

Ladder: Enabling Efficient Low-Precision Deep Learning Computing through Hardware-aware Tensor Transformation

Lei Wang †◇, Lingxiao Ma ◇, Shijie Cao ◇, Quanlu Zhang ◇, Jilong Xue ◇, Yining Shi ‡◇, Ningxin Zheng ◇, Ziming Miao ◇, Fan Yang ◇, Ting Cao ◇, Yuqing Yang ◇, Mao Yang◇

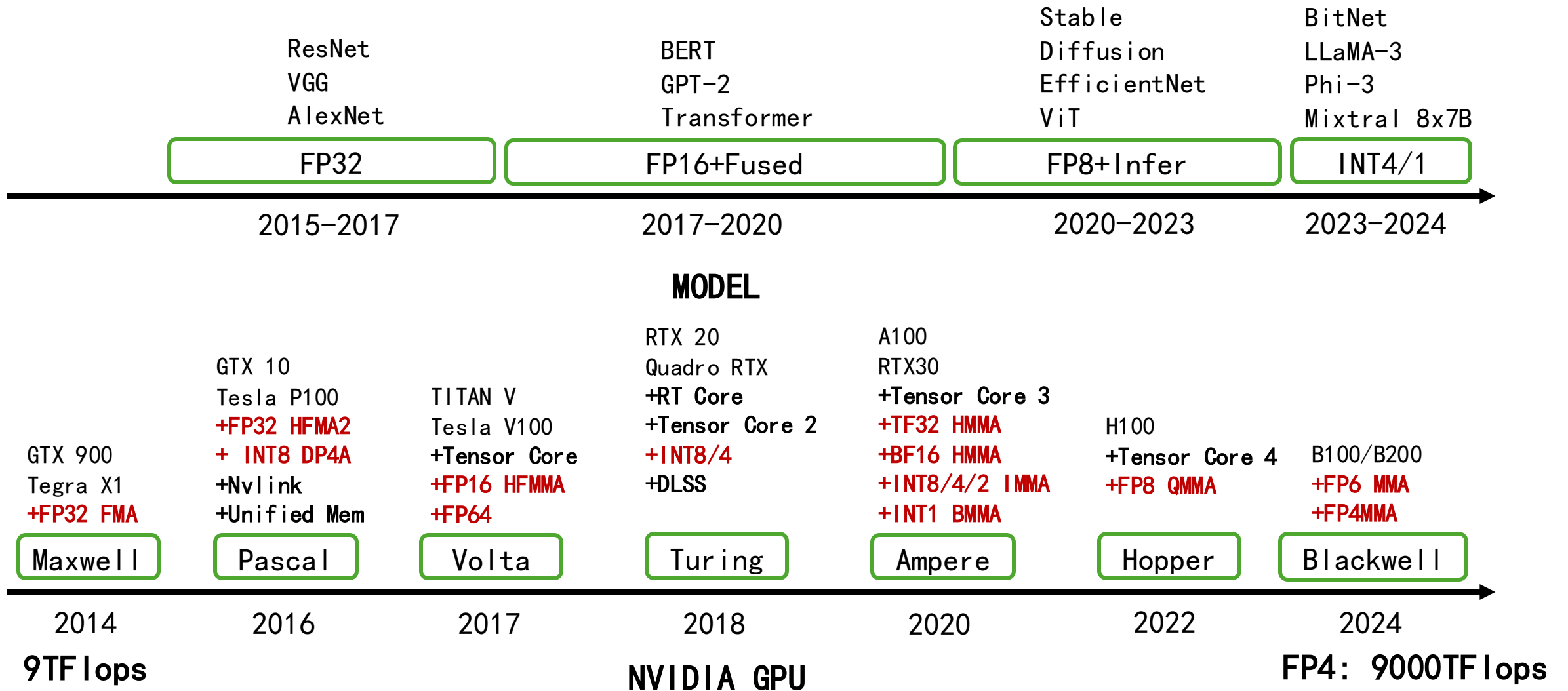
† University of Chinese Academy of Sciences

‡ Peking University

◇ Microsoft Research

Presenter: Chengru Yang

Larger Scale, Fewer Bits



Supporting Low-Precision DNN Computing

低精度训练趋势

- 使用FP16, INT8/4/2/1来训练
- 分组精度缩放
 - MXFP, 分组量化
- 混合精度运算
 - FP16 x INT4/NF4, INT8 x INT1

加速器面临困境

- 需使用代理类型: 例如, 用 FP16 模拟 FP4
- 低精度计算低效问题
 - FP16/INT8 矩阵乘法的平均利用率<60%

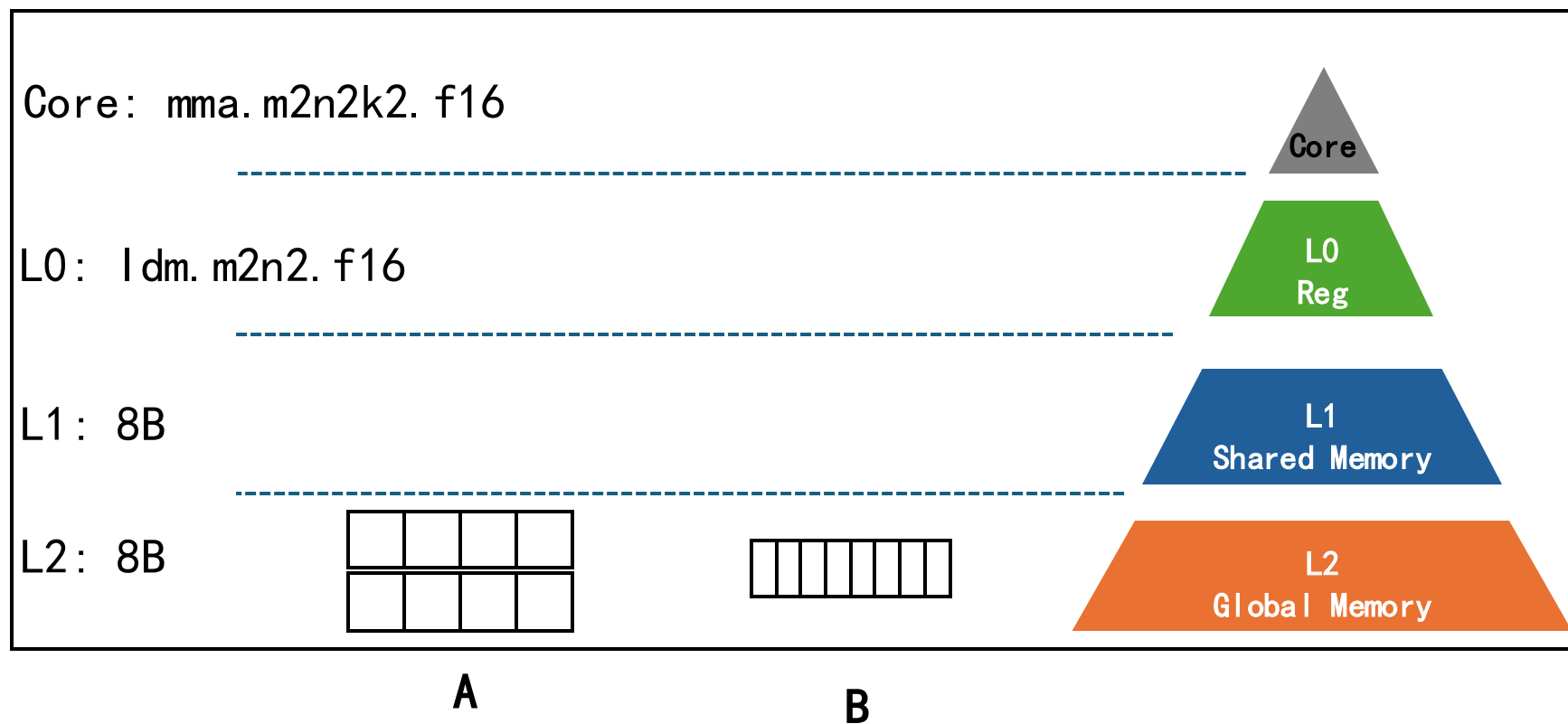
Data Type	W@FP16 / A@FP16			W@INT8 / A@INT8			W@FP8 / A@FP8	W@NF4 / A@FP16
	V100	A100	MI250	V100	A100	MI250	V100/A100/MI250	
cuBLAS	78%	87%	X	X	68%	X	X	X
rocBLAS	X	X	46%	X	X	75%	X	X
AMOS	64%	38%	X	X	45%	X	X	X
TensorIR	67%	56%	22%	X	X	X	X	X
Roller	50%	70%	29%	X	X	X	X	X

MatMul 在不同模型参数精度上和GPU上的利用率

低精度计算不断高速演化使得硬件进化难以跟上

Supporting Low-Precision DNN Computing

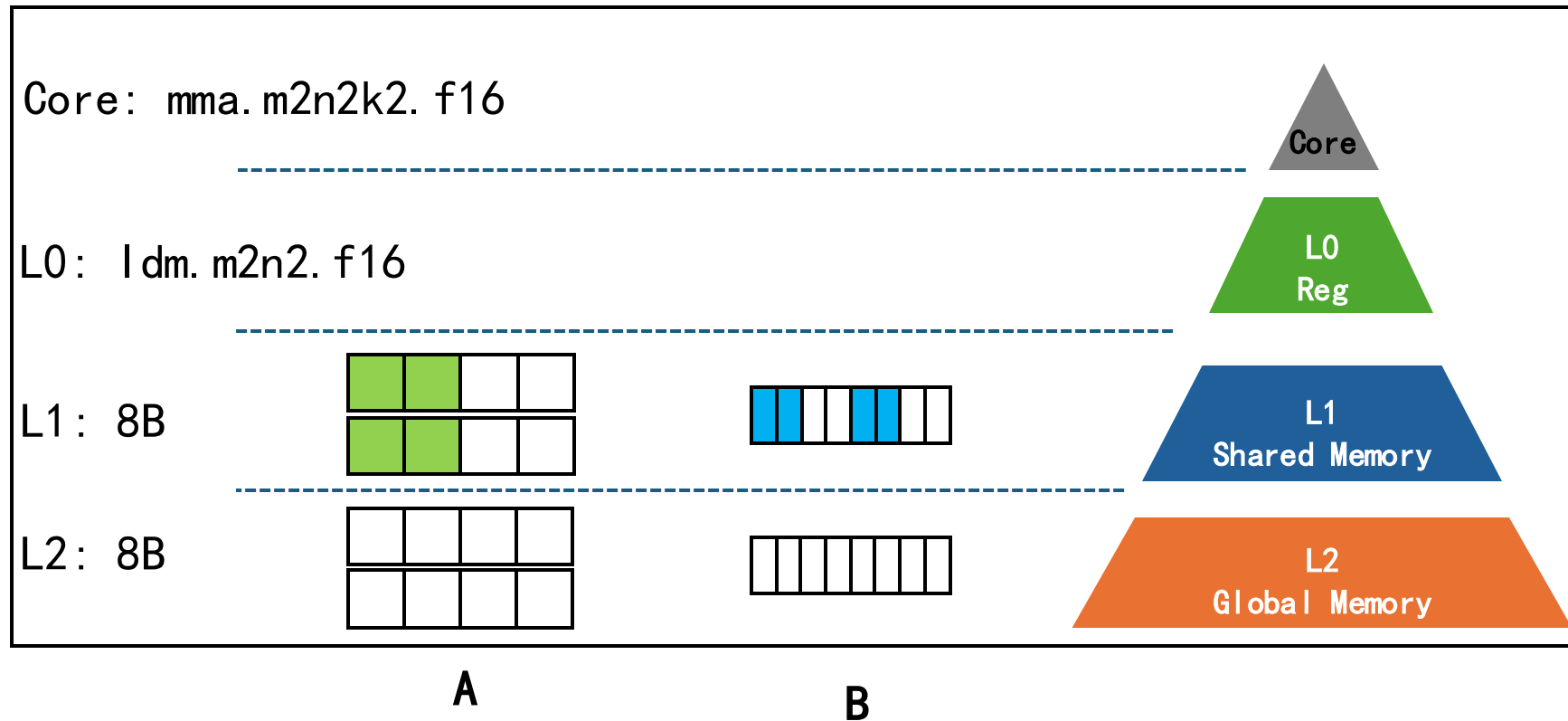
Matrix multiplication: $C[M, N]@FP16 = A[M, K]@FP16 \times B[N, K]@FP8$, $M=2$, $N=2$, $K=4$



