

# **An Introduction to Labs (2)**

**Aaron Zhao, Imperial College London**

# Overview

- Lab 3:
  - Mixed Precision Search
- Lab 4:
  - Software stream: Automated and manual low-level optimizations
  - Hardware stream: Emitting SystemVerilog from Pytorch models

# Mixed Precision Search

# Mixed precision search

- Different blocks of a neural network have different numerical characteristics.
- It naturally makes sense to use different precisions for different blocks.

## Algorithm 1 Transformer layer

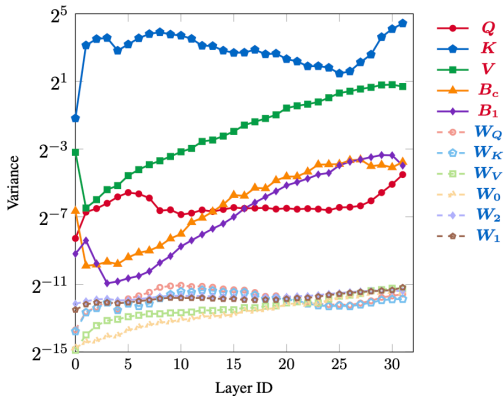
**Require:**  $X$  ▷ Input features

**Require:**  $H$  ▷ Number of heads

```

1:  $X_n \leftarrow \text{LayerNorm}(X)$ 
2: for  $i \in [0, H)$  do
3:    $Q_i \leftarrow X_n W_{Q_i}$  ①
4:    $K_i \leftarrow X_n W_{K_i}$  ②
5:    $V_i \leftarrow X_n W_{V_i}$  ③
6:    $A_i \leftarrow \frac{Q_i K_i^T}{\sqrt{d_k}}$  ④
7:    $\hat{A}_i \leftarrow \text{softmax}(A_i, \text{axis} \leftarrow -1)$ 
8:    $B_i \leftarrow \hat{A}_i V_i$  ⑤
9: end for
10:  $B_c \leftarrow \text{concat}(B_0, \dots, B_{H-1})$ 
11:  $B_0 \leftarrow B_c W_0 + b_0$  ⑥
12:  $B_n \leftarrow \text{LayerNorm}(B_0 + X)$ 
13:  $B_1 \leftarrow \text{ReLU}(B_n W_1 + b_1)$  ⑦
14:  $B_2 \leftarrow B_1 W_2 + b_2$  ⑧
15:  $O \leftarrow B_2 + B_0 + X$ 
16: return  $O$ 

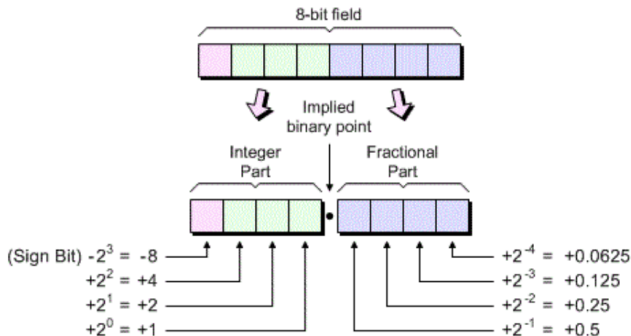
```



# Integer quantization

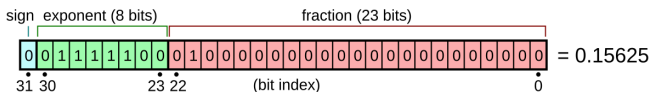
Depending on your background, the definition of integer quantization can be different. We define integer quantization as the process of converting a floating-point number to an integer number with a fixed decimal point (same as fixed-point arithmetic).

These are also known as linear quantization.



# Non-linear quantization

- Float based
  - Sign (Blue)
  - Exponent (Green)
  - Mantissa (Red)
  - Exponent bias (a constant)
- Logarithmic based



**Software stream: Automated and manual low-level optimizations**

## torch.compile

`torch.compile` is a tool that can be used to compile PyTorch models to a lower-level representation. It is currently in the PyTorch ecosystem, and it is a tool that can be used to optimize PyTorch models.

One can see this as an automated way to perform low-level optimizations, replying on optimization flows that are already implemented.

**FX-Graph:** FX-Graph is a symbolic representation of PyTorch models that is designed to be more amenable to optimization

**TorchDynamo:** TorchDynamo is responsible for JIT compiling arbitrary Python code into FX graphs

**TorchInductor:** TorchInductor is responsible for lowering FX graphs into a efficient kernel representations

**Triton:** the language that is used to actually implement these low-level kernel representations.



## Fused kernels

In this lab, we take a look at what is a fused kernel that implements ScaledDotProductAttention.

