hardware:

- Tesla GV100
- 2x Intel Xeon Platinum 8168

Arithmetic Intensity

# PERFORMANCE

## Tensor Contractions



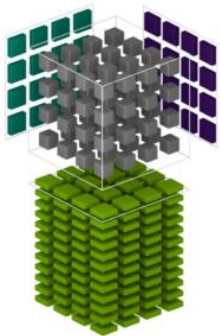
Random tensor contractions:

- 3D to 6D tensors
- FP64 (data) & FP32 (compute)

Hardware:

- Tesla GV100
- 2x Intel Xeon Platinum 8168

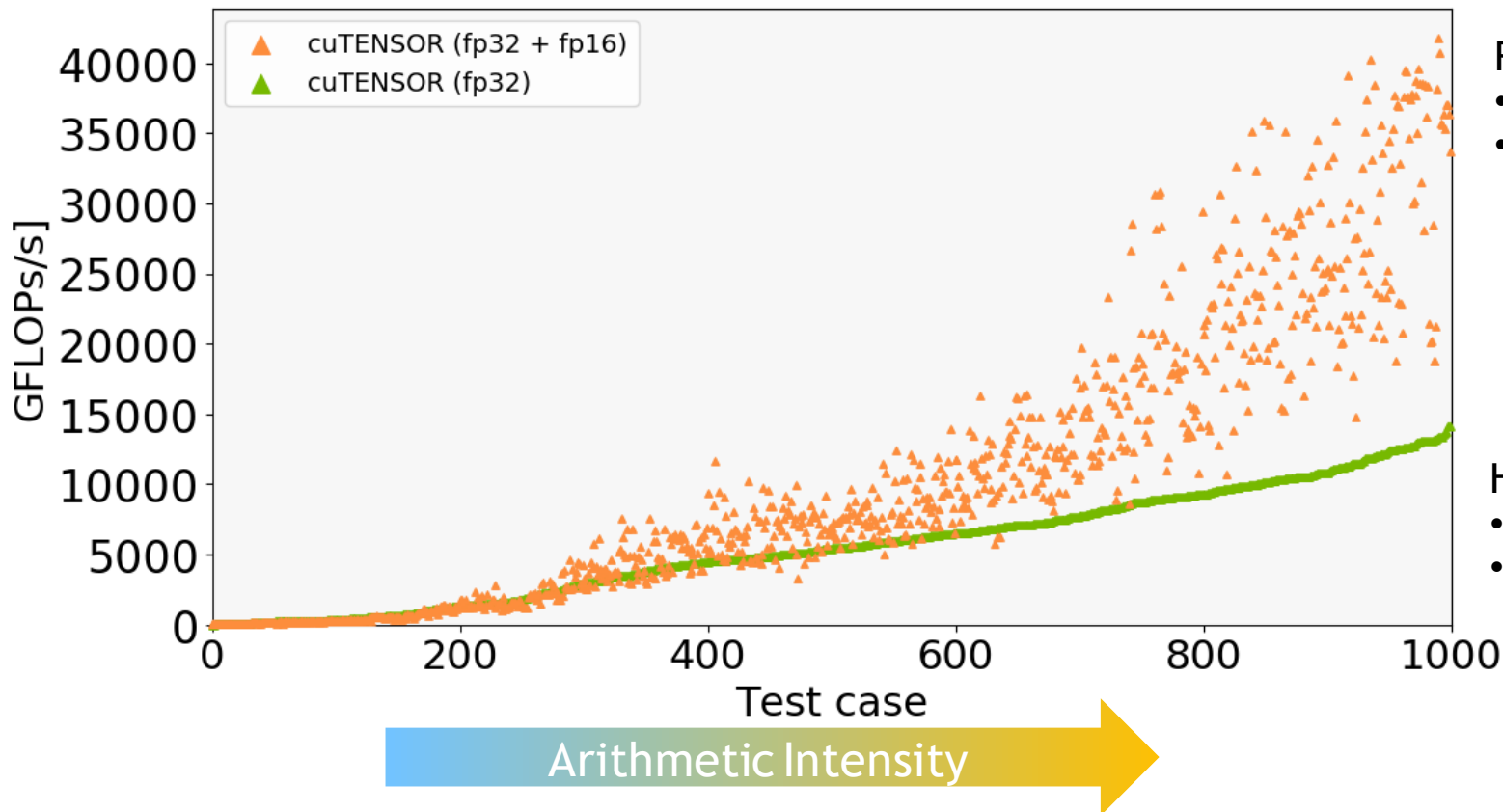
Arithmetic Intensity



# PERFORMANCE

## Tensor Contractions

$$\text{C} = \text{A} * \text{B}$$



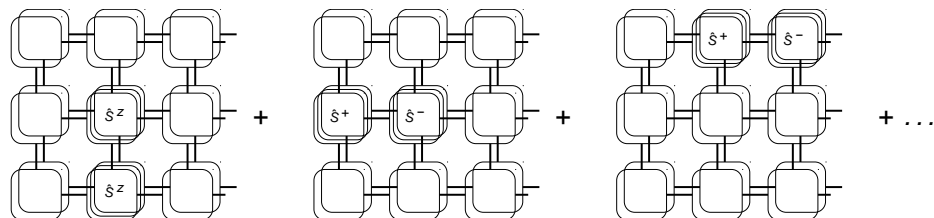
Random tensor contractions:

- 3D to 6D tensors
- FP32 (data) + Tensor Core
  - Accumulation in FP32
  - Inputs truncated to FP16

Hardware:

- Tesla GV100
- 2x Intel Xeon Platinum 8168

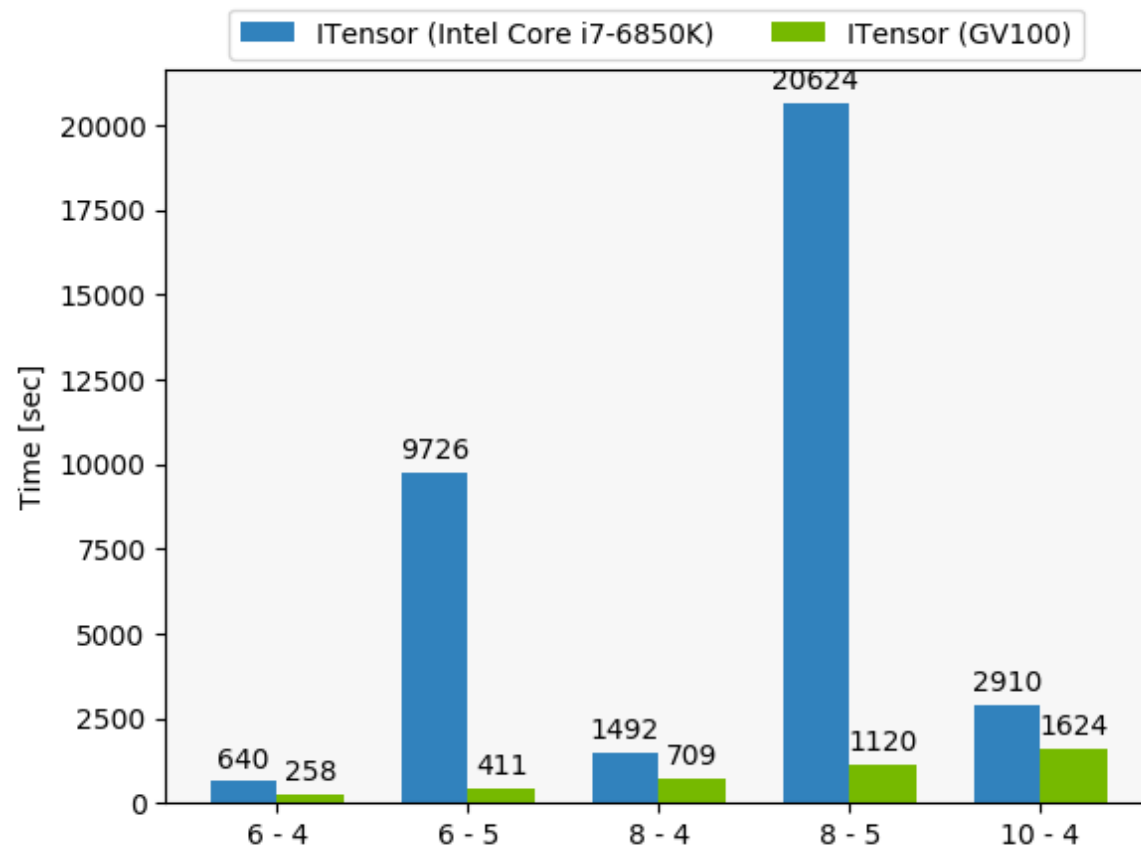
# ITENSOR - PEPS SIMULATION



Speedups: up to 23.6x

## Credits:

- Katharine Hyatt (Flatiron Institute) - <https://github.com/itensor/ITensorsGPU.jl>
- Miles Stoudenmire and Matt Fishman (Flatiron Institute) - <https://github.com/ITensor/ITensors.jl>
- Tim Besard (Ghent University) - <https://github.com/JuliaGPU>