$[X, Y]@Z$
表示：
shape是 $[X, Y]$ 的
每一个元素大小是
 Z 的一组同质元素

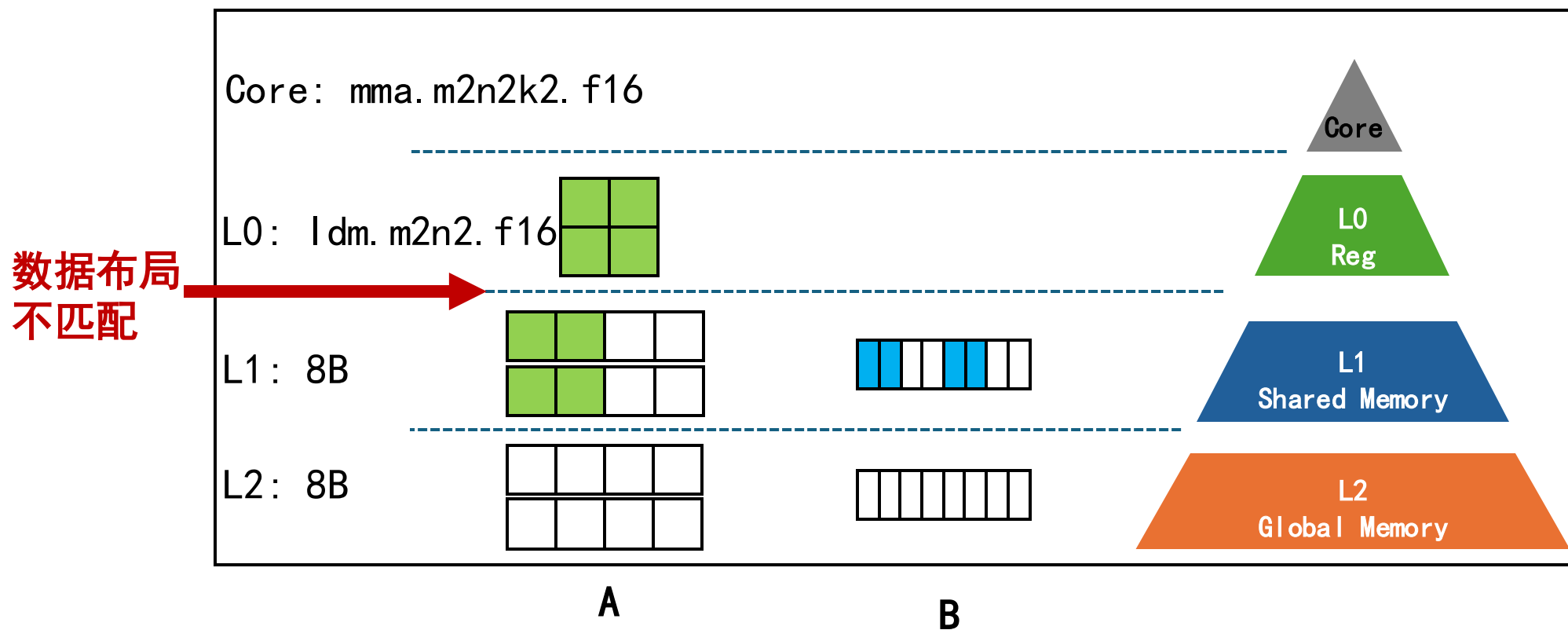
Supporting Low-Precision DNN Computing

Matrix multiplication: $C[M, N]@FP16 = A[M, K]@FP16 \times B[N, K]@FP8$, $M=2$, $N=2$, $K=4$



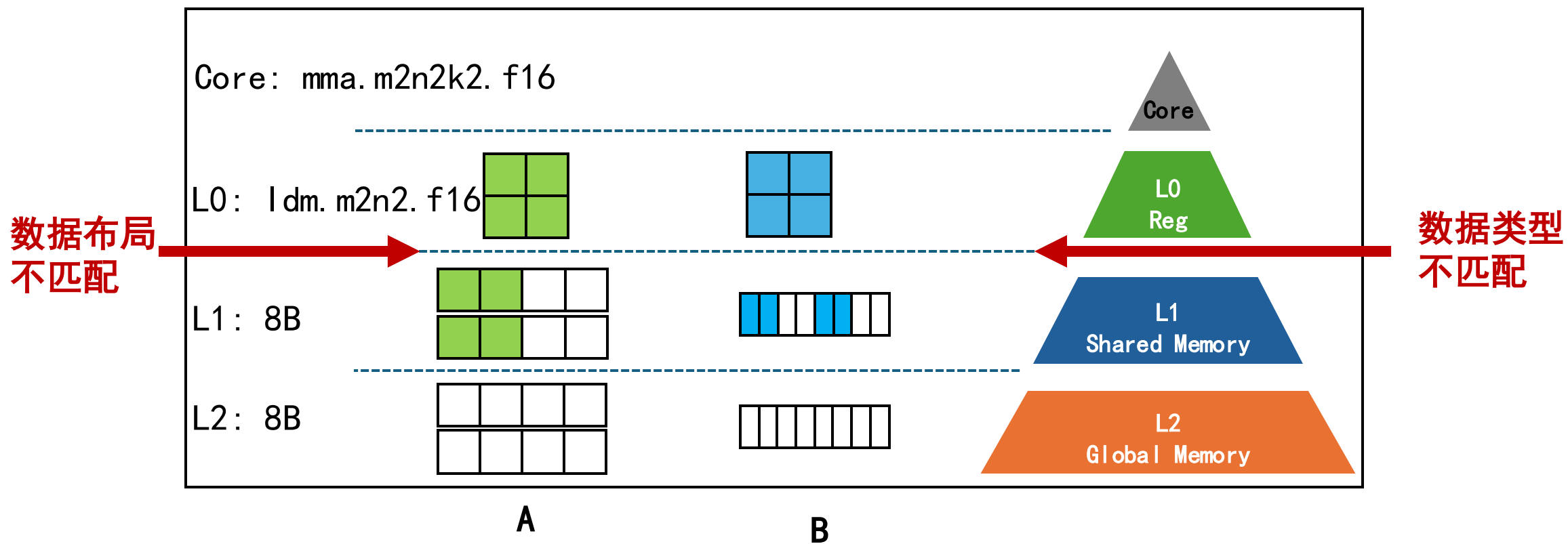
Supporting Low-Precision DNN Computing

Matrix multiplication: $C[M, N]@FP16 = A[M, K]@FP16 \times B[N, K]@FP8$, $M=2$, $N=2$, $K=4$



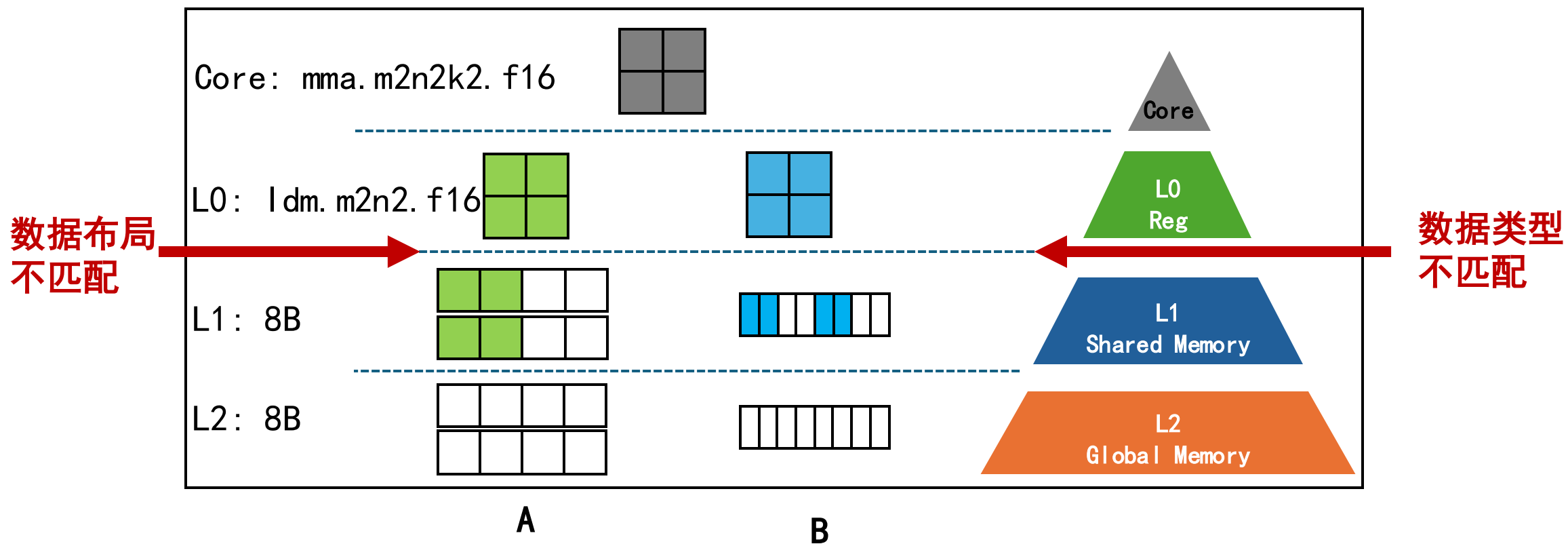
Supporting Low-Precision DNN Computing

Matrix multiplication: $C[M, N]@FP16 = A[M, K]@FP16 \times B[N, K]@FP8$, $M=2, N=2, K=4$



Supporting Low-Precision DNN Computing

Matrix multiplication: $C[M, N]@FP16 = A[M, K]@FP16 \times B[N, K]@FP8$, $M=2, N=2, K=4$



Insights

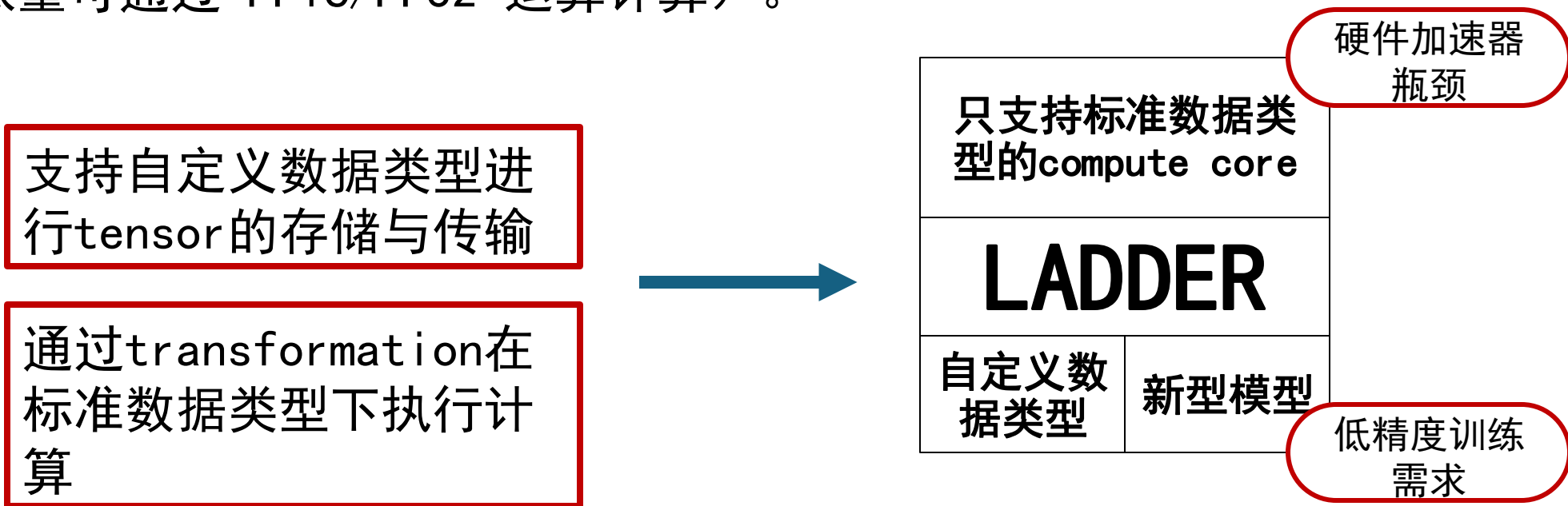
存储与计算分离：硬件加速器虽缺乏对自定义数据类型的计算指令，但其内存系统可通过将数据转换为固定位宽的不透明块来存储任意类型。

