

# **ZigZag-Project: Enabling Fast Architecture-Mapping DSE for Deep Learning Accelerators**

## **ISPASS2023 Tutorial**

Arne Symons, Linyan Mei, Guilherme Paim

Prof. Marian Verhelst

**MICAS Labs, ESAT, KU Leuven, Belgium**



- Introduction (8:30 – 8:50)
- Lab 1: Assess HW performance of DNN layer onto accelerator, with fixed temporal mapping (8:50 – 9:30)
- Lab 2: Automate temporal mapping optimization (9:30 – 10:00)
- Break (10:00 – 10:30)
- Lab 3: Understand the HW architecture definition (10:30 – 11:00)
- Lab 4: Explore layer-fused mappings on multi-core architectures using Stream (11:00 – 11:45)
- Concluding remarks (11:45 – 12:00)

- Project started Dec. 2018
- Many contributors



Prof. Marian Verhelst



Pouya Houshmand



Steven Colleman



Vikram Jain



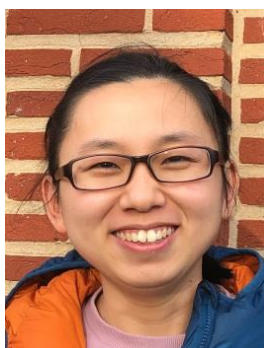
Guilherme Paim



Koen Goetschalckx



Arne Symons



Linyan Mei

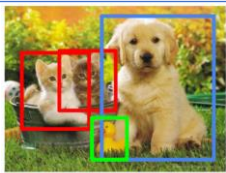


Victor JUNG (ETHz)

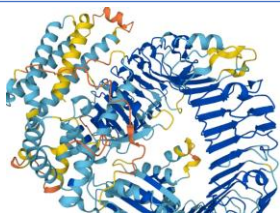


Sebastian Karl (TUM)

## Application

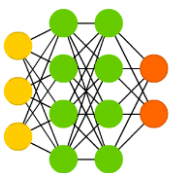


CAT, DOG, DUCK

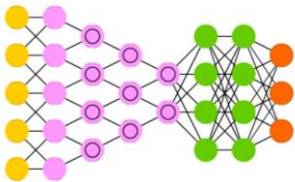


## Algorithm

Deep Feed Forward (DFF)



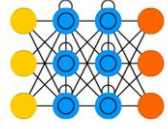
Deep Convolutional Network (DCN)



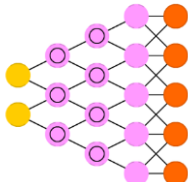
Auto Encoder (AE)



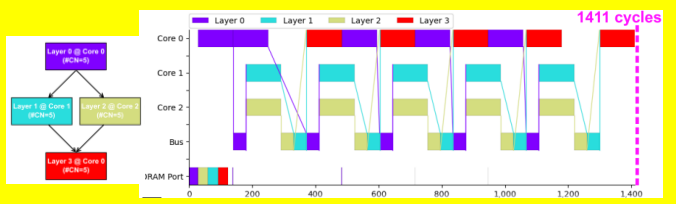
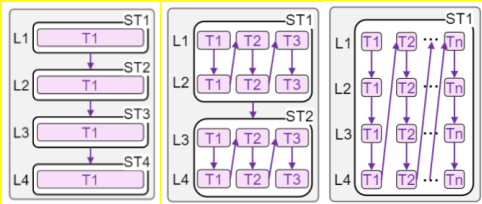
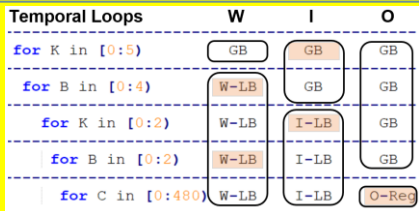
Long / Short Term Memory (LSTM)



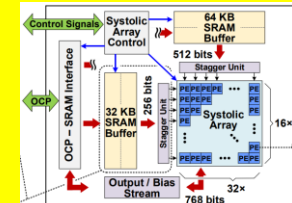
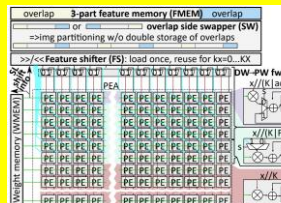
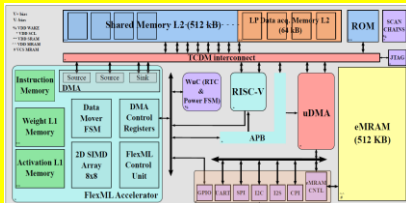
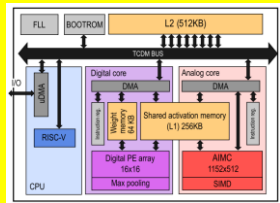
Deconvolutional Network (DN)



