

FLEXIBLE MATRIX MULTIPLICATION

KERNELS ON GPUS

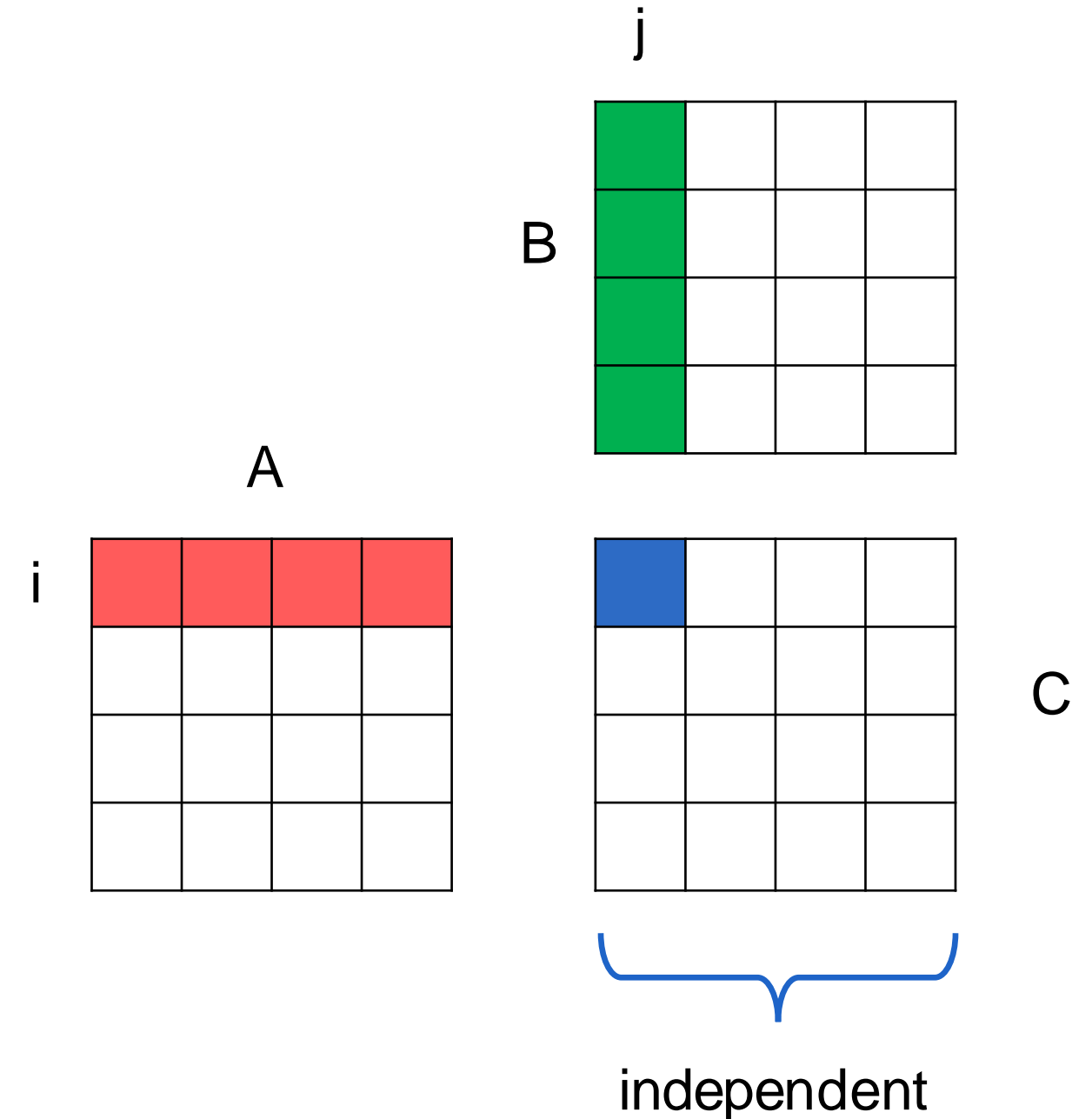
Thomas Faingnaert
1 July 2020

Supervisor: Prof. Dr. Ir. Bjorn De Sutter
Counsellor: Dr. Tim Besard

GENERAL MATRIX MULTIPLICATION (GEMM)

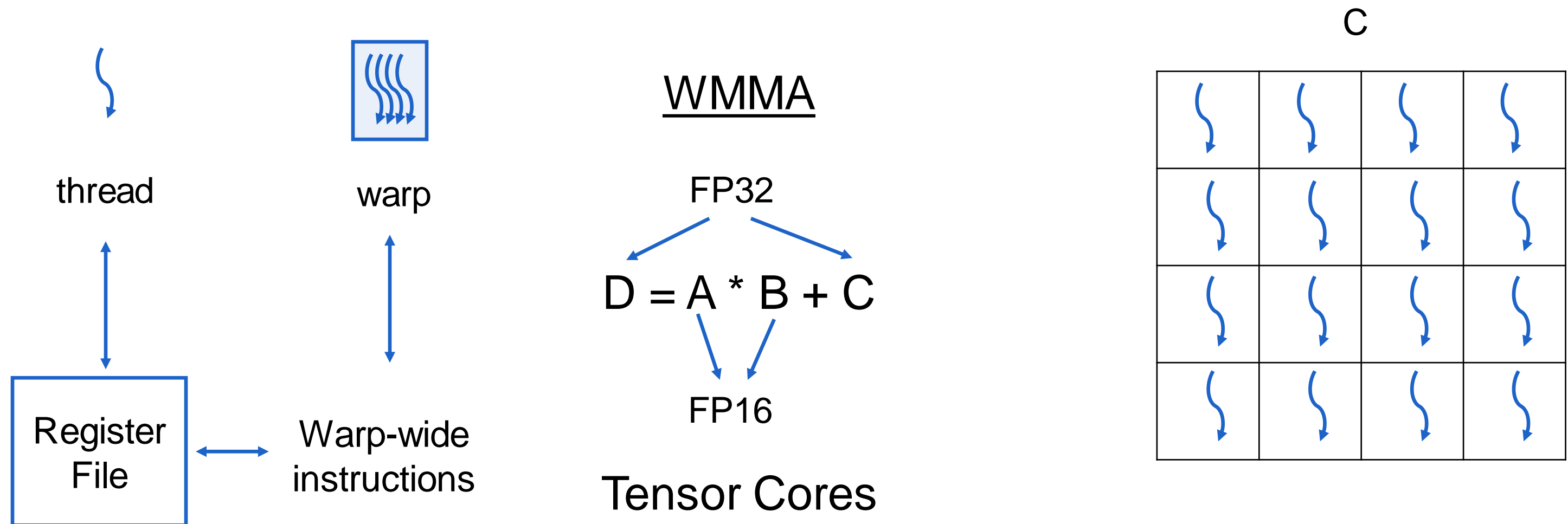
$$C := A * B + C$$

```
for i = 1 : N
  for j = 1 : N
    for k = 1 : N
      C[i, j] += A[i, k] * B[k, j]
    end
  end
end
```