无损类型转换：多数自定义数据类型可无损转换为硬件支持的标准类型（如 NF4 张量可通过 FP16/FP32 运算计算）。

Insights

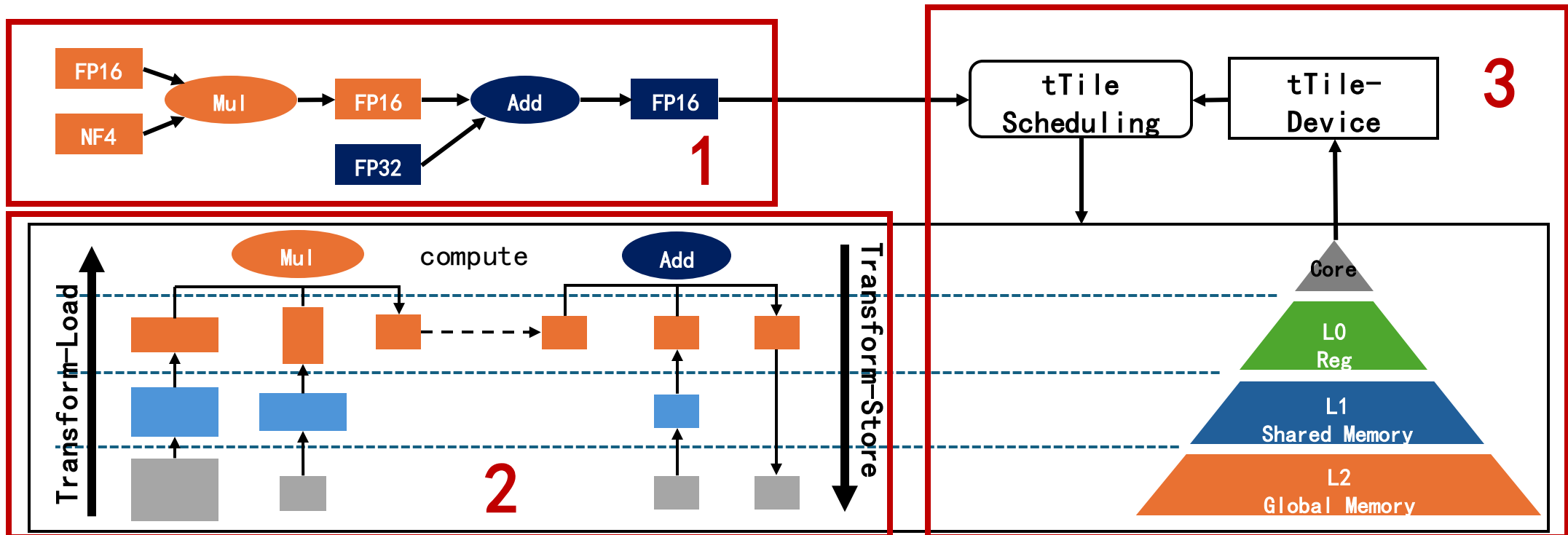
存储与计算分离：硬件加速器虽缺乏对自定义数据类型的计算指令，但其内存系统可通过将数据转换为固定位宽的不透明块来存储任意类型。

无损类型转换：多数自定义数据类型可无损转换为硬件支持的标准类型（如 NF4 张量可通过 FP16/FP32 运算计算）。



LADDER Design

- 1 tTile Abstraction
- 2 tTile Transformation
- 3 Hardware-Aware tTile-Graph Scheduling



tTile Abstraction

tType

```
Class tType {  
    TileShape shape;  
    size_t nElembits;  
    struct metadata;  
    map<TileType, c_func> c_Types;  
};
```

tType表示由一组同质元素组成的数据类型。
shape表示n维数组的形状，
nElemBits表示存储一个元素需要几位，
同时这些元素共享相同的**metadata**，
c_Types表示可以通过**c_func**函数无损的将其转换为另一个**tType**。

举例：

FP16: shape=[1] nElembits=16

NF4: shape=[1] nElembits=4 在metadata中共享value map

OCP-MXFP8: shape=[32] nElembits=8 在metadata中共享 scaling factor

tTile Abstraction

tType

```
Class tType {  
  TileShape shape;  
  size_t nElembits;  
  struct metadata;  
  map<TileType, c_func> c_Types;  
};
```

tTile

```
Class tTile {  
  TileShape shape;  
  tType type;  
  struct metadata;  
};
```

tTile是一组具有相同数据类型`type`和n维数组形状布局的同质元素。
`shape`表示n维tensor的形状，`type`表示对应的tType，**tTile**中的元素共享`metadata`。

举例：

[16, 16]@FP16: 16x16形状，每一个元素为FP16

[16, 2]@16B: 16x2形状，每一个元素大小为16B

[32]@4B: 32形状，每一个元素大小为4B

tTile Abstraction

tType

```
Class tType {  
  TileShape shape;  
  size_t nElembits;  
  struct metadata;  
  map<TileType, c_func> c_Types;  
};
```

tTile

```
Class tTile {  
  TileShape shape;  
  tType type;  
  struct metadata;  
};
```

tTile-Operator

```
Class tTile-Operator {  
  TensorExpr expr;  
  TileShape shape;  
  vector<tTile> get_input_tTiles();  
  vector<tTile> get_output_tTiles();  
  void compute();  
};
```

tTile-Operator表示对形状为 **shape** 的元素的张量计算任务。

get_input_tTiles() 和 **get_output_tTiles**() 分别返回此计算任务得输入和输出tTile,
compute() 执行张量表达式**expr**中定义得输入和输出tTile得计算。

tTile Abstraction

tType

```
Class tType {  
  TileShape shape;  
  size_t nElembits;  
  struct metadata;  
  map<TileType, c_func> c_Types;  
};
```

tTile

```
Class tTile {  
  TileShape shape;  
  tType type;  
  struct metadata;  
};
```

tTile-Operator

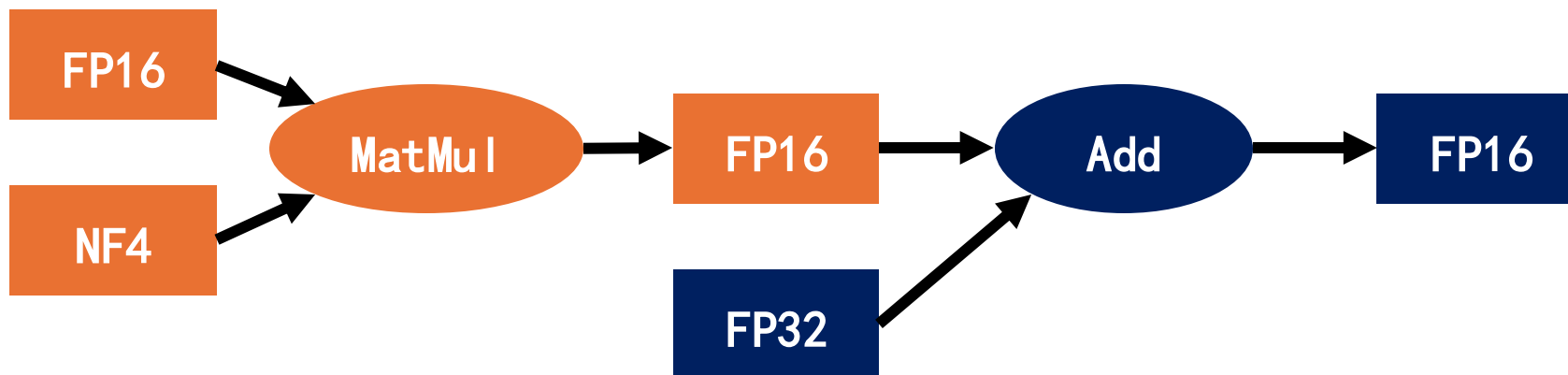
```
Class tTile-Operator {  
  TensorExpr expr;  
  TileShape shape;  
  vector<tTile> get_input_tTiles();  
  vector<tTile> get_output_tTiles();  
  void compute();  
};
```

例子:

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16)$, $M=32$, $N=32$, $K=63$
FP16 **tensor A**和NF4 **tensor B**的MatMul, 并且以FP32作为累加类型, 输出一个FP16类型的张量C[32, 32]

tTile Graph

通过基于 tTile 的细粒度表示，DNN 模型可以被表示为一个细粒度的 tTile-graph



例子:

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16)$, $M=32$, $N=32$, $K=63$
FP16 **tensor A**和NF4 **tensor B**的MatMul, 并且以FP32作为累加类型, 输出一个FP16类型的张量C[32, 32]

tTile-based Hardware Abstraction

LADDER将硬件加速器抽象为一个由多层tTile描述的结构，tTile-device。

每一层可能是内存或者计算单元：

1. 其需求通过一个特定粒度的形状来描述，表示为tTile。
2. 对应的特定粒度使用tType来表示。

下面以NVIDIA A100为例：



L2
Global Memory

Transaction: 32B
tTile: [32]@1B

Global Memory以32B的粒度进行访问

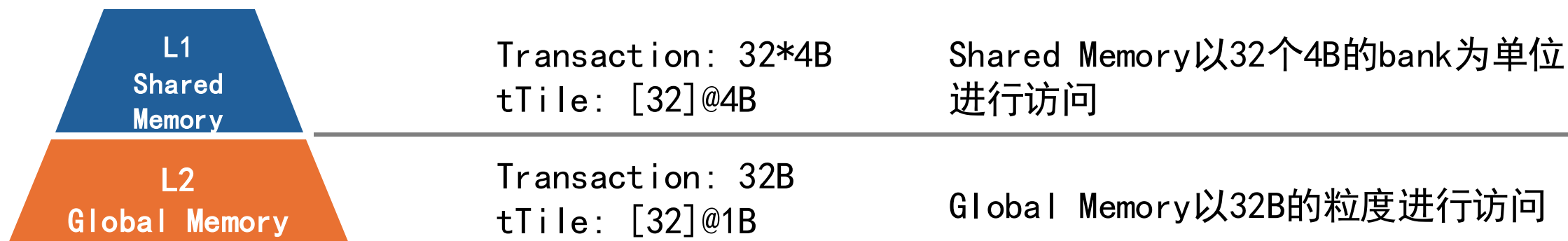
tTile-based Hardware Abstraction

LADDER将硬件加速器抽象为一个由多层tTile描述的结构，tTile-device。

每一层可能是内存或者计算单元：

1. 其需求通过一个特定粒度的形状来描述，表示为tTile。
2. 对应的特定粒度使用tType来表示。

下面以NVIDIA A100为例：



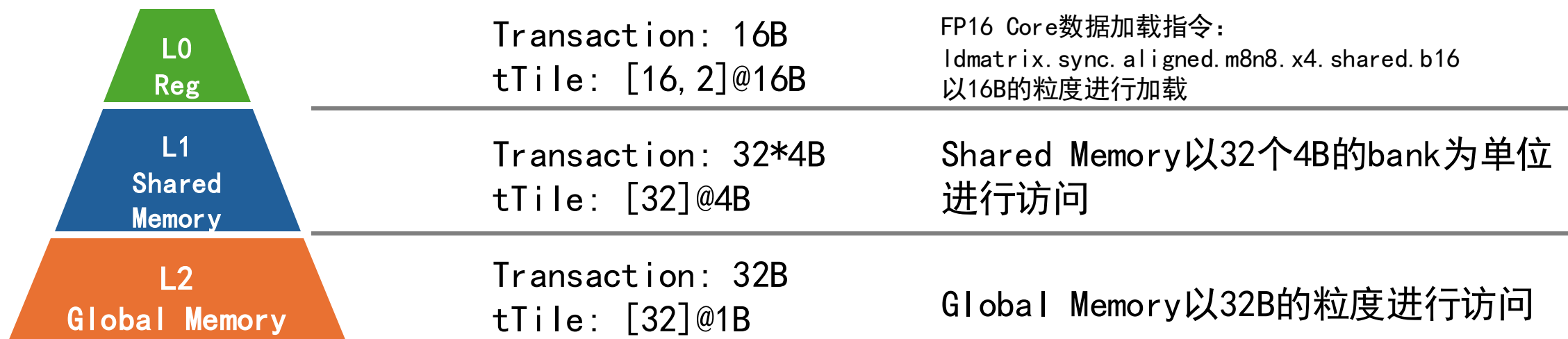
tTile-based Hardware Abstraction

LADDER将硬件加速器抽象为一个由多层tTile描述的结构，tTile-device。

每一层可能是内存或者计算单元：

1. 其需求通过一个特定粒度的形状来描述，表示为tTile。
2. 对应的特定粒度使用tType来表示。

下面以NVIDIA A100为例：



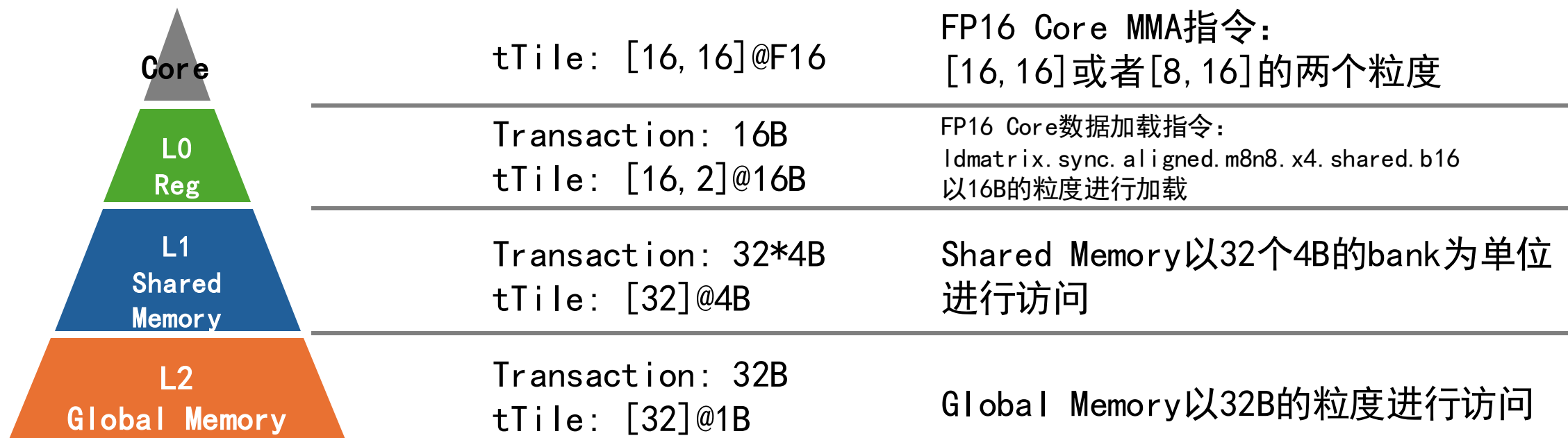
tTile-based Hardware Abstraction

LADDER将硬件加速器抽象为一个由多层tTile描述的结构，tTile-device。

每一层可能是内存或者计算单元：

1. 其需求通过一个特定粒度的形状来描述，表示为tTile。
2. 对应的特定粒度使用tType来表示。

下面以NVIDIA A100为例：



tTile Transformation

```
tTileslice(tTile_input, index, shape, out_shape);  
tTilemap(tTile_input, map_func);  
tTilepad(tTile_input, pad_shape, pad_value);  
tTileconvert(tTile_input, new_tType);
```