\*results are based on cuTENSOR's early-access version

# Tensor Reductions

# TENSOR REDUCTION

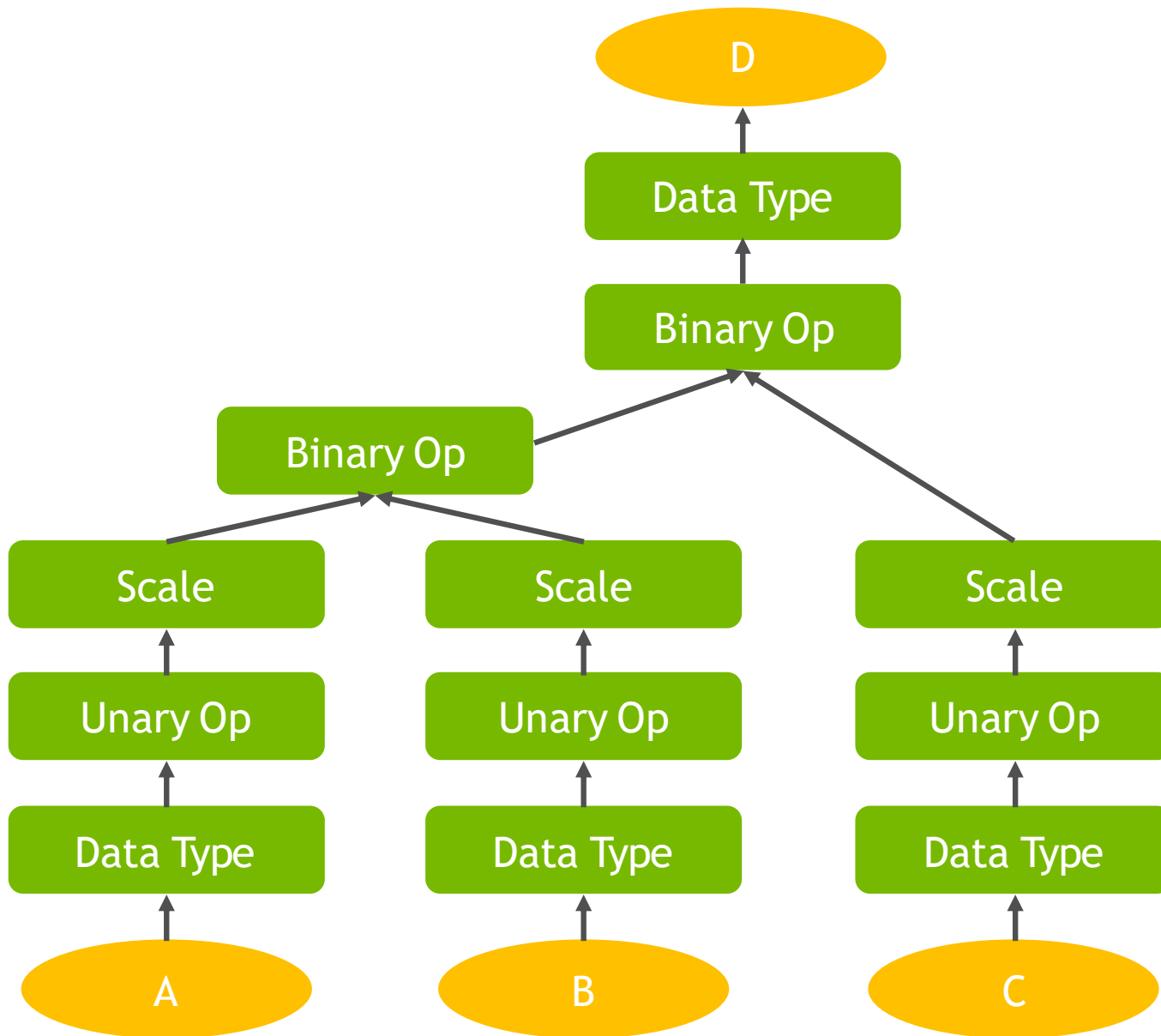
## Examples

$$\text{D} = \sum (\text{A}) + \text{C}$$

- $D_{w,h,n} = \alpha \text{ Reduce}(A_{c,w,h,n}, \text{ADD})$  // Reduction over mode c
- $D_{w,n} = \alpha \text{ Reduce}(A_{c,w,h,n}, \text{ADD})$  // Reduction over mode c and h
- $D_{w,n} = \alpha \text{ Reduce}(\text{RELU}(A_{c,w,h,n}), \text{ADD})$  // RELU + Reduction over mode c and h
- $D_{w,n} = \alpha \text{ Reduce}(\text{RELU}(A_{c,w,h,n}), \text{MAX})$  // RELU + Max-reduction over mode c and h

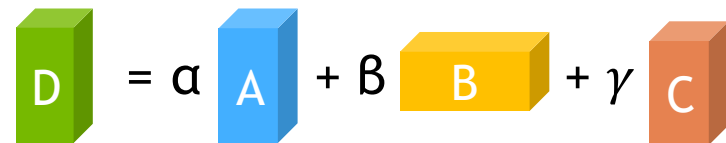


# Element-wise Operations



# ELEMENT-WISE TENSOR OPERATIONS

## Examples


$$D = \alpha A + \beta B + \gamma C$$

- $D_{w,h,c,n} = \alpha A_{c,w,h,n}$
- $D_{w,h,c,n} = \alpha A_{c,w,h,n} + \beta B_{c,w,h,n}$
- $D_{w,h,c,n} = \min(\alpha A_{c,w,h,n}, \beta B_{c,w,h,n})$
- $D_{w,h,c,n} = \alpha A_{c,w,h,n} + \beta B_{w,h,c,n} + \gamma C_{w,h,c,n}$
- $D_{w,h,c,n} = \alpha \text{RELU}(A_{c,w,h,n}) + \beta B_{w,h,c,n} + \gamma C_{w,h,c,n}$
- $D_{w,h,c,n} = \text{FP32}(\alpha \text{RELU}(A_{c,w,h,n}) + \beta B_{w,h,c,n} + \gamma C_{w,h,c,n})$

Enables users to fuse multiple element-wise calls.