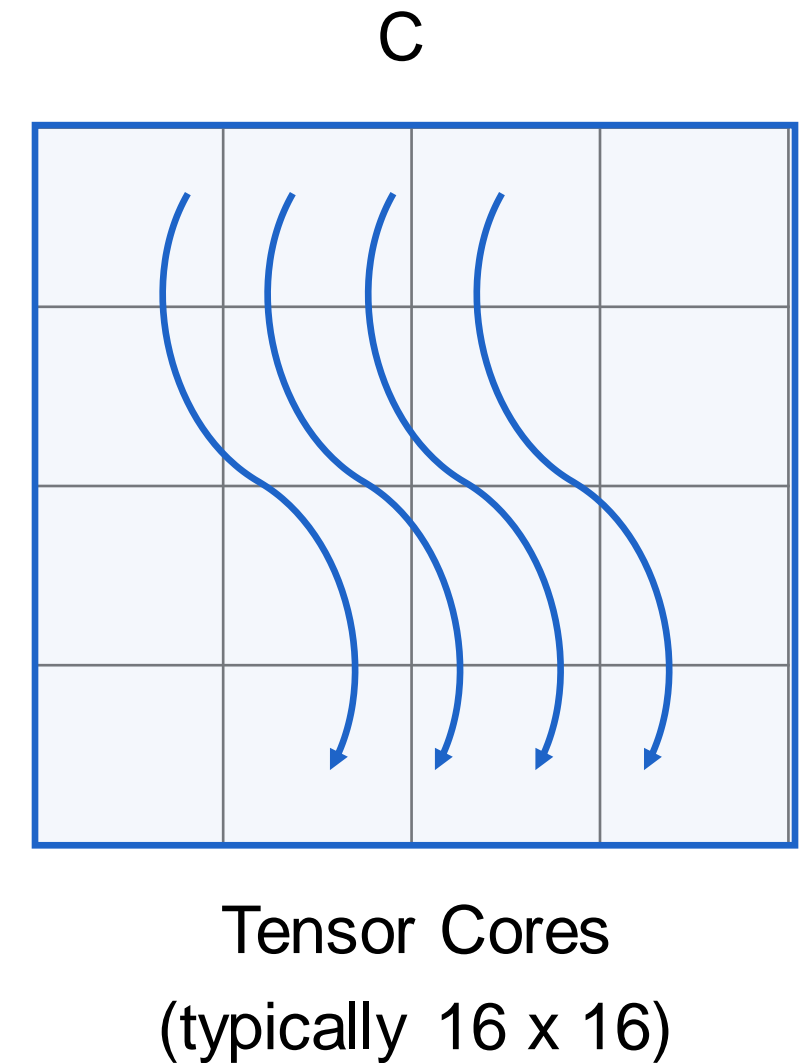
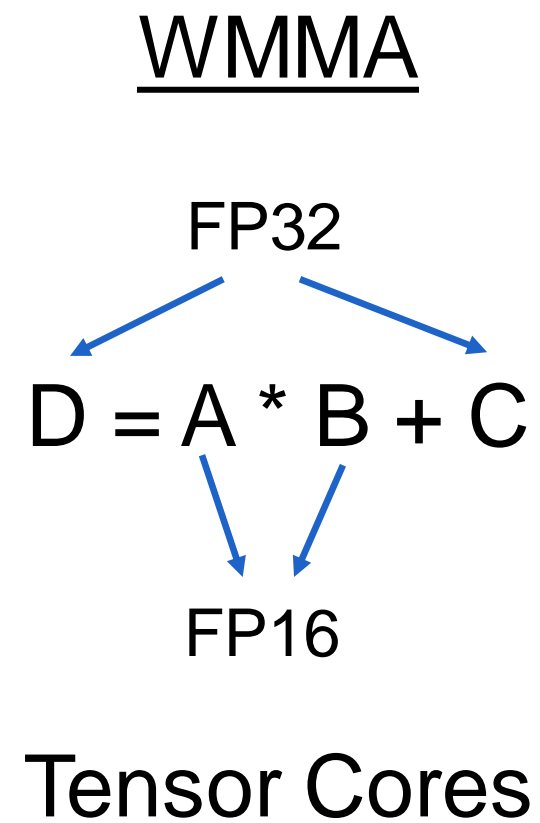
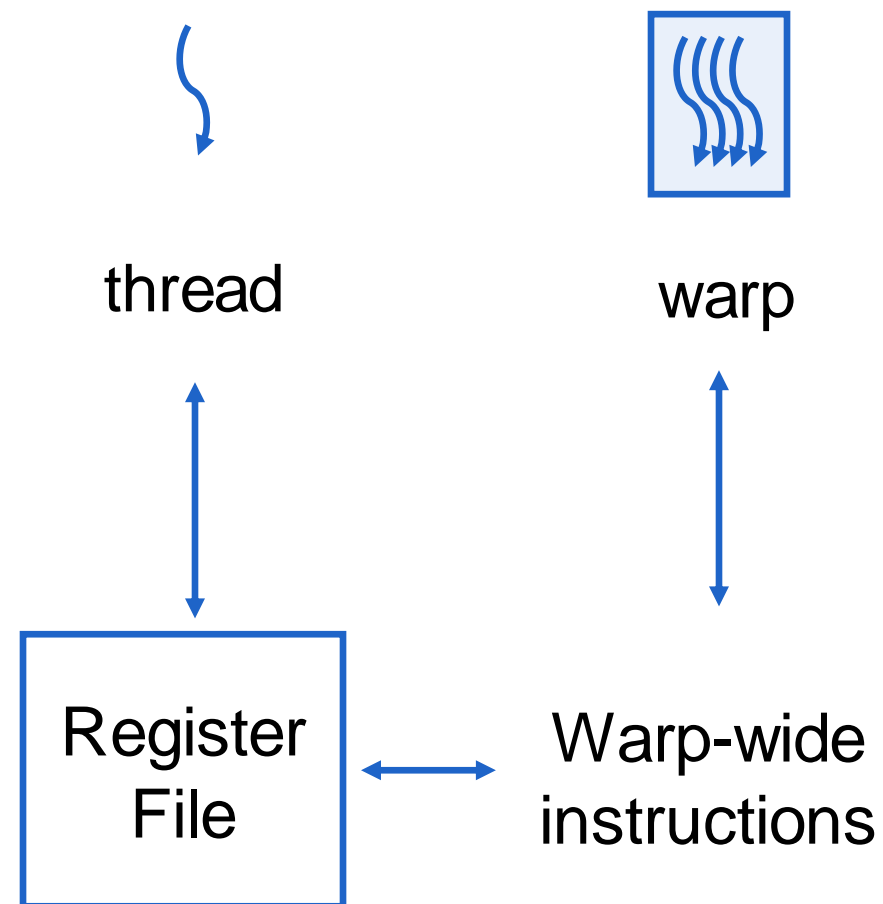
GENERAL MATRIX MULTIPLICATION ON GPUS