- **Slice**

Slice 原语从 tTile_input 的地址索引中切出一组形状为 shape 的元素，并将其作为新的 tTile 返回，新 tTile 的形状为 out_shape。Slice 原语通常用于表示数据分块 (tiling)。

- **Map**

Map 原语修改 tTile 中元素的布局。

- **Pad**

Pad 原语在 tTile_input 的每个边界上填充 pad_value，填充的形状由 pad_shape 指定。

- **Convert**

Convert 原语将 tTile_input 的 tType 转换为给定的 new_tType。

Hardware-Aware tTile-Graph Scheduling

Algorithm 1: Hint-based layer-wise scheduling

输入：表示为 tTile-graph 的 DNN 模型 g 和 tTile device

```
Data:  $g$ : tTile-graph;  $D$ : tTile-device
Result:  $g_{ret}$ : scheduled tTile-graph
1 Function GetDeviceHint( $g, D$ ):
2    $D = \text{SelectDeviceConfig}(g, D)$ ;
3    $\text{HintShape} = \text{None}$ ,  $\text{HintGranularity} = \text{None}$ ;
4   for  $layer \in D.layers$  do
5      $\text{HintGranularity} = \text{LCM}(\text{HintGranularity}, layer.tTile.type)$ ;
6   for  $layer \in D.layers$  do
7      $layer.tTile = \text{convert}(layer.tTile, \text{HintGranularity})$ ;
8      $\text{HintShape} = \text{LCM}(\text{HintShape}, layer.tTile.shape)$ ;
9   for  $layer \in D.layers$  do
10     $layer.tTile.shape = \text{HintShape}$ ;
11  return  $D$ ;
12 Function ScheduleTransform( $op, D, l_{id}$ ):
13   $tTile_h = op.tTile[l_{id}-1]$ ;
14   $tTile_l = op.tTile[l_{id}]$ ;
15   $\text{ScheduleSlice}(tTile_l, tTile_h)$ ;
16  if  $\text{LCM}(tTile_l.shape, tTile_h.shape) \neq tTile_l.shape$  then
17     $\text{SchedulePad}(tTile_l, tTile_h, D)$ ;
18  if  $tTile_l.type \neq tTile_h.type$  then
19     $\text{ScheduleConvert}(tTile_l, tTile_h, D)$ ;
20  if  $nBits(tTile_h.shape[-1]) \neq nBits(D.layers[l_{id}].shape[-1])$  then
21     $\text{ScheduleMap}(tTile_l, tTile_h, D)$ ;
22  return  $op.transform[l_{id}-1]$ ;
23 Function ScheduleConnectedGraph( $g, D$ ):
24   $D = \text{GetDeviceHint}(g, D)$ ;
25  for  $l_{id}$  in  $\text{length}(D.layers)$  do
26    for  $op \in g[l_{id}]$  do
27       $op.tTile[l_{id}] = \text{ScheduleTiling}(op, D, l_{id})$ ;
28      if  $l_{id} > 0$  then
29         $op.transform[l_{id}] = \text{ScheduleTransform}(op, D, l_{id})$ ;
30   $g = \text{ProfileAndSelect}(g)$ ;
31  return  $g$ ;
32 Function Schedule( $g, D$ ):
33   $g = \text{ExtractConnectedGraph}(g, D)$ ;
34  for  $g_{conn} \in g$  do
35     $g_{conn} = \text{ScheduleConnectedGraph}(g_{conn}, D)$ ;
36  return  $g$ ;
```

Hardware-Aware tTile-Graph Scheduling

Algorithm 1: Hint-based layer-wise scheduling

输入：表示为 tTile-graph 的 DNN 模型 g 和 tTile device

```
Data:  $g$ : tTile-graph;  $D$ : tTile-device  
Result:  $g_{ret}$ : scheduled tTile-graph  
1 Function GetDeviceHint( $g, D$ ):  
2    $D = \text{SelectDeviceConfig}(g, D)$ ;  
3    $\text{HintShape} = \text{None}$ ,  $\text{HintGranularity} = \text{None}$ ;  
4   for  $layer \in D.layers$  do  
5      $\text{HintGranularity} = \text{LCM}(\text{HintGranularity}, layer.tTile.type)$ ;  
6   for  $layer \in D.layers$  do  
7      $layer.tTile = \text{convert}(layer.tTile, \text{HintGranularity})$ ;  
8      $\text{HintShape} = \text{LCM}(\text{HintShape}, layer.tTile.shape)$ ;  
9   for  $layer \in D.layers$  do  
10     $layer.tTile.shape = \text{HintShape}$ ;  
11  return  $D$ ;  
12 Function ScheduleTransform( $op, D, l_{id}$ ):  
13   $tTile_h = op.tTile[l_{id}-1]$ ;  
14   $tTile_l = op.tTile[l_{id}]$ ;  
15   $\text{ScheduleSlice}(tTile_l, tTile_h)$ ;  
16  if  $\text{LCM}(tTile_l.shape, tTile_h.shape) \neq tTile_l.shape$  then  
17     $\text{SchedulePad}(tTile_l, tTile_h, D)$ ;  
18  if  $tTile_l.type \neq tTile_h.type$  then  
19     $\text{ScheduleConvert}(tTile_l, tTile_h, D)$ ;  
20  if  $nBits(tTile_h.shape[-1]) \neq nBits(D.layers[l_{id}].shape[-1])$  then  
21     $\text{ScheduleMap}(tTile_l, tTile_h, D)$ ;  
22  return  $op.transform[l_{id}-1]$ ;  
23 Function ScheduleConnectedGraph( $g, D$ ):  
24   $D = \text{GetDeviceHint}(g, D)$ ;  
25  for  $l_{id}$  in  $\text{length}(D.layers)$  do  
26    for  $op \in g[l_{id}]$  do  
27       $op.tTile[l_{id}] = \text{ScheduleTiling}(op, D, l_{id})$ ;  
28      if  $l_{id} > 0$  then  
29         $op.transform[l_{id}] = \text{ScheduleTransform}(op, D, l_{id})$ ;  
30   $g = \text{ProfileAndSelect}(g)$ ;  
31  return  $g$ ;  
32 Function Schedule( $g, D$ ):  
33   $g = \text{ExtractConnectedGraph}(g, D)$ ;  
34  for  $g_{conn} \in g$  do  
35     $g_{conn} = \text{ScheduleConnectedGraph}(g_{conn}, D)$ ;  
36  return  $g$ ;
```

将 g 分为 g_{conn} 子图，每个子图从最底层内存加载 tTile 到 Core，再将结果存回最低层。可能是一个 tTile-OP 或一组可融合的 tTile-OP

Hardware-Aware tTile-Graph Scheduling

Algorithm 1: Hint-based layer-wise scheduling

Data: g : tTile-graph; D : tTile-device

Result: g_{ret} : scheduled tTile-graph

```
1 Function GetDeviceHint( $g, D$ ):  
2    $D = \text{SelectDeviceConfig}(g, D)$ ;  
3    $\text{HintShape} = \text{None}$ ,  $\text{HintGranularity} = \text{None}$ ;  
4   for  $layer \in D.layers$  do  
5      $\text{HintGranularity} = \text{LCM}(\text{HintGranularity}, layer.tTile.type)$ ;  
6   for  $layer \in D.layers$  do  
7      $layer.tTile = \text{convert}(layer.tTile, \text{HintGranularity})$ ;  
8      $\text{HintShape} = \text{LCM}(\text{HintShape}, layer.tTile.shape)$ ;  
9   for  $layer \in D.layers$  do  
10     $layer.tTile.shape = \text{HintShape}$ ;  
11  return  $D$ ;
```

```
12 Function ScheduleTransform( $op, D, l_{id}$ ):  
13    $tTile_h = op.tTile[l_{id}-1]$ ;  
14    $tTile_l = op.tTile[l_{id}]$ ;  
15    $\text{ScheduleSlice}(tTile_l, tTile_h)$ ;  
16   if  $\text{LCM}(tTile_l.shape, tTile_h.shape) \neq tTile_l.shape$  then  
17      $\text{SchedulePad}(tTile_l, tTile_h, D)$ ;  
18   if  $tTile_l.type \neq tTile_h.type$  then  
19      $\text{ScheduleConvert}(tTile_l, tTile_h, D)$ ;  
20   if  $nBits(tTile_h.shape[-1]) \neq nBits(D.layers[l_{id}].shape[-1])$  then  
21      $\text{ScheduleMap}(tTile_l, tTile_h, D)$ ;  
22   return  $op.transform[l_{id}-1]$ ;
```

```
23 Function ScheduleConnectedGraph( $g, D$ ):  
24    $D = \text{GetDeviceHint}(g, D)$ ;  
25   for  $l_{id}$  in  $\text{length}(D.layers)$  do  
26     for  $op \in g[l_{id}]$  do  
27        $op.tTile[l_{id}] = \text{ScheduleTiling}(op, D, l_{id})$ ;  
28       if  $l_{id} > 0$  then  
29          $op.transform[l_{id}] = \text{ScheduleTransform}(op, D, l_{id})$ ;  
30    $g = \text{ProfileAndSelect}(g)$ ;  
31  return  $g$ ;
```

```
32 Function Schedule( $g, D$ ):  
33    $g = \text{ExtractConnectedGraph}(g, D)$ ;  
34   for  $g_{conn} \in g$  do  
35      $g_{conn} = \text{ScheduleConnectedGraph}(g_{conn}, D)$ ;  
36  return  $g$ ;
```

输入：表示为tTile-graph的DNN模型 g 和tTile device

计算Device Hint:

选择Core (比如优先FP16)

通过最小公倍数 (LCM) 来计算HintGranularity

确保跨层数据对齐, 并且通过LCM来计算HintShape

将 g 分为 g_{conn} 子图, 每个子图从最底层内存加载tTile到Core, 再将结果存回最低层。可能是一个tTile-OP或一组可融合的tTile-OP

Hardware-Aware tTile-Graph Scheduling

Algorithm 1: Hint-based layer-wise scheduling

Data: g : tTile-graph; D : tTile-device

Result: g_{ret} : scheduled tTile-graph

```
1 Function GetDeviceHint( $g, D$ ):
2    $D = \text{SelectDeviceConfig}(g, D)$ ;
3    $\text{HintShape} = \text{None}$ ,  $\text{HintGranularity} = \text{None}$ ;
4   for  $layer \in D.layers$  do
5      $\text{HintGranularity} = \text{LCM}(\text{HintGranularity}, layer.tTile.type)$ ;
6   for  $layer \in D.layers$  do
7      $layer.tTile = \text{convert}(layer.tTile, \text{HintGranularity})$ ;
8      $\text{HintShape} = \text{LCM}(\text{HintShape}, layer.tTile.shape)$ ;
9   for  $layer \in D.layers$  do
10     $layer.tTile.shape = \text{HintShape}$ ;
11  return  $D$ ;
```

```
12 Function ScheduleTransform( $op, D, l_{id}$ ):
13    $tTile_h = op.tTile[l_{id}-1]$ ;
14    $tTile_l = op.tTile[l_{id}]$ ;
15    $\text{ScheduleSlice}(tTile_l, tTile_h)$ ;
16   if  $\text{LCM}(tTile_l.shape, tTile_h.shape) \neq tTile_l.shape$  then
17      $\text{SchedulePad}(tTile_l, tTile_h, D)$ ;
18   if  $tTile_l.type \neq tTile_h.type$  then
19      $\text{ScheduleConvert}(tTile_l, tTile_h, D)$ ;
20   if  $nBits(tTile_h.shape[-1]) \neq nBits(D.layers[l_{id}].shape[-1])$  then
21      $\text{ScheduleMap}(tTile_l, tTile_h, D)$ ;
22  return  $op.transform[l_{id}-1]$ ;
```

```
23 Function ScheduleConnectedGraph( $g, D$ ):
24    $D = \text{GetDeviceHint}(g, D)$ ;
25   for  $l_{id}$  in  $\text{length}(D.layers)$  do
26     for  $op \in g[l_{id}]$  do
27        $op.tTile[l_{id}] = \text{ScheduleTiling}(op, D, l_{id})$ ;
28       if  $l_{id} > 0$  then
29          $op.transform[l_{id}] = \text{ScheduleTransform}(op, D, l_{id})$ ;
30    $g = \text{ProfileAndSelect}(g)$ ;
31  return  $g$ ;
```

```
32 Function Schedule( $g, D$ ):
33    $g = \text{ExtractConnectedGraph}(g, D)$ ;
34   for  $g_{conn} \in g$  do
35      $g_{conn} = \text{ScheduleConnectedGraph}(g_{conn}, D)$ ;
36  return  $g$ ;
```