# cuTENSOR's API

# CONTRACTION OPERATION

API



```
cutensorStatus_t cutensorInitContractionDescriptor ( const cutensorHandle_t *handle,
    cutensorContractionDescriptor_t *desc,
    const cutensorTensorDescriptor_t *descA, const int modeA[], const uint32_t alignmentRequirementA,
    const cutensorTensorDescriptor_t *descB, const int modeB[], const uint32_t alignmentRequirementB,
    const cutensorTensorDescriptor_t *descC, const int modeC[], const uint32_t alignmentRequirementC,
    const cutensorTensorDescriptor_t *descD, const int modeD[], const uint32_t alignmentRequirementD,
    cutensorComputeType_t typeCompute );

cutensorStatus_t cutensorInitContractionFind ( const cutensorHandle_t *handle,
    cutensorContractionFind *find, const cutensorAlgo_t algo );

cutensorStatus_t cutensorInitContractionPlan ( const cutensorHandle_t *handle,
    cutensorContractionPlan_t *plan, const cutensorContractionDescriptor_t *desc,
    const cutensorContractionFind *find, uint64_t workspaceSize );

cutensorStatus_t cutensorContraction( const cutensorHandle_t *handle,
    const cutensorContractionPlan_t *plan,
    const void *alpha, const void *A, const void *B,
    const void *beta, const void *C, void *D,
    void *workspace, uint64_t workspaceSize, cudaStream_t stream );
```