Massively parallel processing



GENERAL MATRIX MULTIPLICATION ON GPU

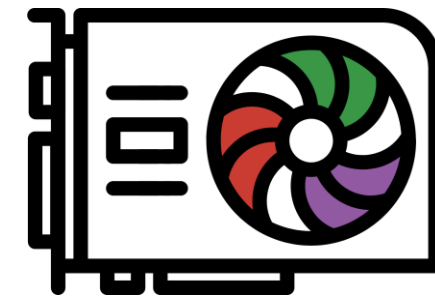
Massively parallel processing



TOOL FLOW



LLVM



CUDAAnative.jl

WMMA API



NVIDIA Turing GPU

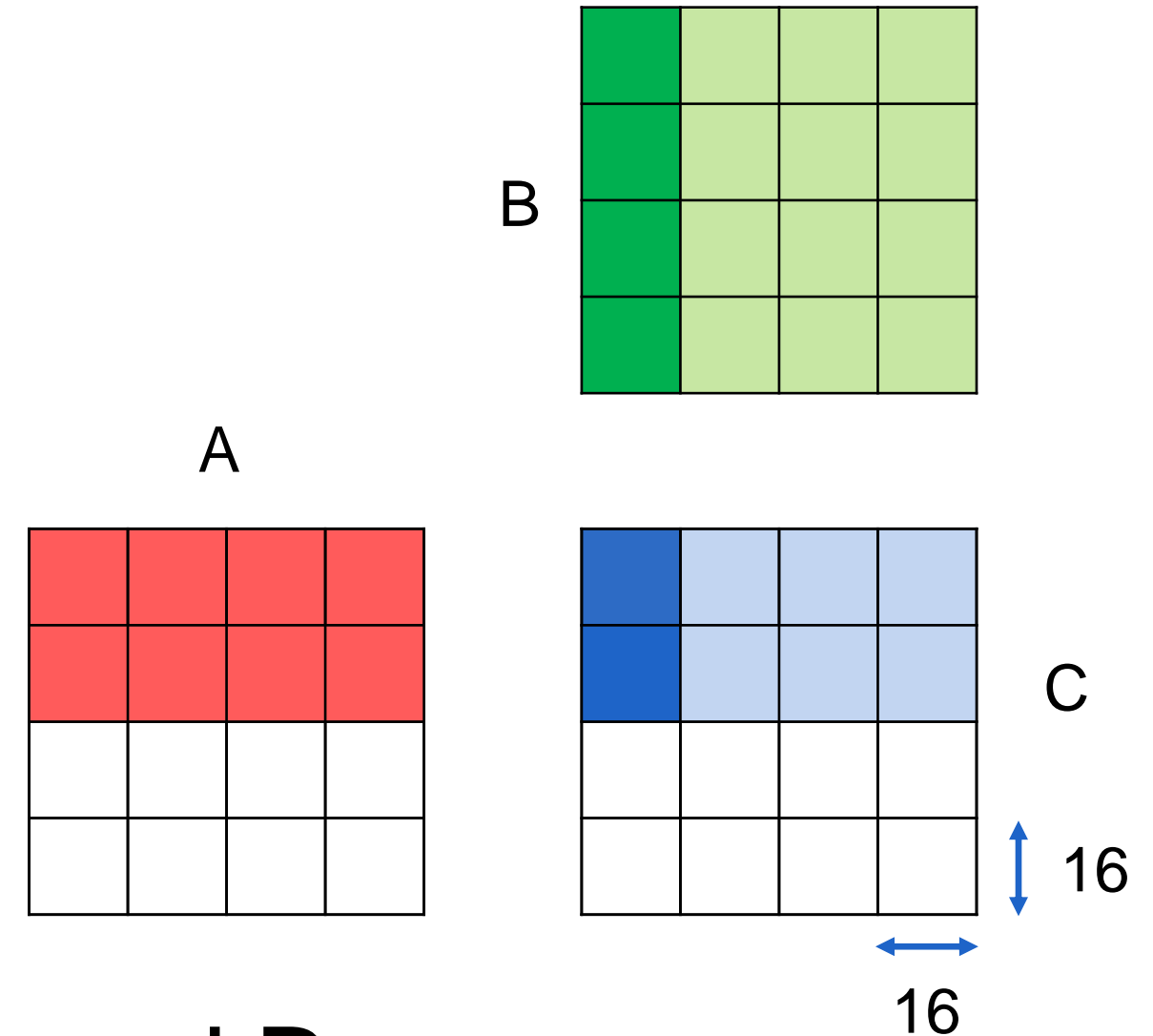
Tensor Cores

Codegen
changes

GENERAL MATRIX MULTIPLICATION ON GPU

$$C := A * B + C$$

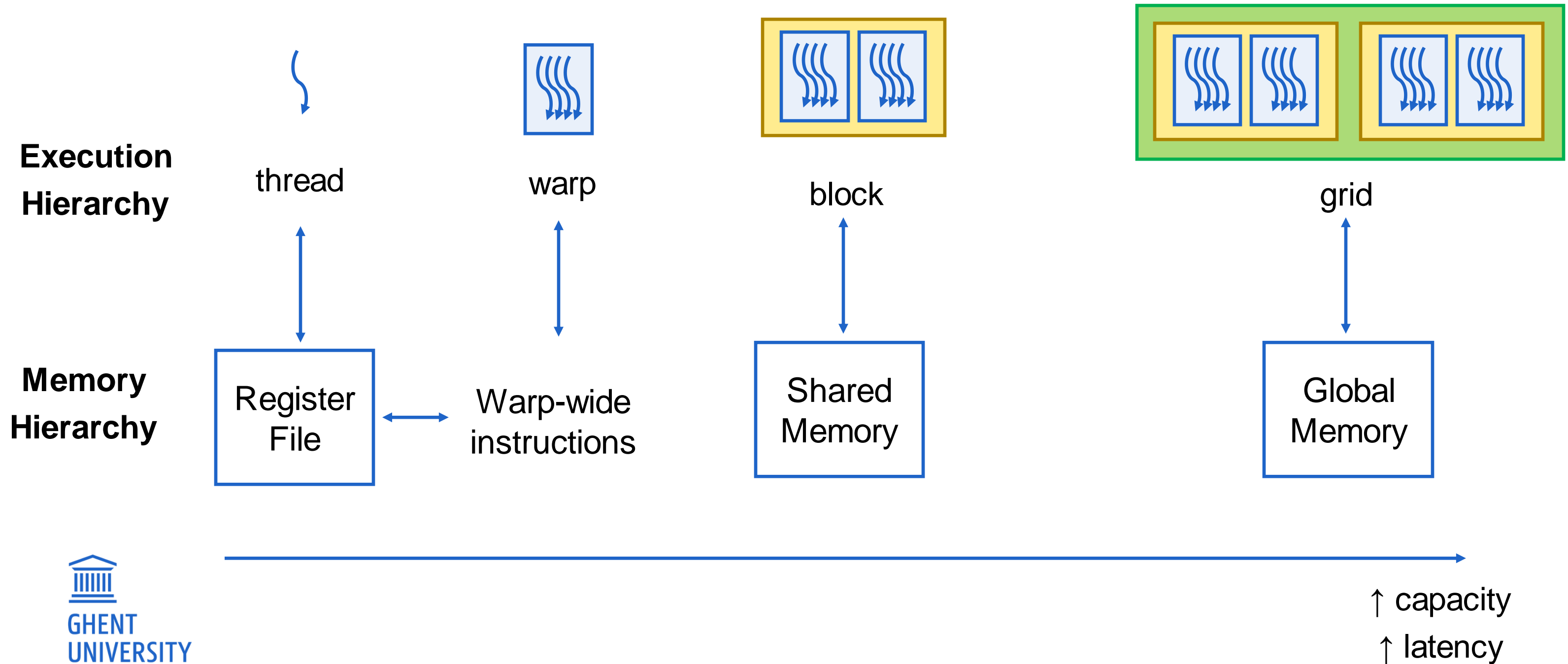
```
for i = 1 : 16 : N
  for j = 1 : 16 : N
    for k = 1 : 16 : N
      # calculate a 16 x 16 matrix mult.
    end
  end
end
```