输入：表示为tTile-graph的DNN模型 g 和tTile device

计算Device Hint:

选择Core (比如优先FP16)

通过最小公倍数 (LCM) 来计算HintGranularity

确保跨层数据对齐, 并且通过LCM来计算HintShape

输入算子 op , Device hint D , 当前层 l_{id}

输出当前层调度的tTile Transformation

将 g 分为 g_{conn} 子图, 每个子图从最底层内存加载tTile到Core, 再将结果存回最低层。可能是一个tTile-OP或一组可融合的tTile-OP

Hardware-Aware tTile-Graph Scheduling

Algorithm 1: Hint-based layer-wise scheduling

Data: g : tTile-graph; D : tTile-device

Result: g_{ret} : scheduled tTile-graph

```
1 Function GetDeviceHint( $g, D$ ):
2    $D = \text{SelectDeviceConfig}(g, D)$ ;
3    $\text{HintShape} = \text{None}$ ,  $\text{HintGranularity} = \text{None}$ ;
4   for  $layer \in D.layers$  do
5      $\text{HintGranularity} = \text{LCM}(\text{HintGranularity}, layer.tTile.type)$ ;
6   for  $layer \in D.layers$  do
7      $layer.tTile = \text{convert}(layer.tTile, \text{HintGranularity})$ ;
8      $\text{HintShape} = \text{LCM}(\text{HintShape}, layer.tTile.shape)$ ;
9   for  $layer \in D.layers$  do
10     $layer.tTile.shape = \text{HintShape}$ ;
11  return  $D$ ;
```

```
12 Function ScheduleTransform( $op, D, l_{id}$ ):
13    $tTile_h = op.tTile[l_{id}-1]$ ;
14    $tTile_l = op.tTile[l_{id}]$ ;
15    $\text{ScheduleSlice}(tTile_l, tTile_h)$ ;
16   if  $\text{LCM}(tTile_l.shape, tTile_h.shape) \neq tTile_l.shape$  then
17      $\text{SchedulePad}(tTile_l, tTile_h, D)$ ;
18   if  $tTile_l.type \neq tTile_h.type$  then
19      $\text{ScheduleConvert}(tTile_l, tTile_h, D)$ ;
20   if  $nBits(tTile_h.shape[-1]) \neq nBits(D.layers[l_{id}].shape[-1])$  then
21      $\text{ScheduleMap}(tTile_l, tTile_h, D)$ ;
22   return  $op.transform[l_{id}-1]$ ;
```

```
23 Function ScheduleConnectedGraph( $g, D$ ):
24    $D = \text{GetDeviceHint}(g, D)$ ;
25   for  $l_{id}$  in  $\text{length}(D.layers)$  do
26     for  $op \in g[l_{id}]$  do
27        $op.tTile[l_{id}] = \text{ScheduleTiling}(op, D, l_{id})$ ;
28       if  $l_{id} > 0$  then
29          $op.transform[l_{id}] = \text{ScheduleTransform}(op, D, l_{id})$ ;
30    $g = \text{ProfileAndSelect}(g)$ ;
31   return  $g$ ;
```

```
32 Function Schedule( $g, D$ ):
33    $g = \text{ExtractConnectedGraph}(g, D)$ ;
34   for  $g_{conn} \in g$  do
35      $g_{conn} = \text{ScheduleConnectedGraph}(g_{conn}, D)$ ;
36   return  $g$ ;
```

输入：表示为tTile-graph的DNN模型 g 和tTile device

计算Device Hint:

选择Core (比如优先FP16)

通过最小公倍数 (LCM) 来计算HintGranularity

确保跨层数据对齐, 并且通过LCM来计算HintShape

输入算子 op , Device hint D , 当前层 l_{id}

输出当前层调度的tTile Transformation

拼接子图

完善层间调度

性能分析

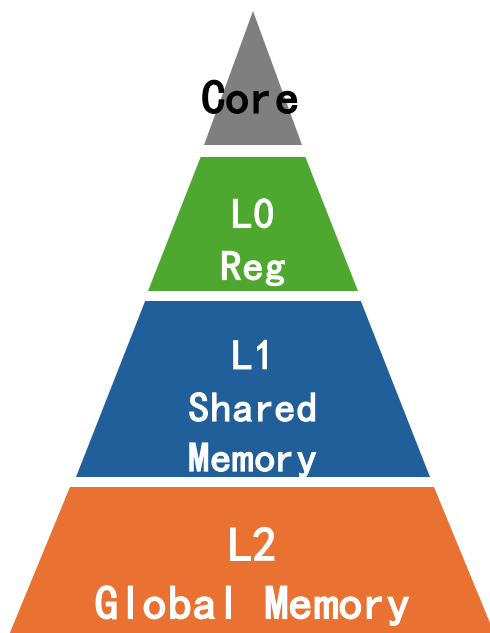
LADDER 提供了一种以 延迟为导向的策略, 目标是最小化端到端延迟

将 g 分为 g_{conn} 子图, 每个子图从最底层内存加载tTile到Core, 再将结果存回最低层。可能是一个tTile-OP或一组可融合的tTile-OP

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$



tTile: $[16, 16]@F16$

Transaction: **16B**

tTile: $[16, 2]@16B$

Transaction: $32 * 4B$

tTile: $[32]@4B$

Transaction: $32B$

tTile: $[32]@1B$

HintGranularity: Idmatrix使用的**16B**

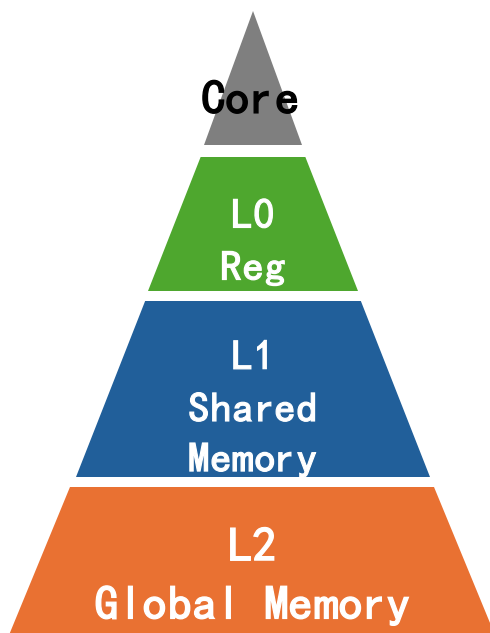
tTile-device for NVIDIA A100

计算Device Hint

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}(A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16)$,
 $M=32, N=32, K=63$



tTile-device for NVIDIA A100

tTile: $[16, 16]@F16$

Transaction: 16B
tTile: $[16, 2]@16B$

Transaction: $32*4B$
tTile: $[32]@4B$

Transaction: $32B$
tTile: $[32]@1B$

HintGranularity: Idmatrix使用的 $16B$

HintShape:
计算目标是使得L1和L2的Transaction都能对齐。

$LCM(128, 32) = 128B$

$\text{innerdim} = 128B / 16B = 8$

$\text{outerdim} = 16*16*2B(FP16) / 128B = 4$

最终算子分块为 $[4, 8]$ 的倍数且粒度为 $16B$

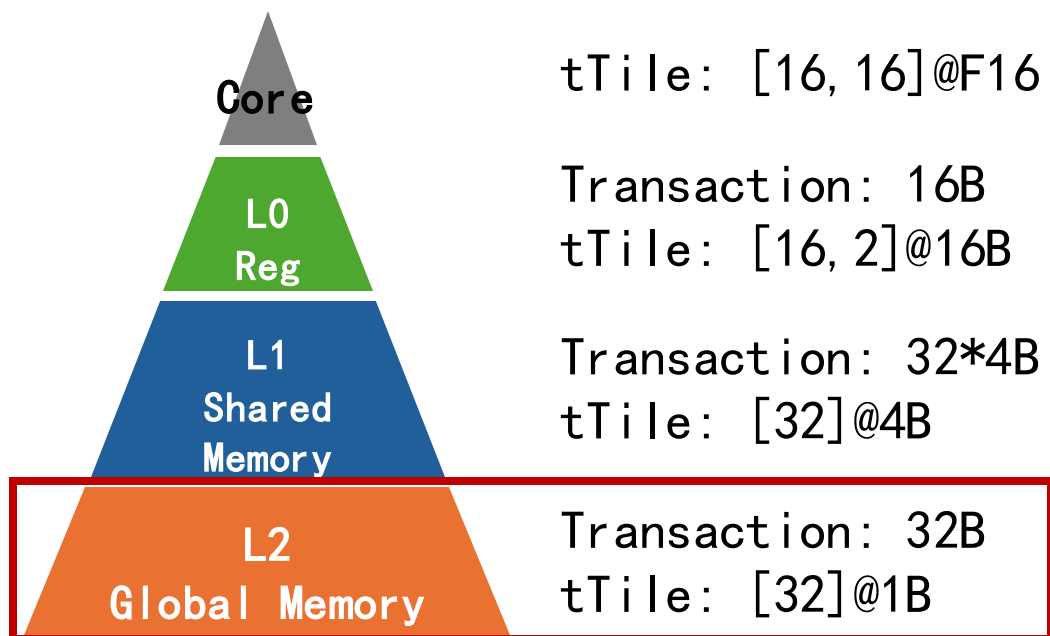
计算Device Hint

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}(A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16)$,
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)



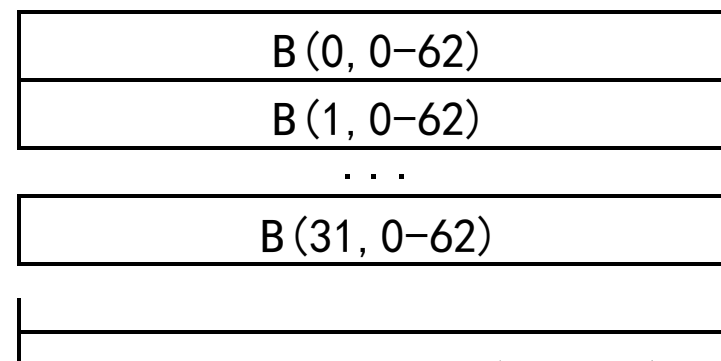
tTile-device for NVIDIA A100

Tensor大小: $63*4\text{bit}*32 = 1008B$

Hint size = 512B

sliceB: $[16, 63]@NF4$

padB: $[16, 64]@NF4$



$63*4\text{bit} = 252B$ 并非对齐

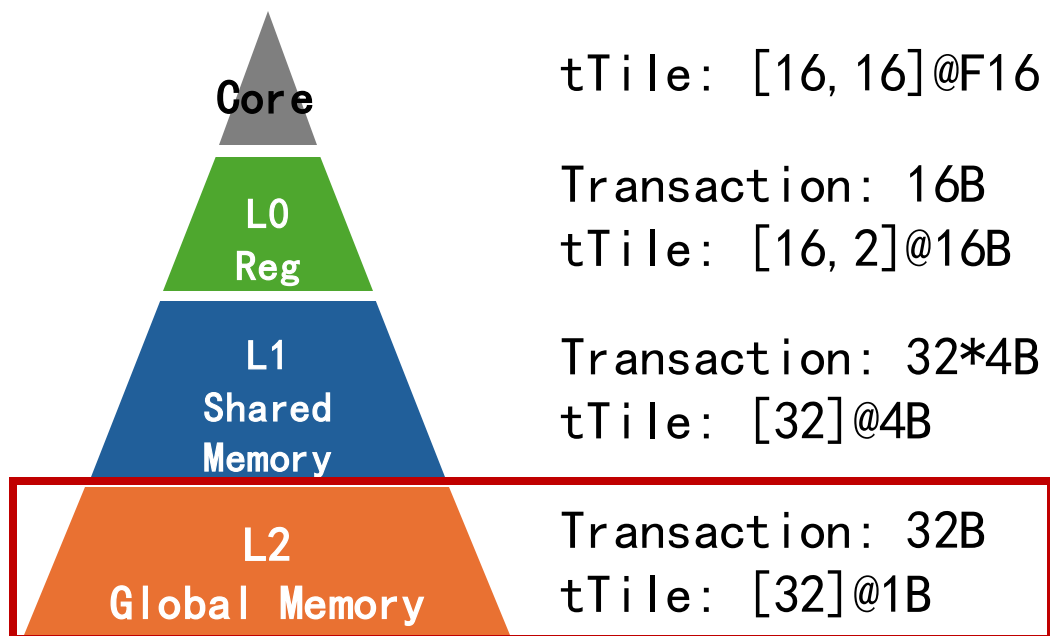
$B[32, 63]@NF4$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}(A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16)$,
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)