# CONTRACTION OPERATION

API



- $D_{m,n,u} = \alpha A_{m,k,u} * B_{k,n}$

```
cutensorInitContractionDescriptor ( &handle, &desc,  
                                     descA, { 'm', 'k', 'u' }, 256,  
                                     descB, { 'k', 'n' },      256,  
                                     descC, { 'm', 'n', 'u' }, 256,  
                                     descD, { 'm', 'n', 'u' }, 256,  
                                     CUTENSOR_R_MIN_F32 );
```

```
cutensorInitContractionFind ( &handle, &find, ... );
```

```
cutensorInitContractionPlan ( &handle, &plan, &desc, &find,... );
```

```
cutensorContraction ( handle,  
                      &plan,  
                      alpha, A, B,  
                      beta,  C, D,  
                      workspace, workspaceSize, stream );
```

# CONTRACTION OPERATION

## API

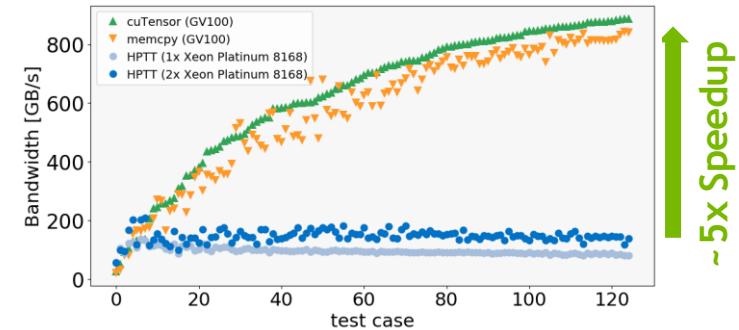
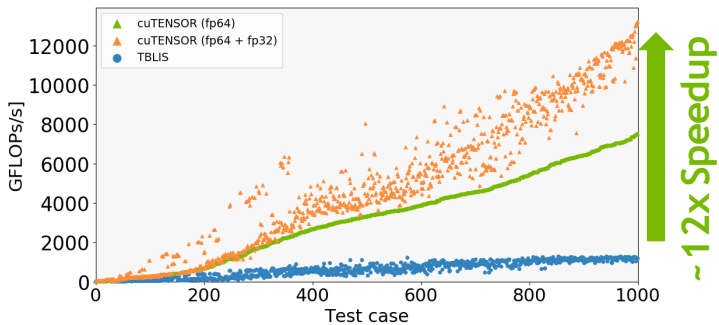


- $$D_{m,n,u} = \alpha A_{m,k,u} * B_{k,n}$$

```
cutensorInitContractionDescriptor ( &handle, &desc,           // Define Problem
                                   descA, { 'm', 'k', 'u' }, 256,
                                   descB, { 'k', 'n' },      256,
                                   descC, { 'm', 'n', 'u' }, 256,
                                   descD, { 'm', 'n', 'u' }, 256,
                                   CUTENSOR_R_MIN_F32 );
cutensorInitContractionFind ( &handle, &find, ... );           // Define search space
cutensorInitContractionPlan ( &handle, &plan, &desc, &find, ... ); // Initialize plan
cutensorContraction ( handle,                                     // Execute contraction
                      &plan,
                      alpha, A, B,
                      beta, C, D,
                      workspace, workspaceSize, stream );
```

# CONCLUSION

cuTENSOR: A CUDA library for high-performance tensor primitives



$$D = \sum (A * B) + C$$

$$D = \alpha A + \beta B + \gamma C$$

Available at: <https://developer.nvidia.com/cuTENSOR>

*Your feedback is highly appreciated.*





# REDUCTION OPERATION

API

$$\text{D} = \sum (\text{A}) + \text{C}$$

```
cutensorStatus_t cutensorReduction ( const cutensorHandle_t *handle,  
    const void *alpha, const void *A, const cutensorTensorDescriptor_t *descA, const int modeA[],  
    const void *beta,  const void *C, const cutensorTensorDescriptor_t *descC, const int modeC[],  
                                void *D, const cutensorTensorDescriptor_t *descD, const int modeD[],  
    cutensorOperator_t opReduce, cutensorComputeType_t typeCompute, cudaStream_t stream );
```