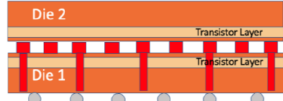
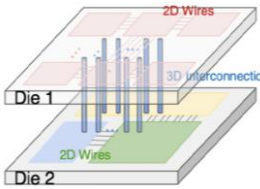
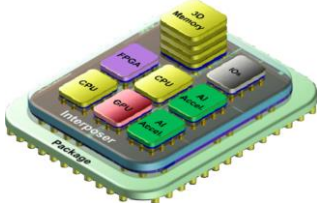
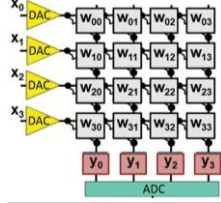
## Mapping/ Scheduling

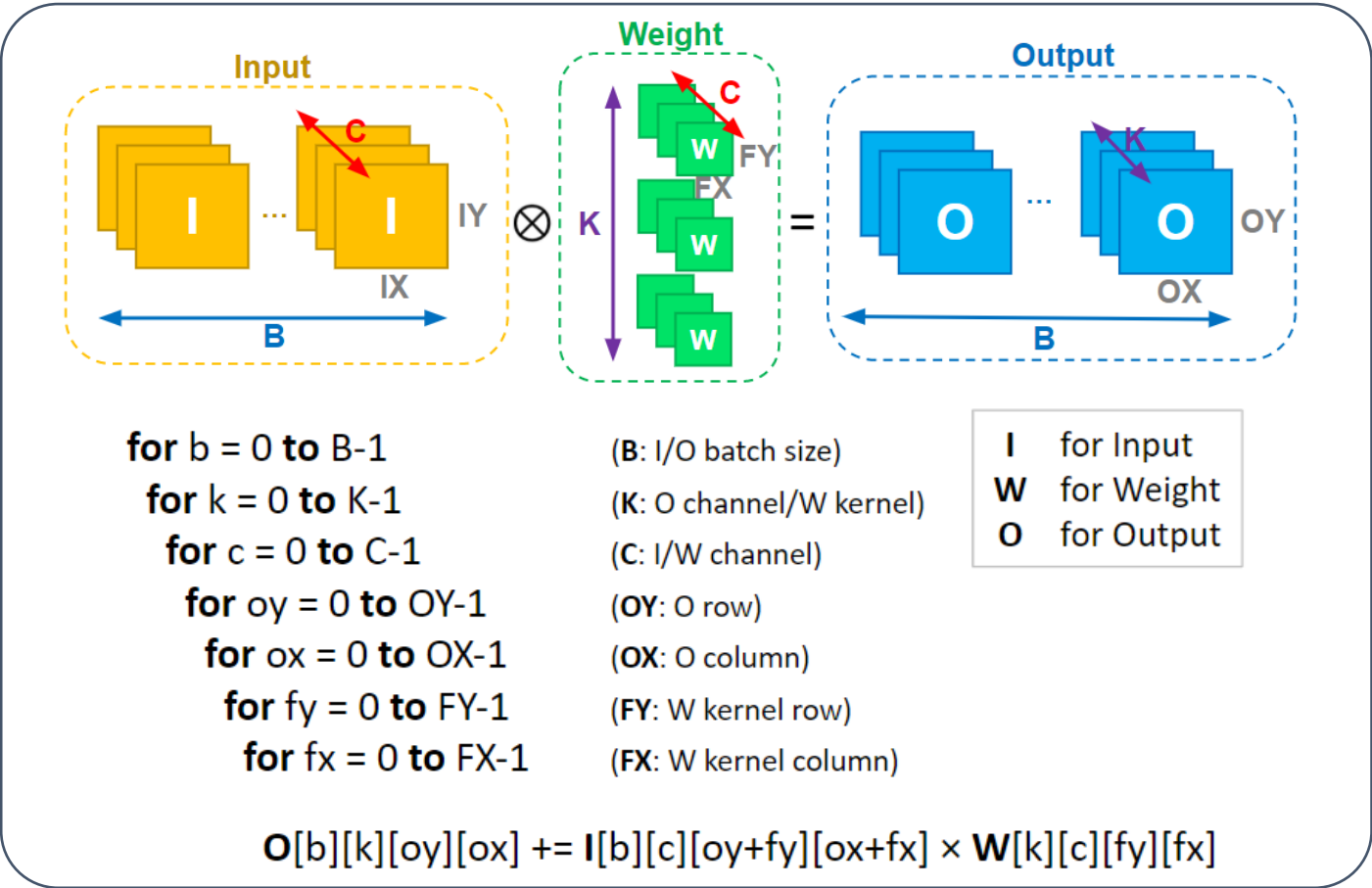


## HW Arch.



## Technology





	B	K	C	OY	OX	FY	FX
W	×	✓	✓	×	×	✓	✓
I	✓	×	✓	? <sup>IY</sup>	? <sup>IX</sup>	? <sup>IY</sup>	? <sup>IX</sup>
O	✓	✓	×	✓	✓	×	×

- ✓ relevant (r)
- × irrelevant (ir)
- ? partially relevant (pr)
- ?<sup>IX/IY</sup> partially relevant to IX/IY

**A DNN Conv2D layer:**

**3D** operand (**W/I/O**) space.

**7D** nested for-loop  
MAC operation.

Each Operand has its own  
(ir)relevant loop dimensions.

- **r** loops contribute to **data size**.
- **ir** loops contribute to **data reuse**.
- **pr** loops contribute to both **data size** and **data reuse**.



Workload	I Batch size	O channel	I / W channel	O row	O column	W row	W column
Conv2D (right fig.)	B	K	C	OY	OX	FY	FX
Conv1D	B	K	C	1	OX	1	FX
Depthwise Conv2D*	B	1	1	OY	OX	FY	FX
Pointwise Conv2D	B	K	C	OY	OX	1	1
Matrix-Vector Multi.	1	K	C	1	1	1	1
Matrix-Matrix Multi.	B	K	C	1	1	1	1

\* Repeat 'C' or 'K' times to finish one Depthwise Conv2D layer (C = K).

Conv2D

for b = 0 to B-1  
  for k = 0 to K-1  
    for c = 0 to C-1  
      for oy = 0 to OY-1  
        for ox = 0 to OX-1  
          for fy = 0 to FY-1  
            for fx = 0 to FX-1  
              O[b][k][oy][ox] += I[b][c][oy+fy][ox+fx] × W[k][c][fy][fx]

(B: I/O batch size)  
(K: O channel/W kernel)  
(C: I/W channel)  
(OY: O row)  
(OX: O column)  
(FY: W kernel row)  
(FX: W kernel column)