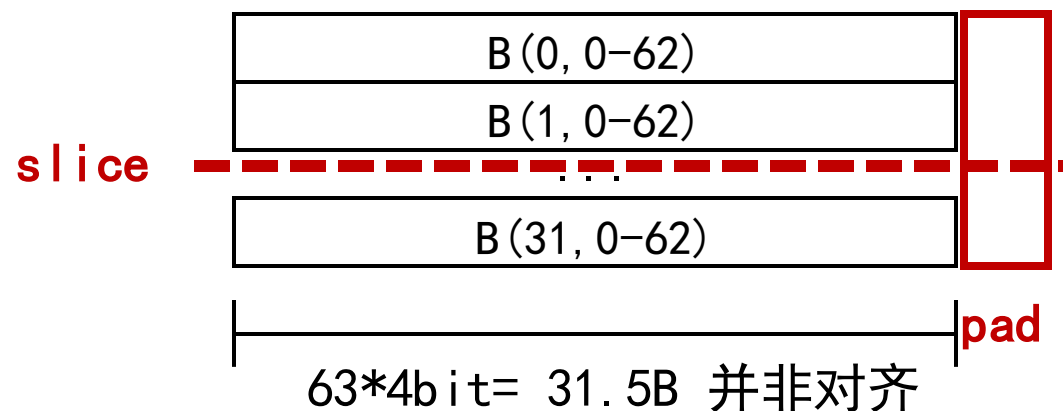
tTile-device for NVIDIA A100

Tensor大小: $63*4\text{bit}*32 = 1008B$

Hint size = 512B

sliceB: $[16, 63]@NF4$

padB: $[16, 64]@NF4$



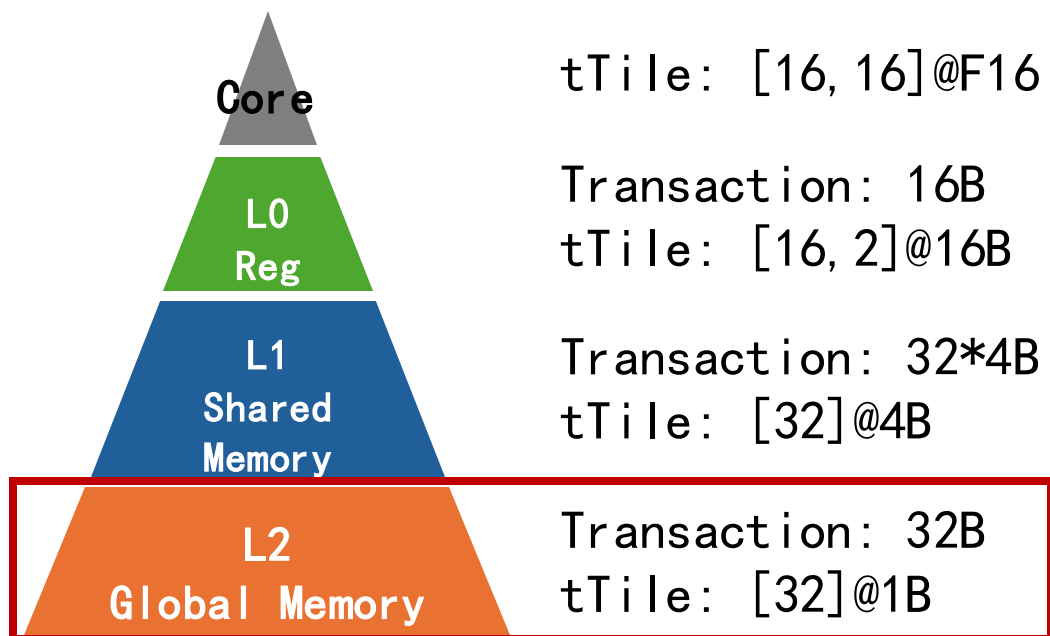
$B[32, 63]@NF4$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)



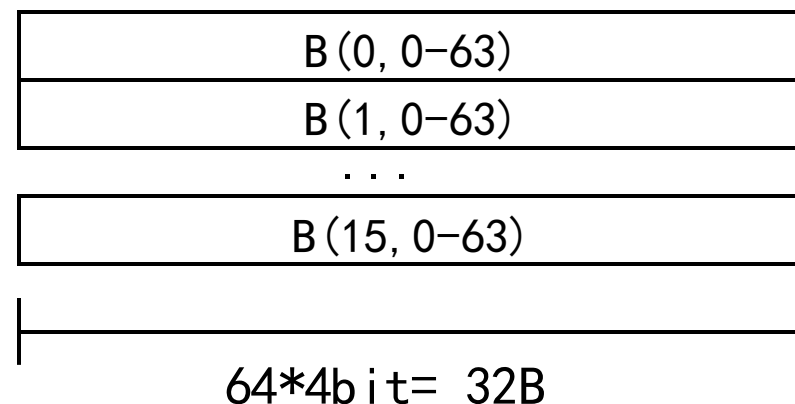
tTile-device for NVIDIA A100

Tensor大小: $64*4\text{bit}*16 = 512B$

Hint size = 512B

$NF4 \neq FP16$

convertB: $[16, 64]@FP16$



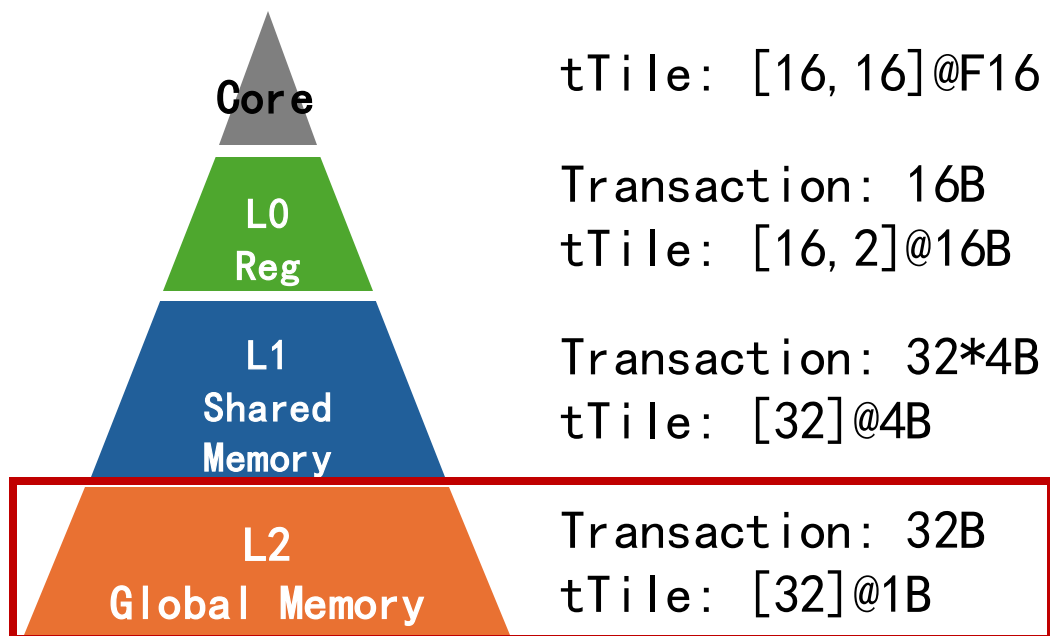
$B[16, 64]@NF4$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

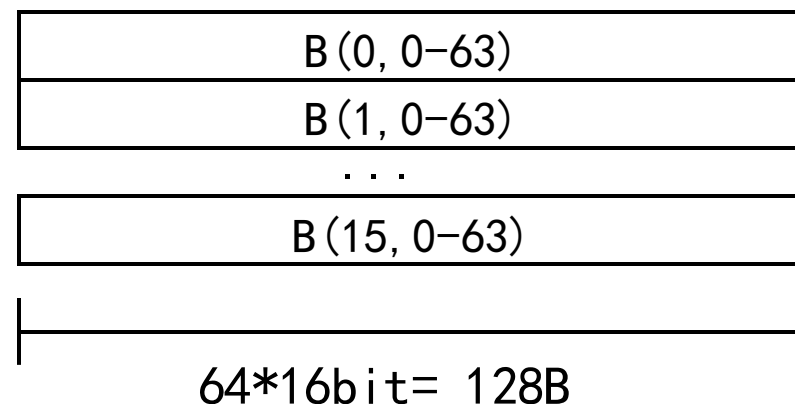
Device Hint: $[4, 8]@16B$ (size: 512B)



tTile-device for NVIDIA A100

Tensor大小: $64*16\text{bit}*16 = 2048B$

Hint size = 512B



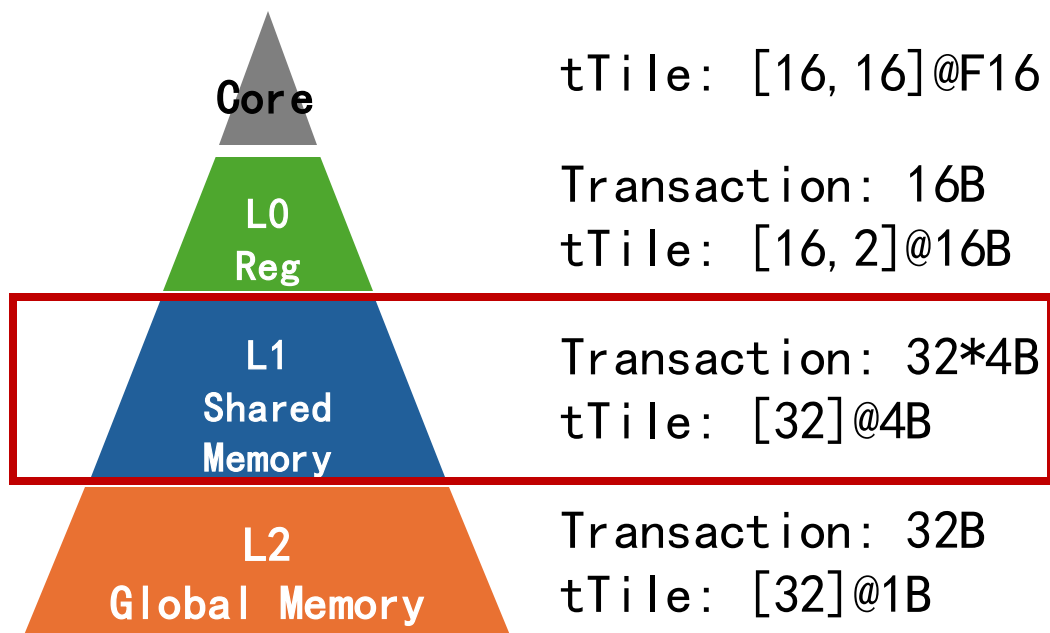
$B[16, 64]@FP16$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

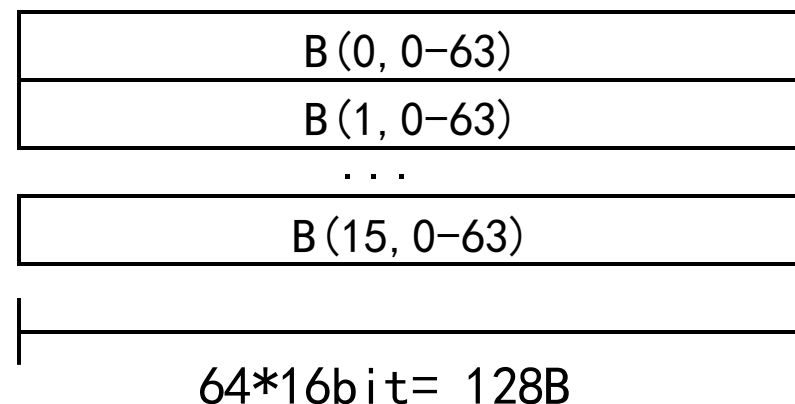
$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)



tTile-device for NVIDIA A100

Tensor大小: $64*16\text{bit}*16 = 2048B$
Hint size = 512B



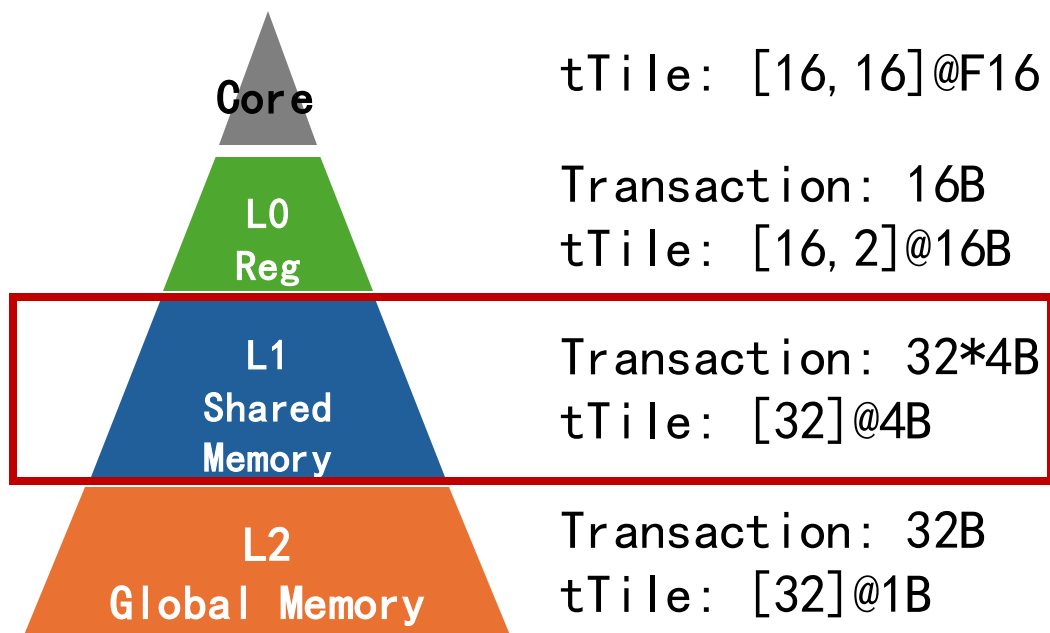
$B[16, 64]@FP16$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)