# REDUCTION OPERATION

API

$$D = \sum (A) + C$$

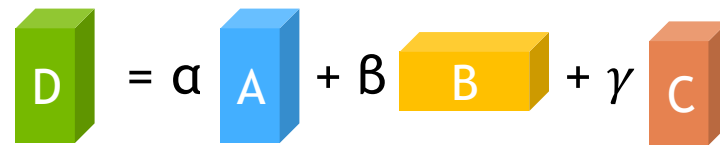
```
cutensorStatus_t cutensorReduction ( const cutensorHandle_t *handle,  
    const void *alpha, const void *A, const cutensorTensorDescriptor_t *descA, const int modeA[],  
    const void *beta,  const void *C, const cutensorTensorDescriptor_t *descC, const int modeC[],  
                                void *D, const cutensorTensorDescriptor_t *descD, const int modeD[],  
    cutensorOperator_t opReduce, cutensorComputeType_t typeCompute, cudaStream_t stream );
```

- $D_{w,h,n} = \alpha \text{Reduce}(A_{c,w,h,n}) + \beta C_{w,h,n}$

```
auto status = cutensorReduction ( handle,  
    alpha, A, descA, { 'c', 'w', 'h', 'n' },  
    beta,  C, descC, { 'w', 'h', 'n' },  
        D, descC, { 'w', 'h', 'n' },  
    CUTENSOR_OP_ADD, CUTENSOR_R_MIN_32F,  
    stream );
```

# ELEMENT-WISE OPERATION

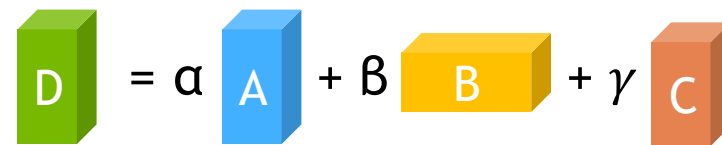
API


$$D = \alpha A + \beta B + \gamma C$$

```
cutensorStatus_t cutensorElementwiseTrinary ( const cutensorHandle_t *handle,  
    const void *alpha, const void *A, const cutensorTensorDescriptor_t *descA, const int modeA[],  
    const void *beta,  const void *B, const cutensorTensorDescriptor_t *descB, const int modeB[],  
    const void *gamma, const void *C, const cutensorTensorDescriptor_t *descC, const int modeC[],  
                                void *D, const cutensorTensorDescriptor_t *descD, const int modeD[],  
    cutensorOperator_t opAB, cutensorOperator_t opABC, cudaDataType_t typeCompute,  
    cudaStream_t stream );
```

# ELEMENT-WISE OPERATION

API


$$D = \alpha A + \beta B + \gamma C$$

```
cutensorStatus_t cutensorElementwiseTrinary ( const cutensorHandle_t *handle,  
    const void *alpha, const void *A, const cutensorTensorDescriptor_t *descA, const int modeA[],  
    const void *beta,  const void *B, const cutensorTensorDescriptor_t *descB, const int modeB[],  
    const void *gamma, const void *C, const cutensorTensorDescriptor_t *descC, const int modeC[],  
    void *D, const cutensorTensorDescriptor_t *descD, const int modeD[],  
    cutensorOperator_t opAB, cutensorOperator_t opABC, cudaDataType_t typeCompute,  
    cudaStream_t stream );
```