Inefficient: large working set to store A and B
→ Exploit data reuse: temporal locality

GENERAL MATRIX MULTIPLICATION ON GPUS

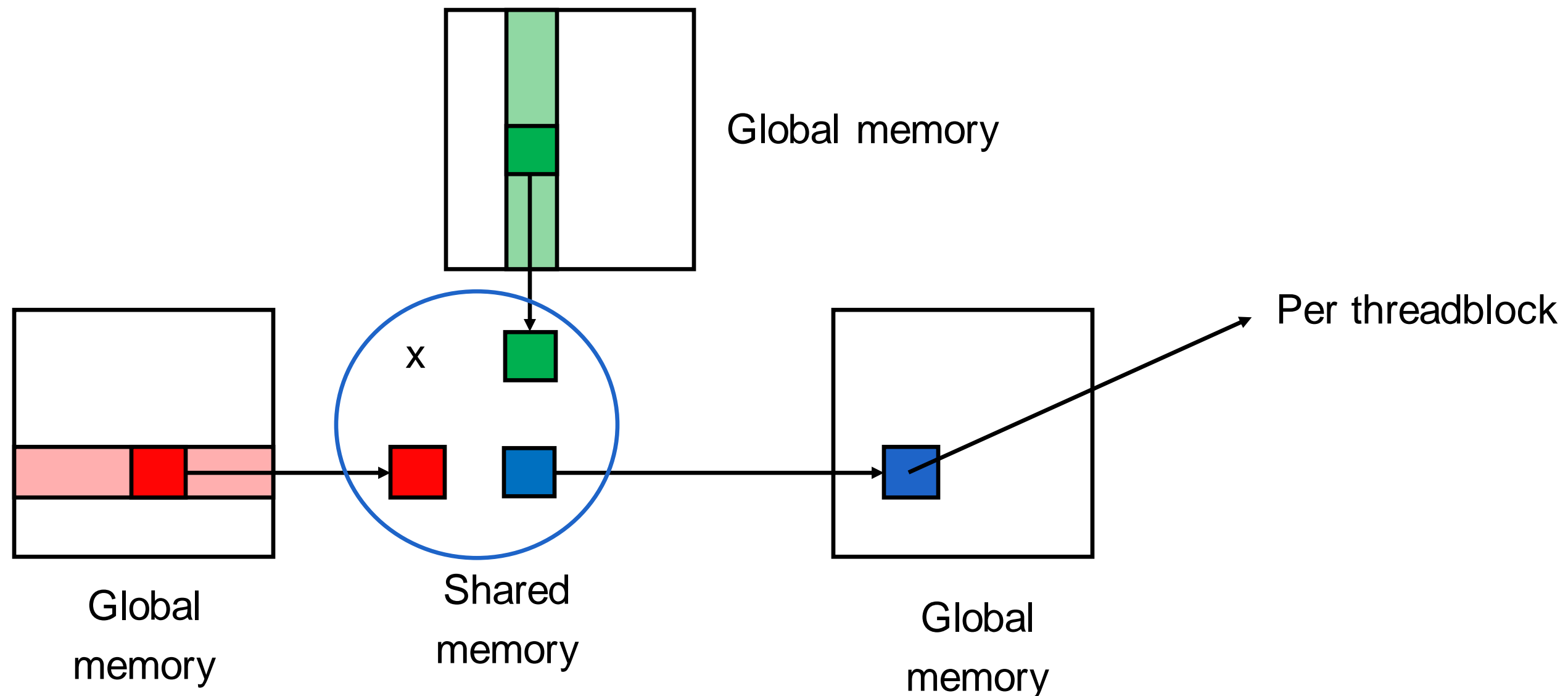
Massively parallel processing



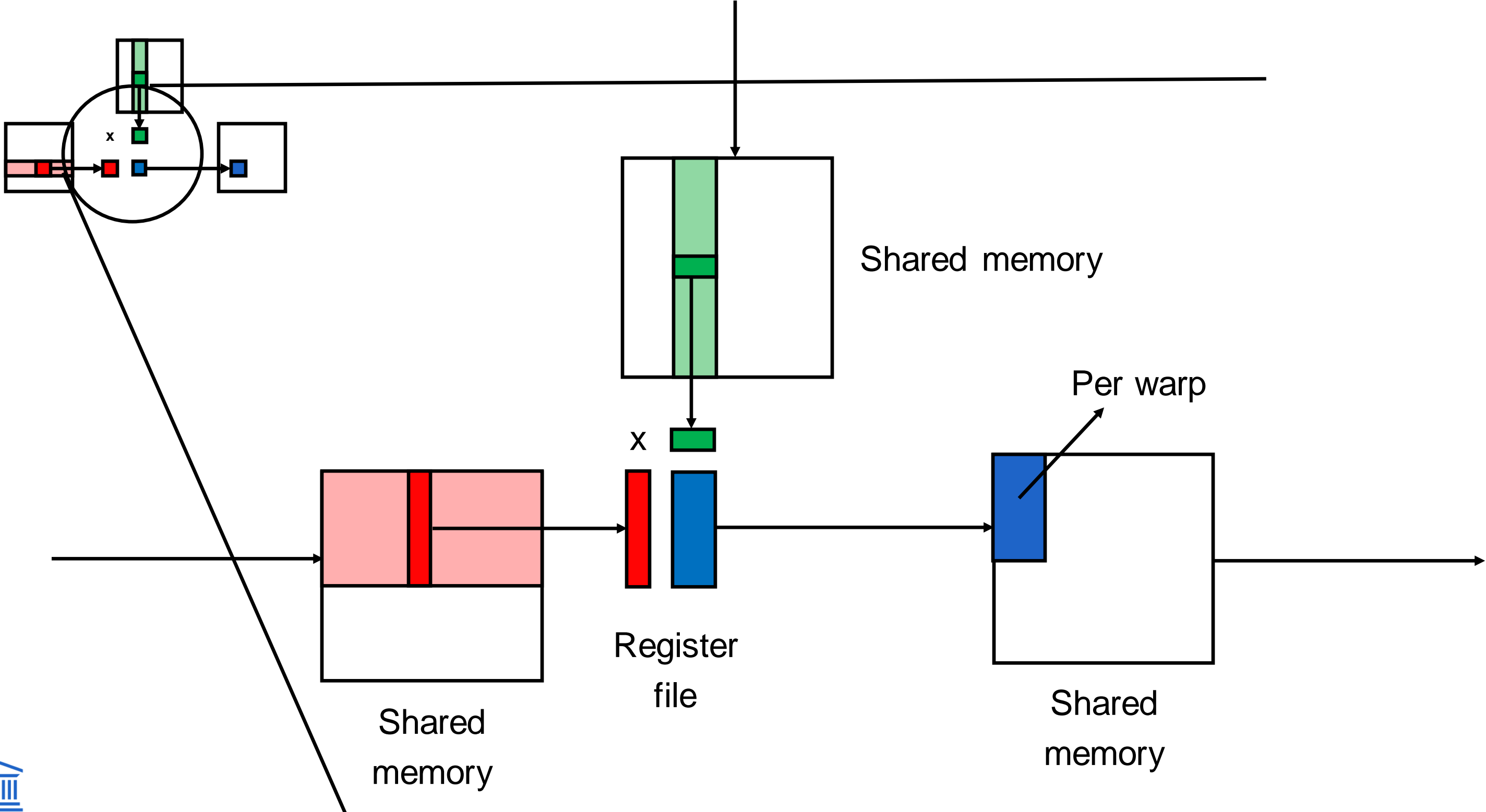
DESIGN OF PERFORMANT GEMM ON GPUS

Small N: perfect temporal locality

Large N: explicit blocking techniques



DESIGN OF PERFORMANT GEMM ON GPUS



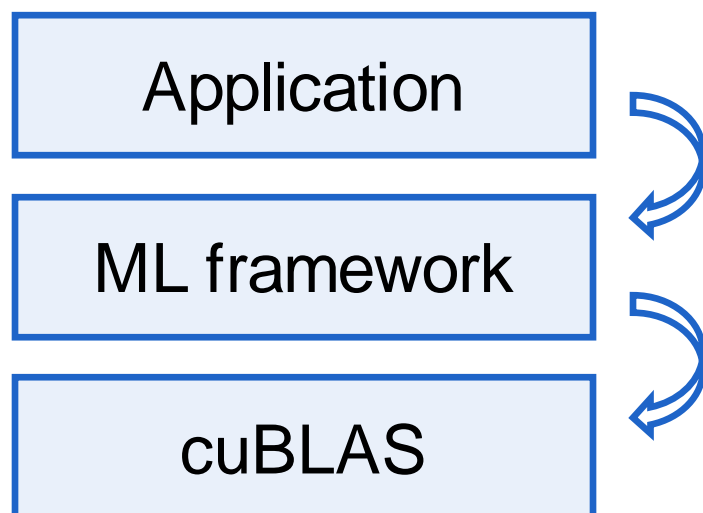
GENERAL MATRIX MULTIPLICATION

$$C := A * B + C$$

```
for block_tile_i = ...  
  ⋮  
  for k = ...  
    C[i, j] += A[i, k] * B[k, j]  
  end  
  ⋮  
end
```

cuBLAS
cuDNN

THE NEED FOR FLEXIBILITY



Problem:

- Only limited set of kernels

THE NEED FOR FLEXIBILITY

Goal: framework for flexible GEMM