I for Input  
W for Weight  
O for Output

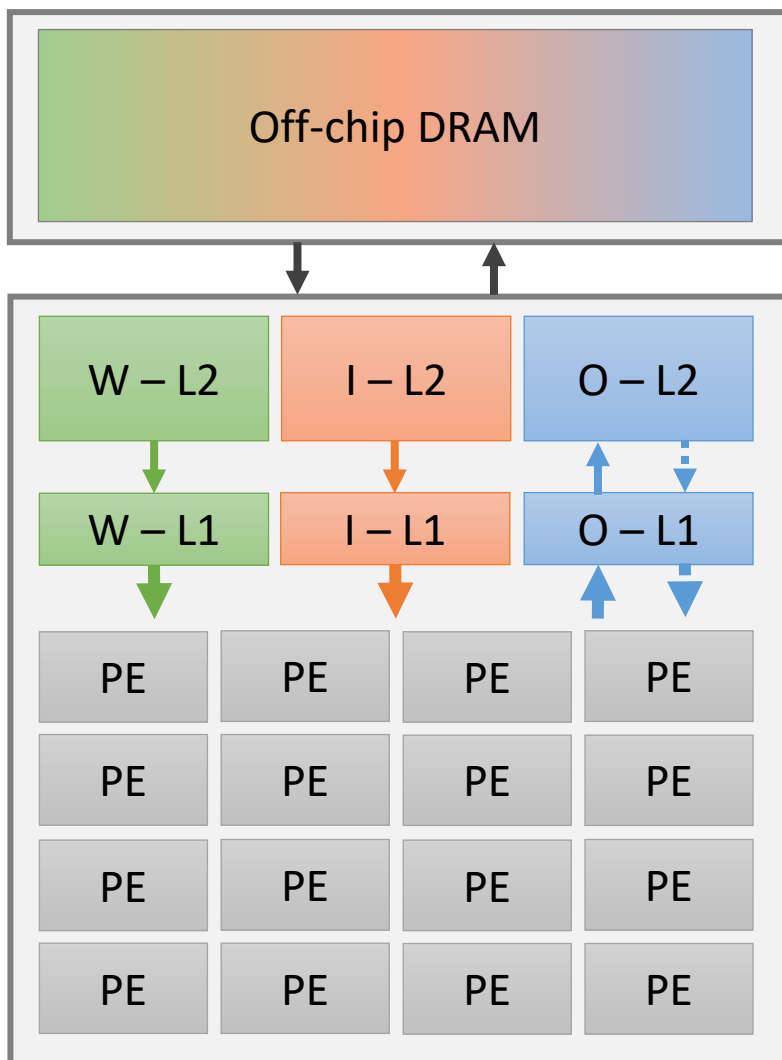
MMM

for b = 0 to B-1  
  for k = 0 to K-1  