Kernel fusion can happen at different levels, we can fuse kernels at the level of the computation graph, or we can fuse kernels at the level of the actual implementation.

In this example in the lab, we fuse at the pythonic level, which is a higher level of abstraction.

# Porting a custom kernel

We aim to port an MXInt kernel in this lab.

What would you do?

- It is a custom kernel that is implemented in CUDA (CUTLASS CUTE)
- CUTLASS is a template library that provides highly tuned matrix multiplication kernels for NVIDIA GPUs using CUDA.
- This requires a matching GPU CC (compute capability).
- You will have to build this kernel from scratch

# Porting a custom kernel

## What is MXInt?

### IEEE Float32 (FP32)

1-bit sign, 8-bit exponent, 23-bit mantissa



### IEEE Float16 (FP16)

1-bit sign, 5-bit exponent, 10-bit mantissa



### MiniFloat / Denormed Minifloat (DMF)

1-bit sign, 4-bit exponent, 3-bit mantissa



### Block Minifloat (BM)

1-bit sign,  $E$ -bit exponent,  $M$ -bit mantissa  
 $B$ -bit shared exponent bias



...



### Block Floating Point (BFP)

1-bit sign,  $M$ -bit mantissa  
 $E$ -bit shared exponent



...



### Block Logarithm (BL)

1-bit sign,  $E$ -bit exponent  
 $B$ -bit shared exp bias



...



# **Hardware stream: Emitting SystemVerilog from Pytorch models**

## Porting a custom kernel (iii)

The goal of this lab is to automatically generate a fully-connected layer in SystemVerilog, and test it using Cocotb.

**Cocotb** is a COroutine based COsimulation TestBench environment for verifying VHDL and SystemVerilog RTL using Python.

We use the Cocotb with the **Verilator** backend.

**Cocotb**: direct testing in Python, no need to write a testbench in SystemVerilog.

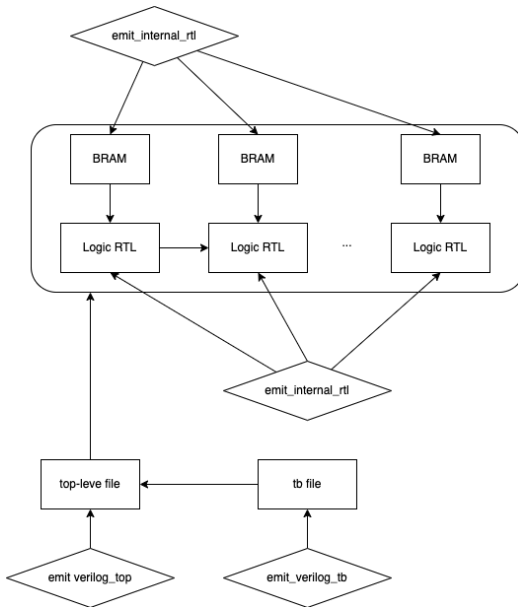
**Verilator**: ‘up-compiler’ SystemVerilog into multithreaded C++, lightening fast, no need to open vendor tools when doing behavior level testing.

## EmitVerilog Pass in MASE

- Classic source to source generation
- Directly generate SystemVerilog from MaseGraph

```
from chop.passes.graph.transforms import (  
    emit_verilog_top_transform_pass,  
    emit_internal_rtl_transform_pass, emit_bram_transform_pass,  
    emit_verilog_tb_transform_pass)
```

## EmitVerilog Pass in MASE (ii)



## EmitVerilog Pass in MASE (iii)

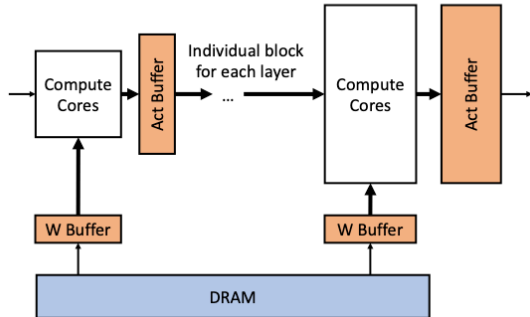
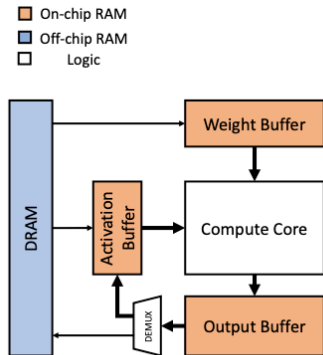
EmitVerilog generates dataflow designs