Which flexibility is needed?

- Literature study
- CUTLASS

THE NEED FOR FLEXIBILITY

$$C := A * B + C$$

```
for block_tile_i = ...  
  ⋮  
  for k = ...  
    C[i, j] += A[i, k] * B[k, j]  
  end  
  ⋮  
end
```

Memory layout?

- N vs. T
- NCHW vs. NHWC
- Tensors

THE NEED FOR FLEXIBILITY

$C := A * B + C$

```
for block_tile_i = ...  
  ⋮  
  for k = ...  
    C[i, j] += A[i, k] * B[k, j]  
  end  
  ⋮  
end
```

Data types?

- FP16, FP32, FP64
- Mixed precision
- Complex numbers
- User-defined type

THE NEED FOR FLEXIBILITY

$C := f(A * B + C)$

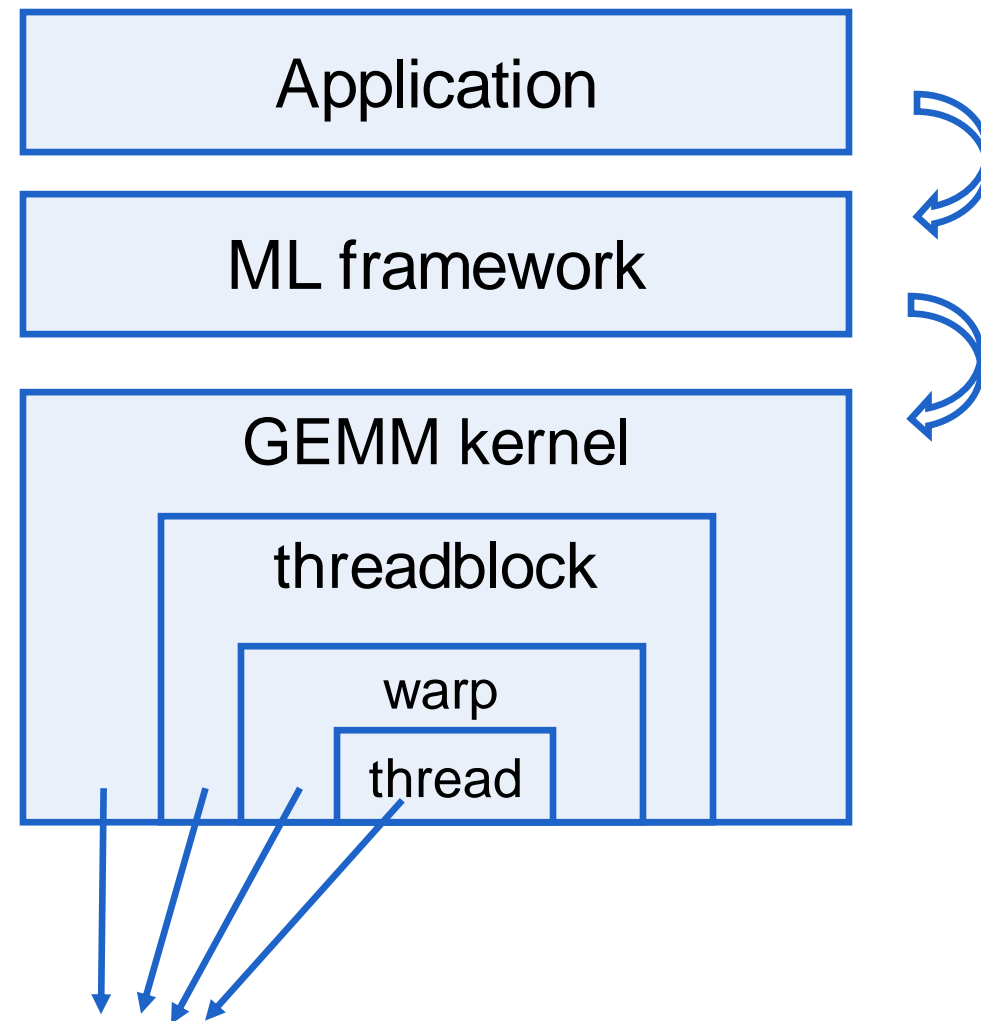
```
for block_tile_i = ...  
  ⋮  
  for k = ...  
    C[i, j] += f(A[i, k] * B[k, j])  
  end  
  ⋮  
end
```