    for c = 0 to C-1  
      O[b][k] += I[b][c] × W[k][c]

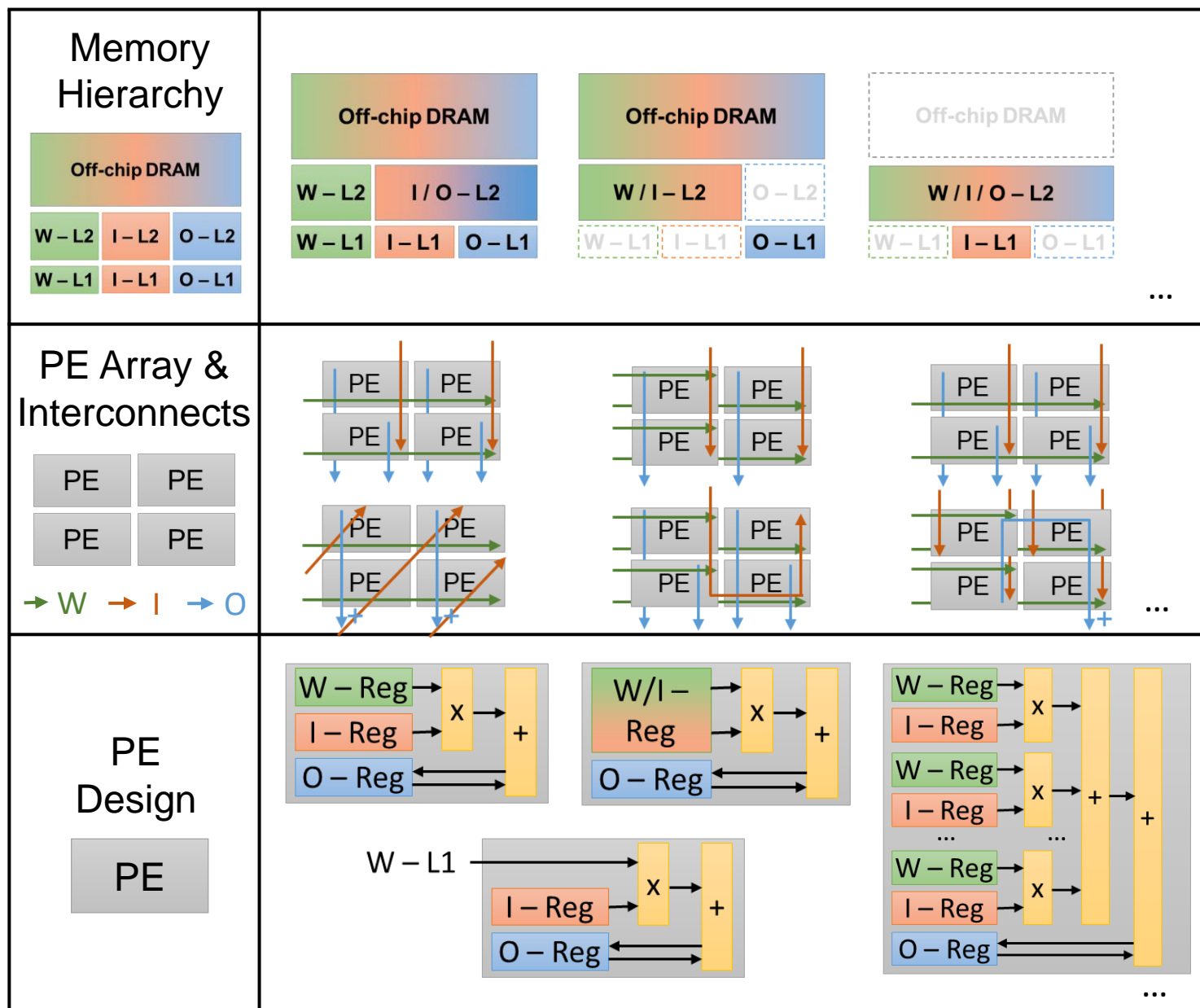
(B: I/O col)  
(K: W/O row)  
(C: W col / I row)