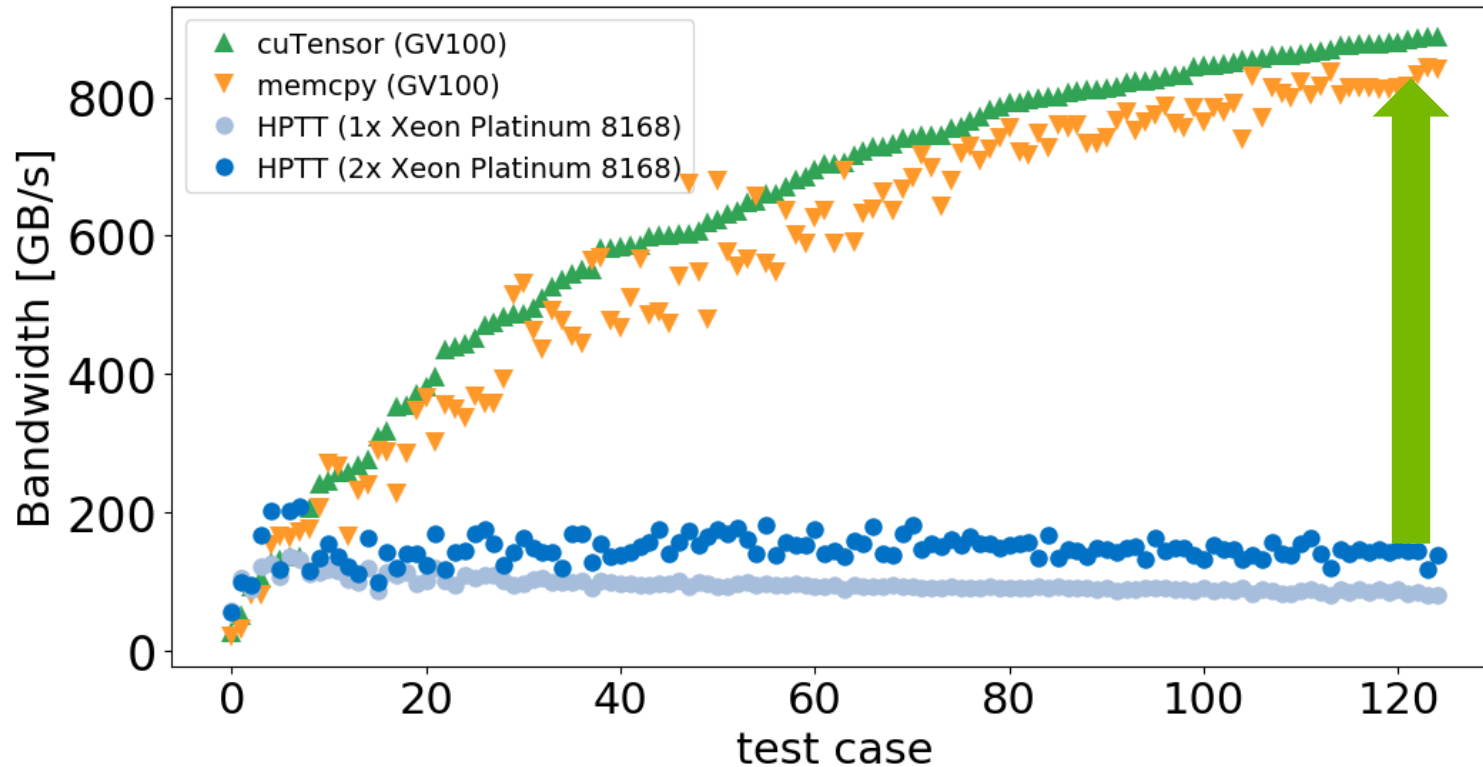
- $D_{w,h,c,n} = \min(\alpha A_{c,w,h,n}, \beta B_{c,w,h}) + \gamma C_{w,h,c,n}$

```
auto status = cutensorElementwiseTrinary ( handle,  
    alpha, A, descA, { 'c', 'w', 'h', 'n' },  
    beta,  B, descB, { 'c', 'w', 'h' },  
    gamma, C, descC, { 'w', 'h', 'c', 'n' },  
    D, descD, { 'w', 'h', 'c', 'n' },  
    CUTENSOR_OP_MIN, CUTENSOR_OP_ADD, CUTENSOR_R_MIN_16F,  
    stream );
```

# PERFORMANCE

## Element-wise Operation

$$\text{C} = \alpha \text{A} + \beta \text{B}$$



~5x over two-socket CPU