Elementwise operations?

- Scaling
- Activation function

FLEXIBLE GEMM APIS

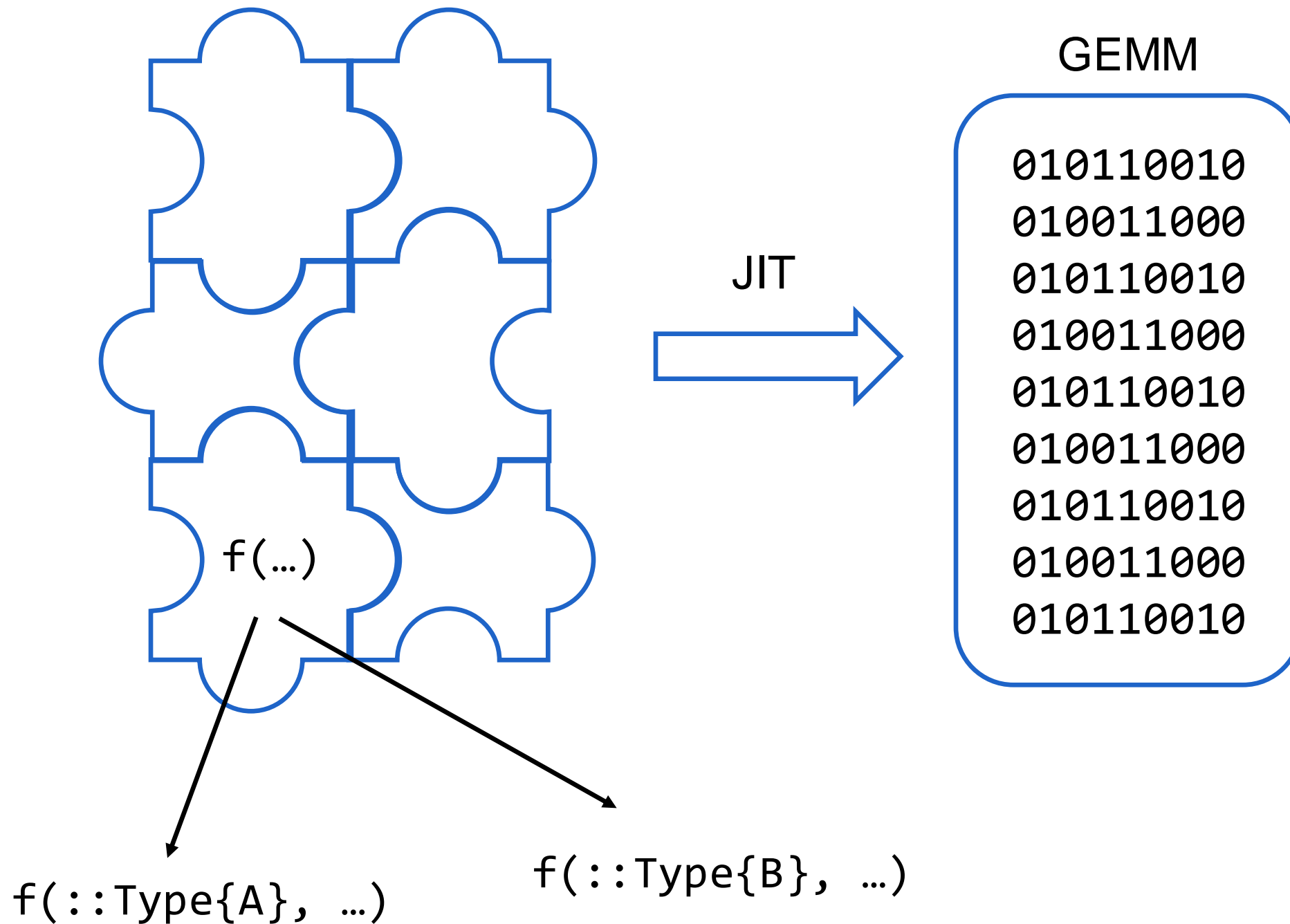
CUTLASS



- Large number of components
- Low-level C++

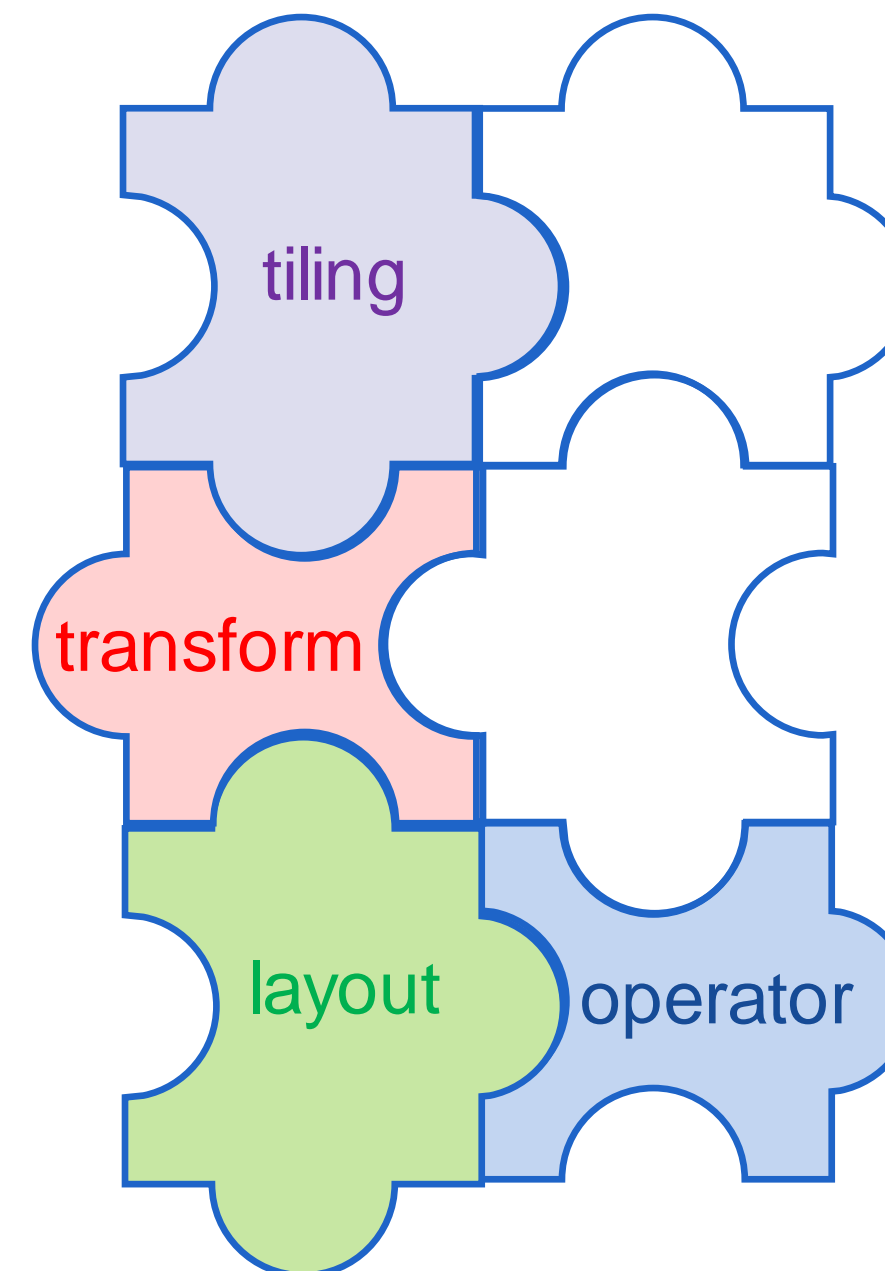