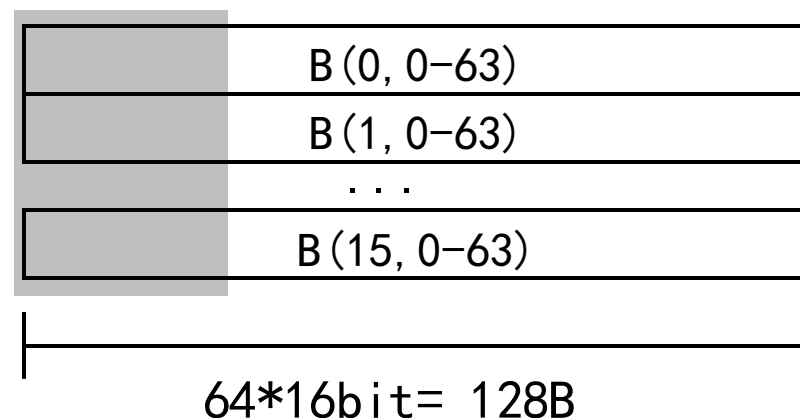


tTile-device for NVIDIA A100

Tensor大小: $64*16\text{bit}*16 = 2048B$

Hint size = 512B

mapB: $[16, 64]@FP16$ (使得切分后能对齐)



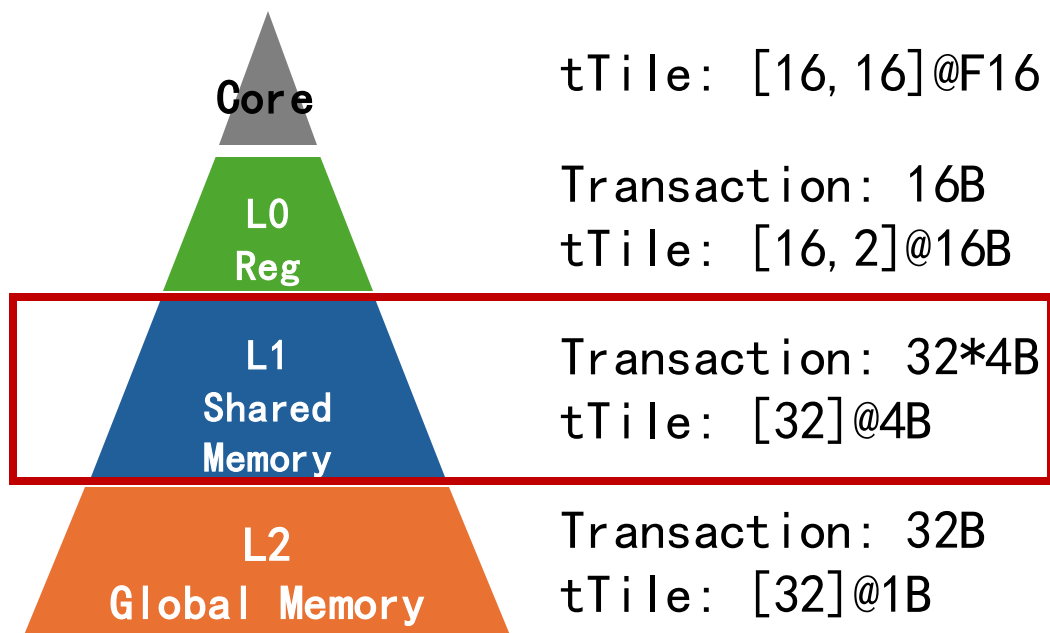
$B[16, 64]@FP16$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

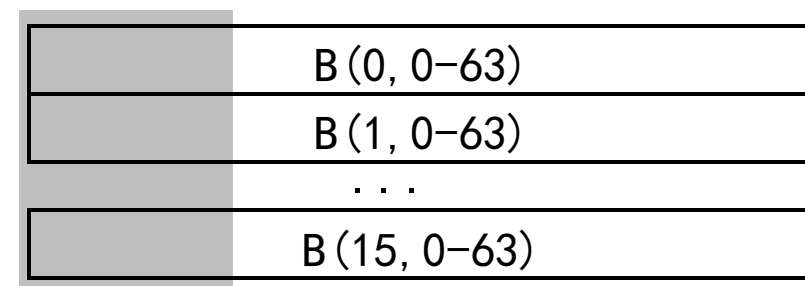
$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)



tTile-device for NVIDIA A100

B(0, 0-15)	B(1, 0-15)	B(2, 0-15)	B(3, 0-15)
B(4, 0-15)	B(5, 0-15)	B(6, 0-15)	B(7, 0-15)
B(8, 0-15)	B(9, 0-15)	B(10, 0-15)	B(11, 0-15)
B(12, 0-15)	B(13, 0-15)	B(14, 0-15)	B(15, 0-15)



$64*16\text{bit} = 128B$

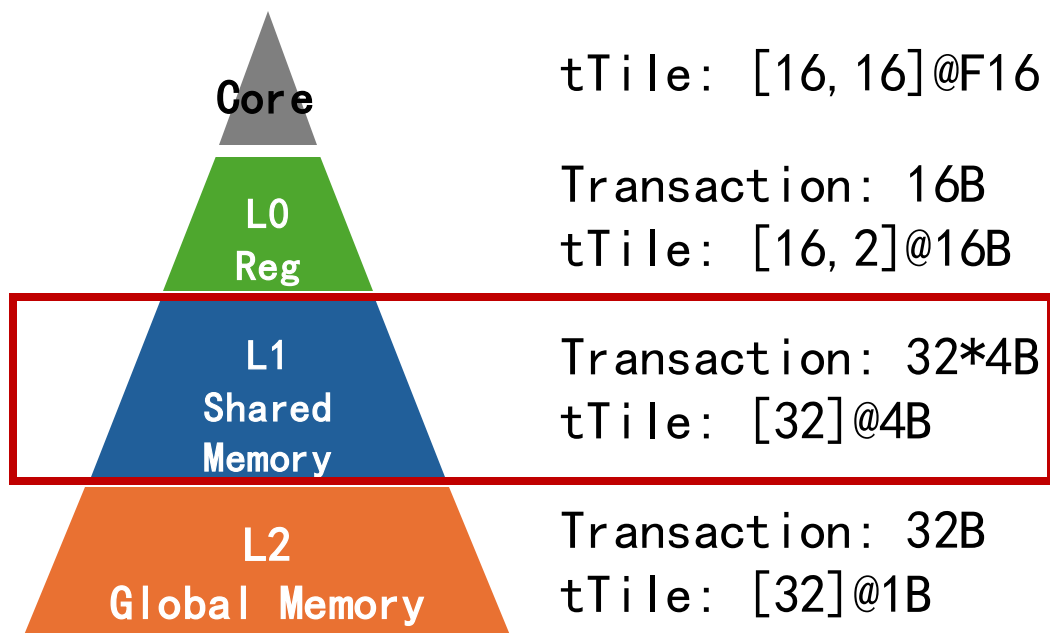
$B[16, 64]@FP16$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)



tTile-device for NVIDIA A100

target tTile: $[16, 16]@FP16$ or $[16, 2]@16B$
sliceB: $[4, 64]@FP16 \rightarrow [16, 16]@FP16$

B(0, 0-15)	B(1, 0-15)	B(2, 0-15)	B(3, 0-15)
B(4, 0-15)	B(5, 0-15)	B(6, 0-15)	B(7, 0-15)
B(8, 0-15)	B(9, 0-15)	B(10, 0-15)	B(11, 0-15)
B(12, 0-15)	B(13, 0-15)	B(14, 0-15)	B(15, 0-15)

$64*16bit = 128B$

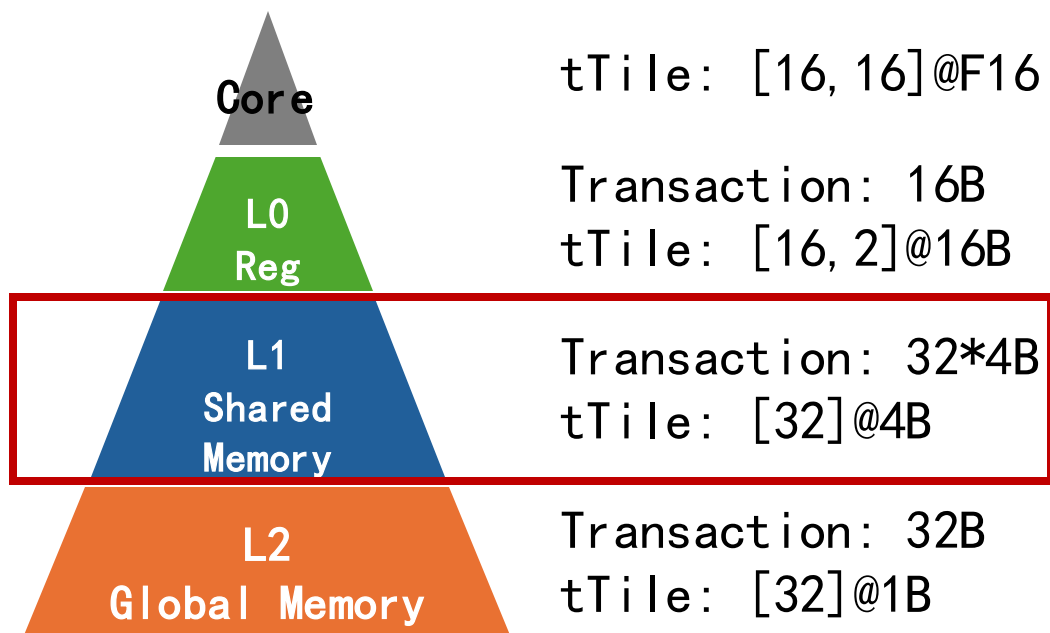
$B[4, 64]@FP16$

Example

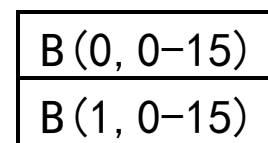
目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

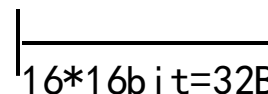
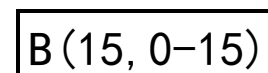
Device Hint: $[4, 8]@16B$ (size: 512B)



tTile-device for NVIDIA A100



...



B(0, 0-15)	B(1, 0-15)	B(2, 0-15)	B(3, 0-15)
B(4, 0-15)	B(5, 0-15)	B(6, 0-15)	B(7, 0-15)
B(8, 0-15)	B(9, 0-15)	B(10, 0-15)	B(11, 0-15)
B(12, 0-15)	B(13, 0-15)	B(14, 0-15)	B(15, 0-15)

$64*16bit = 128B$

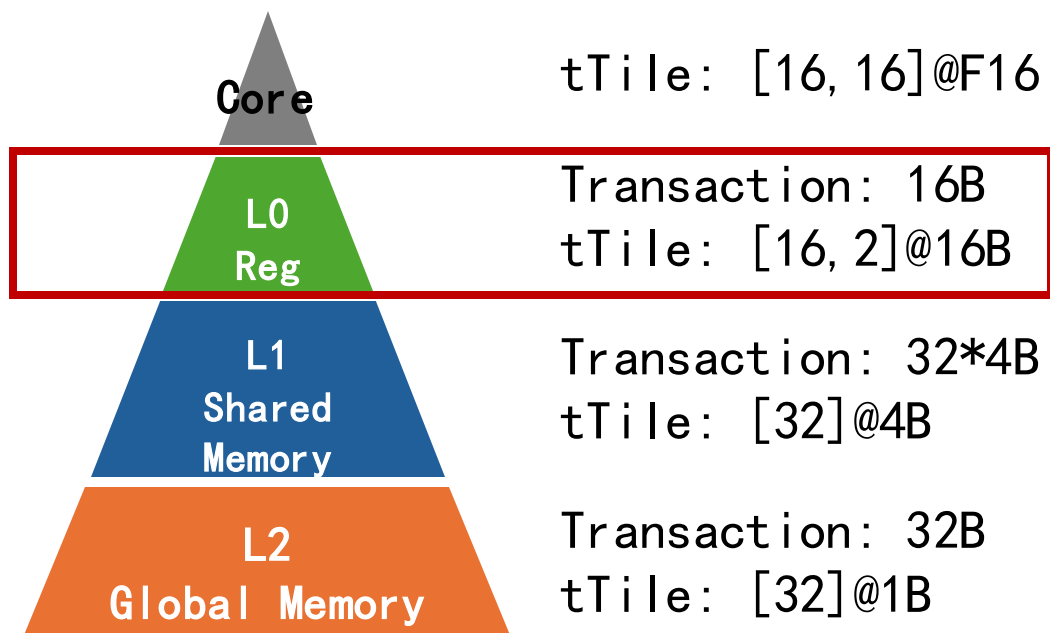
$B[4, 64]@FP16$

Example

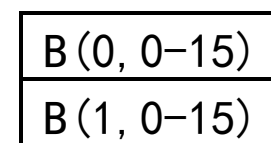
目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16),$
 $M=32, N=32, K=63$

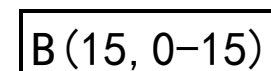
Device Hint: $[4, 8]@16B$ (size: 512B)



tTile-device for NVIDIA A100



...