- Generate functional elements (RTL)
- Generate memory components (BRAM)
- Dataflow accelerator design without making use of the DRAM



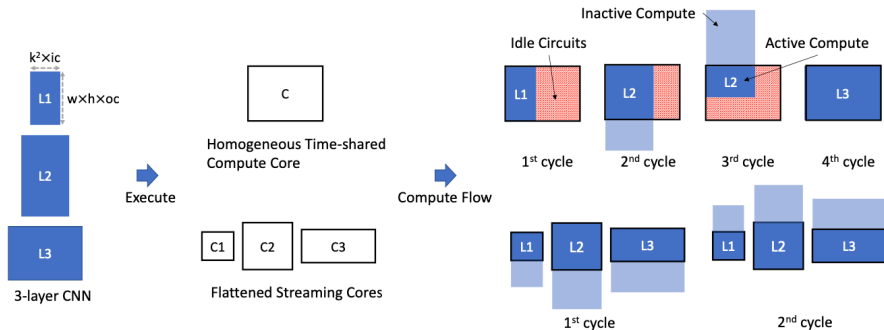
# Dataflow accelerator designs

- A homogeneous Big Compute Core (normal design, ASIC)
- A series of tailored small compute cores (dataflow design, FPGA)



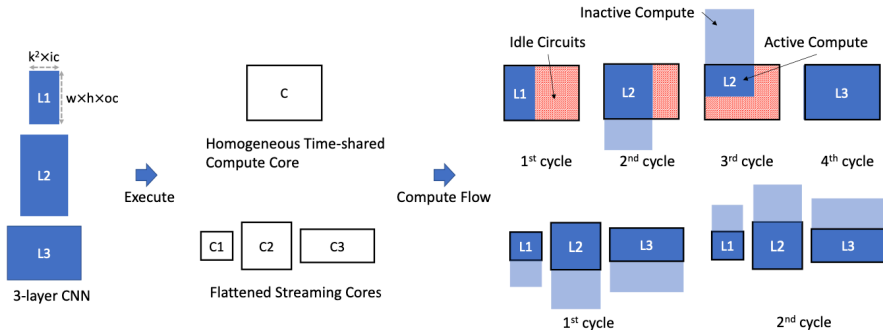
# Dataflow accelerator design: Advantages

- No complex control flow (minimal or no ISA design)
- (Almost) no waste of resources
- (Almost) fixed memory access pattern
- Deeply pipelined



# Dataflow accelerator design: Disadvantages

- Re-program hardware for each new network
- Scalability issues
- If DRAM is utilized, hard to achieve great performance by filling up all pipeline stages



# The compute pattern Simple blocking

```
# Breaking the vector into blocks
for i in range(0, n, block_size):
    # calculate end val considering the last block which can be
    # smaller than block_size
    end_val_i = min(i + block_size, n)
    # Retrieving block of a
    sub_a = a[i:end_val_i]
    # Retrieving corresponding elements from the vector
    sub_b = b[i:end_val_i]
    # multiplication, actual hardware dimension is (block_size, 1)
    result += np.dot(sub_a, sub_b)
```

# The compute pattern

Vector a: N elements

Vector b: N elements



**$N * N$  multipliers in hardware**  
**Finish in 1 clock cycle**

Vector a: N elements

Vector b: N elements

Block size: M

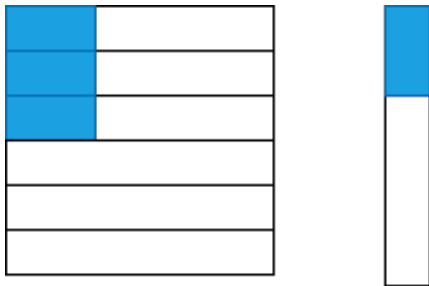


**M multipliers in hardware**  
**Finish in  $N/M$  clock cycle**

## The compute pattern (ii)

- $N \gg M$ , this gives you a chance to do a trade-off between resources and latency simply by changing  $M$ .
- Blocking can happen in a 2D shape!

## 2D Blocking



**M \* M multipliers in hardware**  
**Finish in  $N*N/(M*M)$  clock cycle**

- Parallel Multipliers ( $M^2$ ).
- $M$  Adder Trees ( $\log_2(M)$ ).
- $M$  Accumulators ( $M$ ).