Templates specialised on element datatype, memory layout, ...

FLEXIBLE GEMM APIS IN JULIA



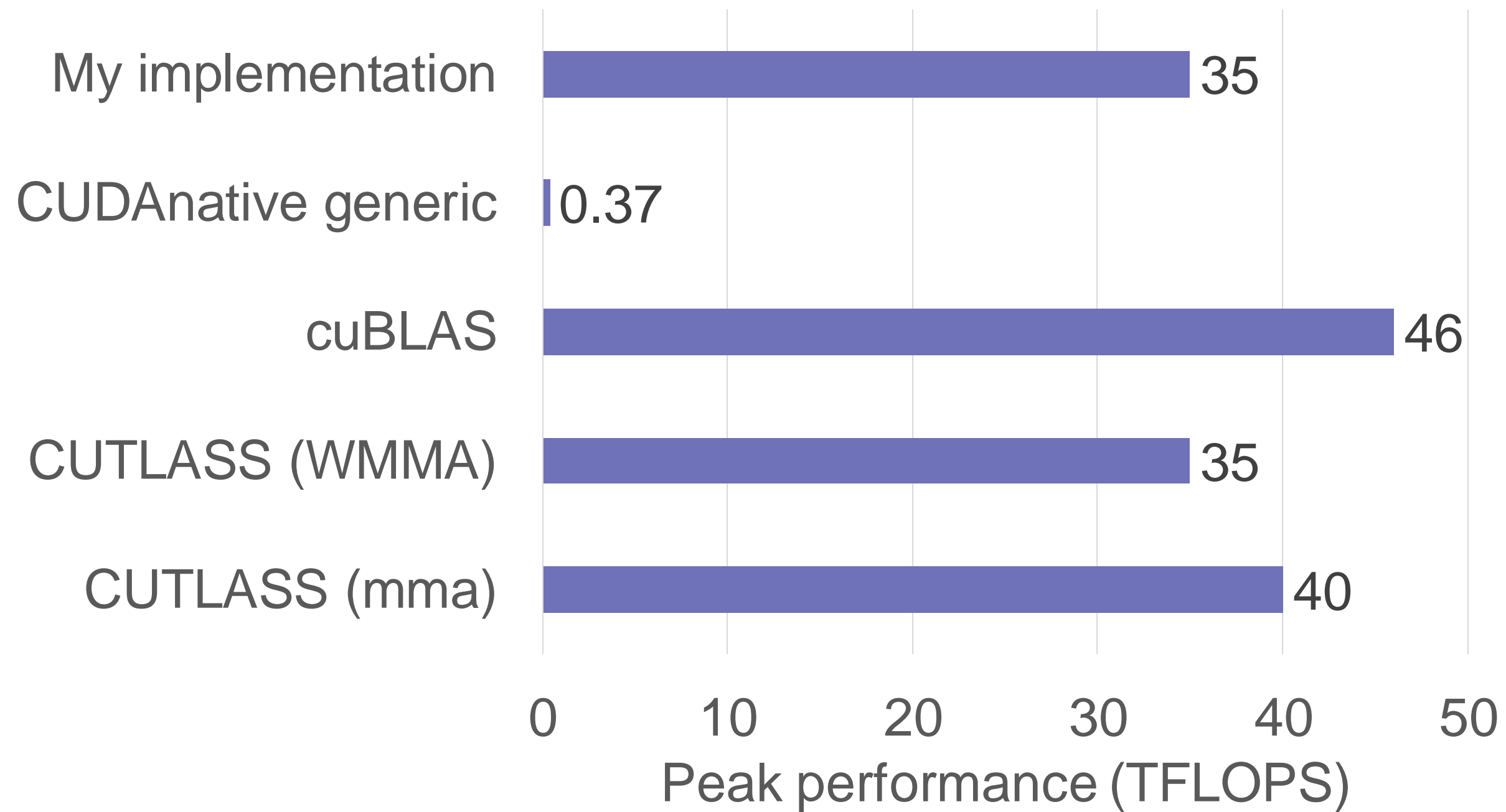
FLEXIBLE GEMM APIS IN JULIA

```
for block_tile_i = ...  
  ⋮  
  for k = ...  
    C[i, j] += f(A[i, k] * B[k, j])  
  end  
  ⋮  
end
```