$B(0, 0-15)$	$B(1, 0-15)$	$B(2, 0-15)$	$B(3, 0-15)$
$B(4, 0-15)$	$B(5, 0-15)$	$B(6, 0-15)$	$B(7, 0-15)$
$B(8, 0-15)$	$B(9, 0-15)$	$B(10, 0-15)$	$B(11, 0-15)$
$B(12, 0-15)$	$B(13, 0-15)$	$B(14, 0-15)$	$B(15, 0-15)$

$64 * 16 \text{ bit} = 128B$

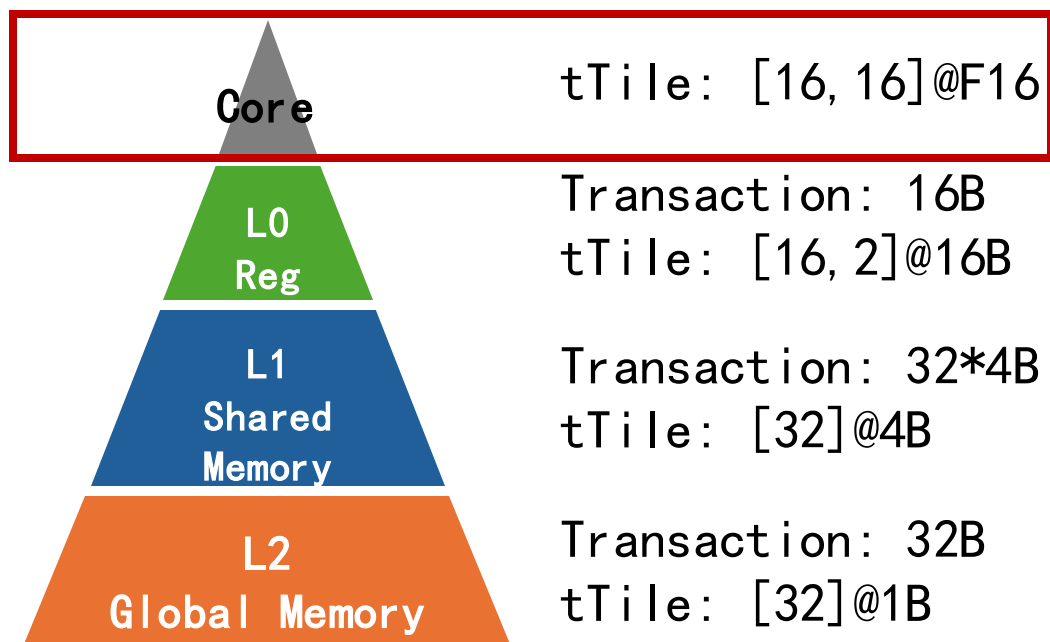
$B[4, 64]@FP16$

Example

目标: $A[32, 63]@FP16$ 和 $B[32, 63]@NF4$ 的MMA

$C = \text{compute}((M, N), \text{lambda } i, j: (\text{sum}((A[i, k]@FP16 * B[j, k]@NF4)@FP32)@FP32)@FP16)$,
 $M=32, N=32, K=63$

Device Hint: $[4, 8]@16B$ (size: 512B)

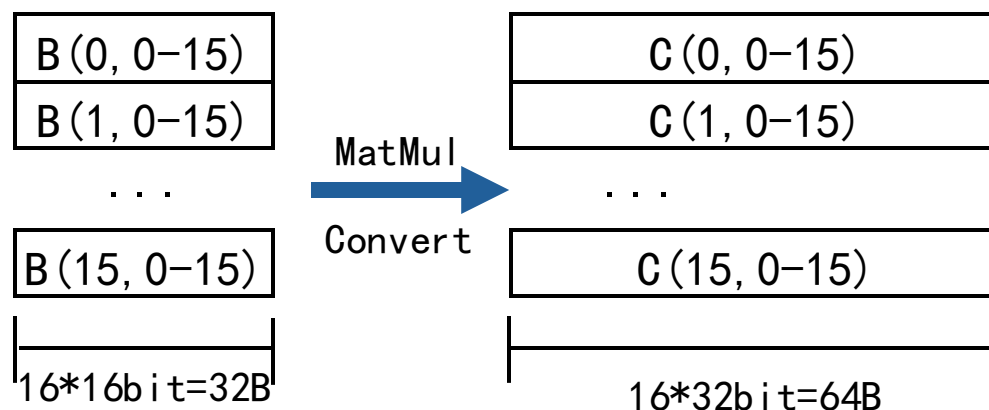


tTile-device for NVIDIA A100

进入Tensor Core进行运算

将结果C:

convertC: $[16, 16]@FP32$

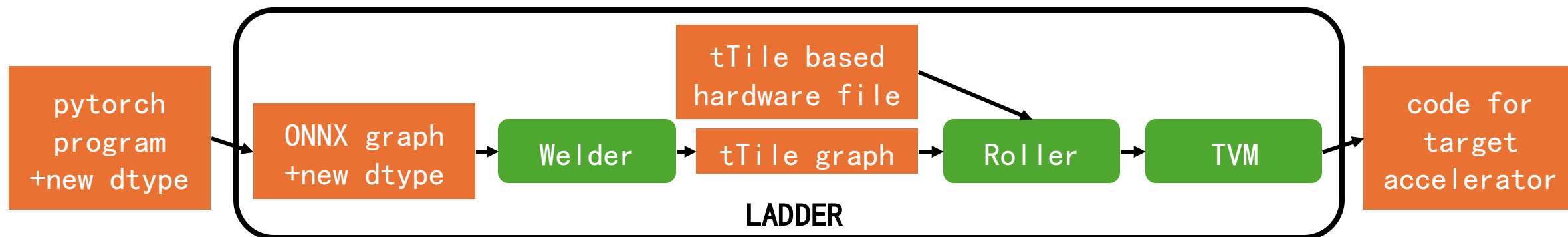


$B[16, 16]@FP16$

$C[16, 16]@FP32$

Implementation

LADDER基于 TVM / Welder / Roller三项工作实现，支持低精度、非标准数据类型。



Welder: 端到端图优化，如算子融合，内存布局转换

Roller: tTile配置自动搜索

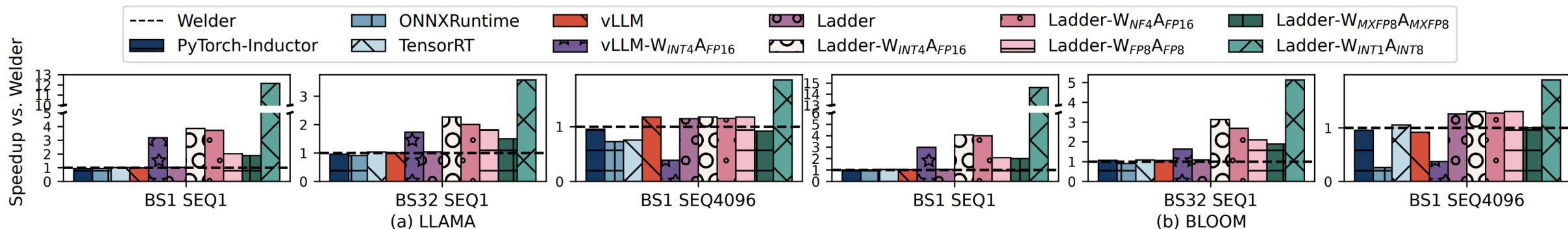
TVM: 内核调度与代码生成（LOP3优化等）

支持NVIDIA CUDA GPU和AMD ROCm GPU

技术细节：支持非 2^n 数据如3-bit，PTX指令控制，low bit转换优化等

Evaluation: End-to-end result on A100

LLM Inference

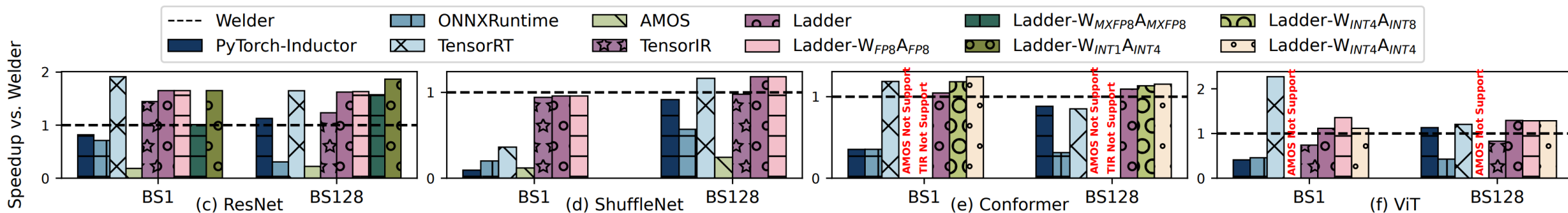


W@FP16 / A@FP16: **1.1x/1.1x** avg. speedup over Welder/TensorRT

W@INT4 / A@FP16(GPTQ): **2.3x** avg. speedup over vLLM

W@INT1 / A@INT8(BitNet): **up to 8.8x** speedup

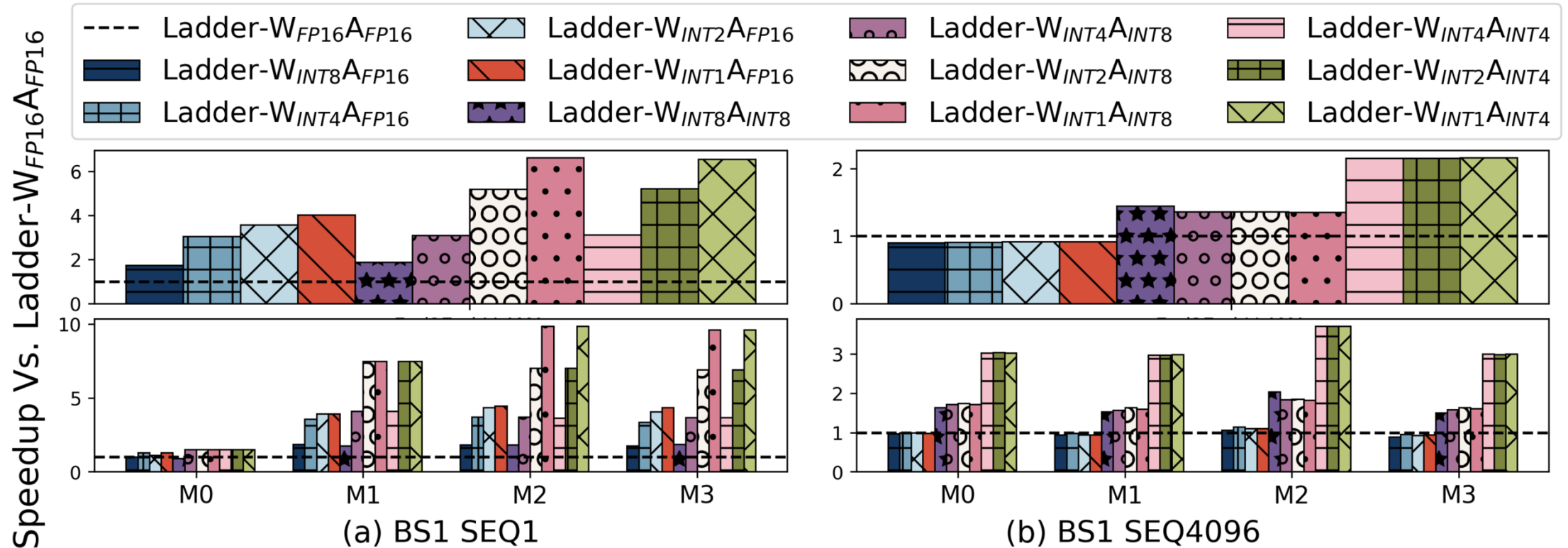
DNN Inference



W@FP16 / A@FP16: **1.4x/1.7x** avg. speedup over Welder/TensorRT

Enable low-precision DNN computing with up to **3x** speedup over Ladder- W@FP16 / A@FP16

Evaluation – Scaling Bit Width



BS1 SEQ1: bounded by memory bw., up to **6.4x** speedup (10.1x speedup on kernel)

BS1 SEQ4096: bounded by tensor core, up to **2.4x** speedup (3.7x speedup on kernel)

LADDER

Background:

- **低精度计算**不断演化
- 软硬件适配难以同步发展

LADDER提出了**TypeTile**抽象和**tensor transformation**的方法:

- 表示不断进化的新数据类型
- 为低精度计算提供系统化的机制
- 充分发挥模型和硬件的性能

谢谢!