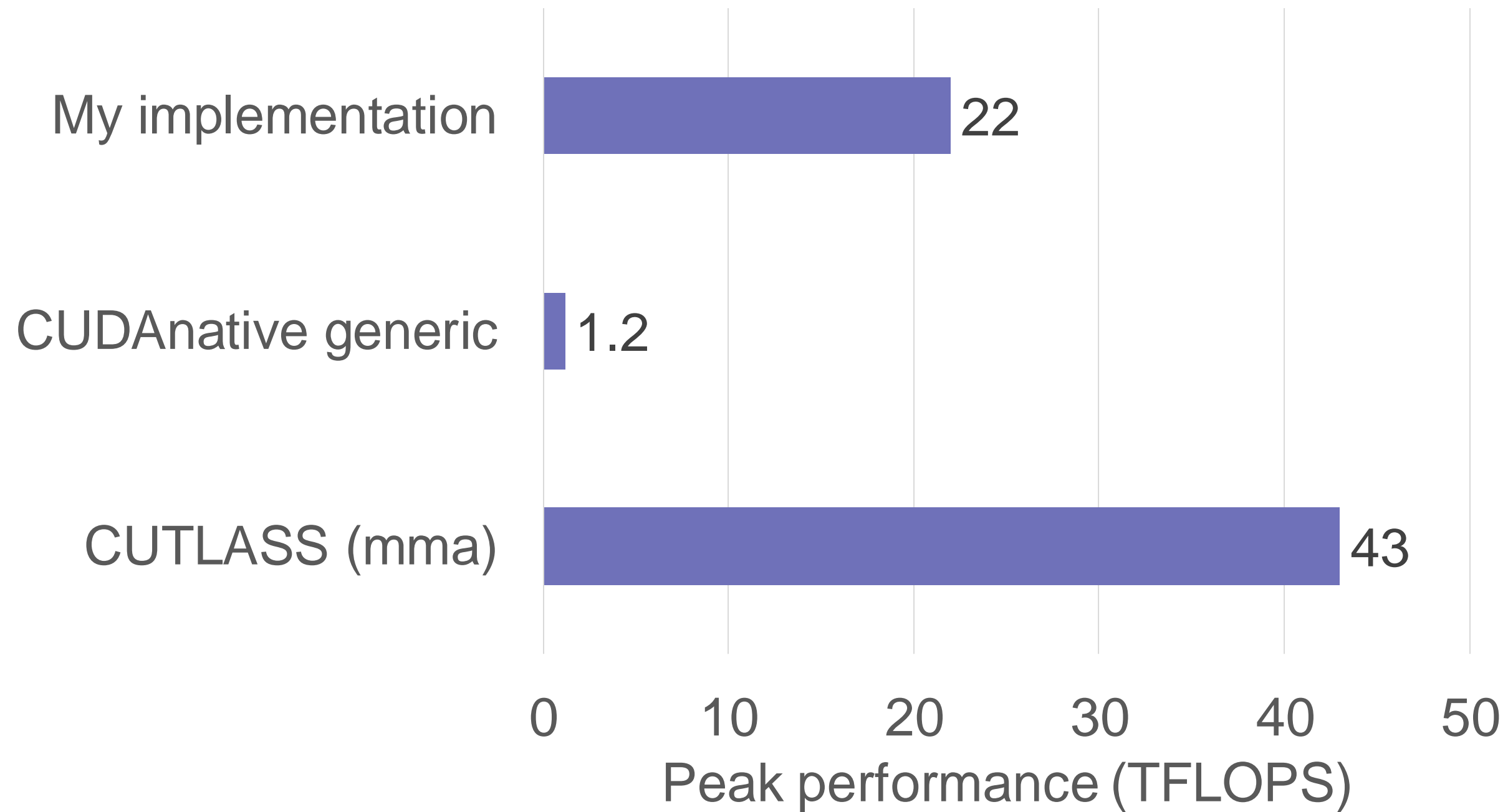
EVALUATION

1. Mixed-precision GEMM (WMMA)



EVALUATION

2. Complex mixed-precision GEMM (WMMA)

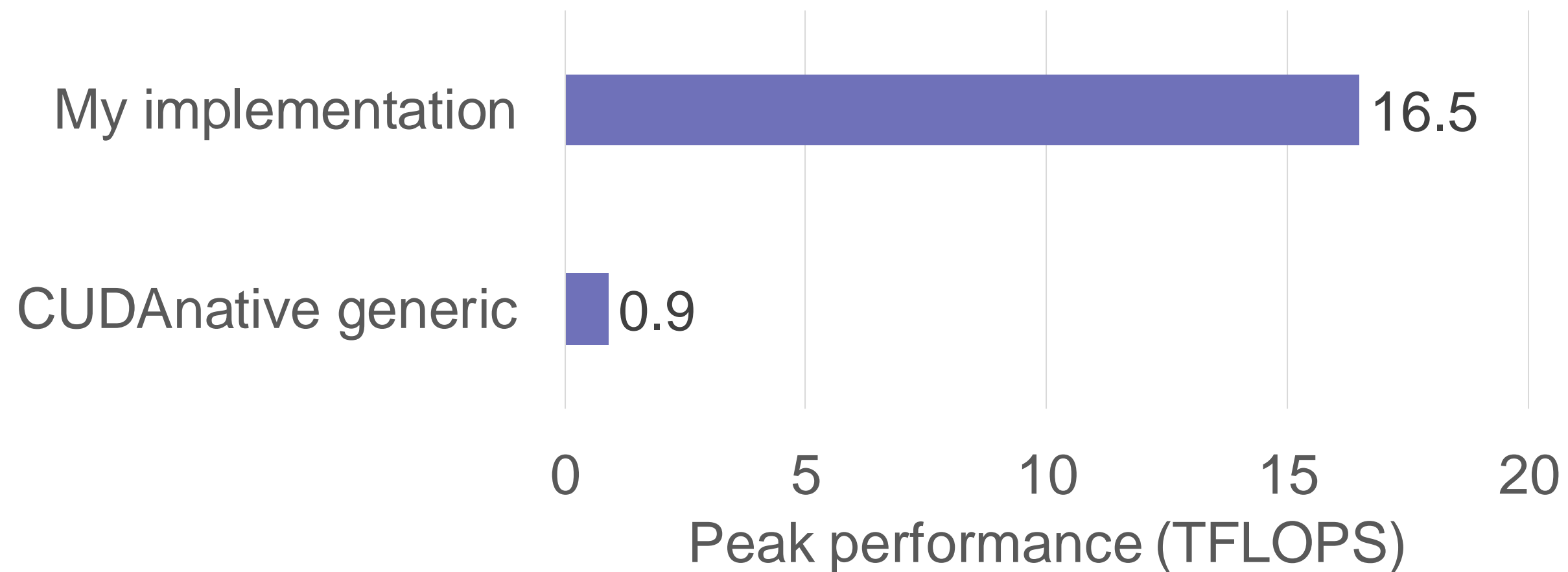


EVALUATION

3. Dual mixed-precision GEMM (WMMA)

$$(a + b\varepsilon) * (c + d\varepsilon) = ac + (ad + bc)\varepsilon$$

Application: automatic differentiation



QUESTIONS?