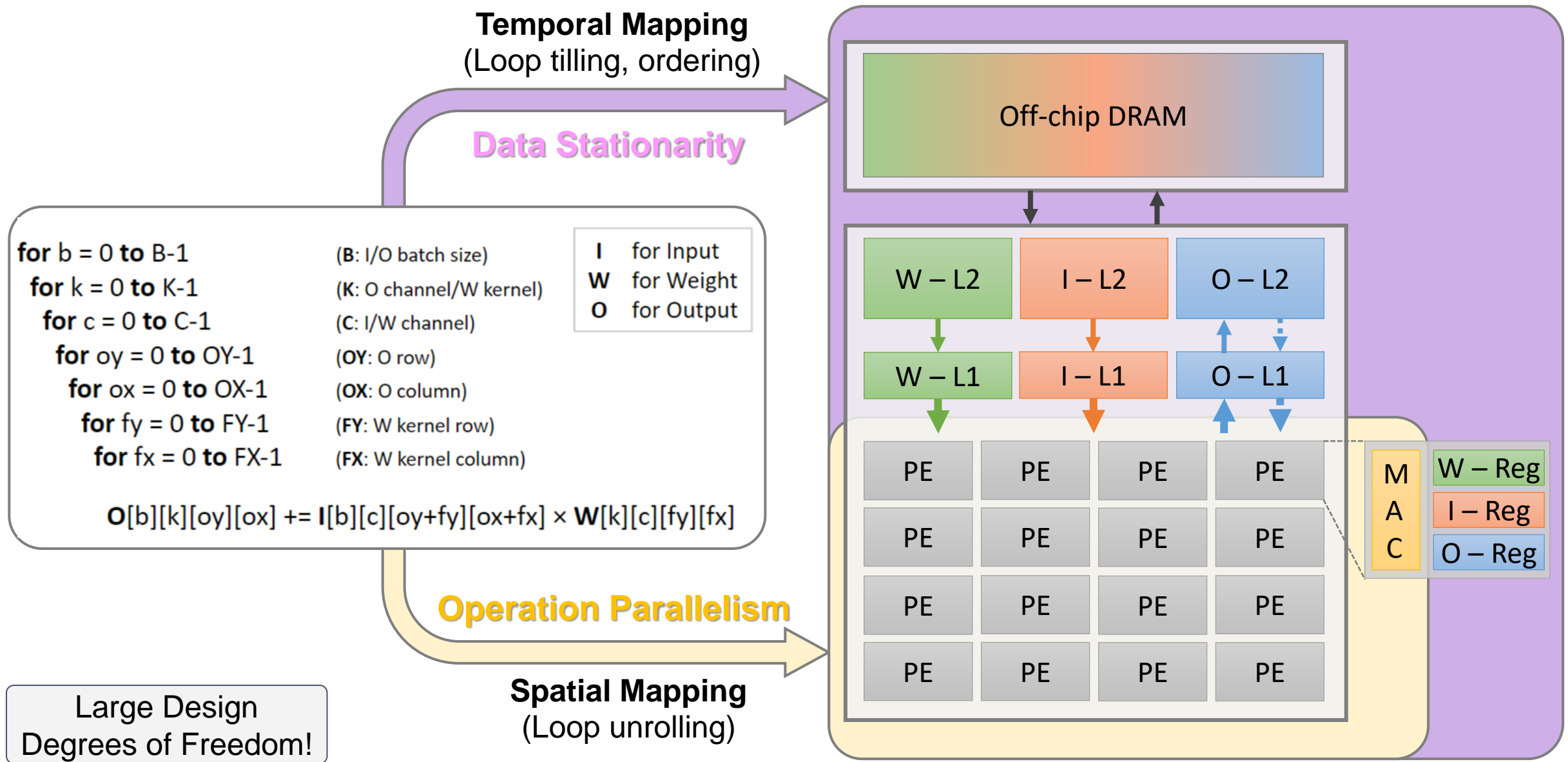
Most **ML workloads** fit into the regular **nested for-loop** format.

**No data dependency** between each for-loop.

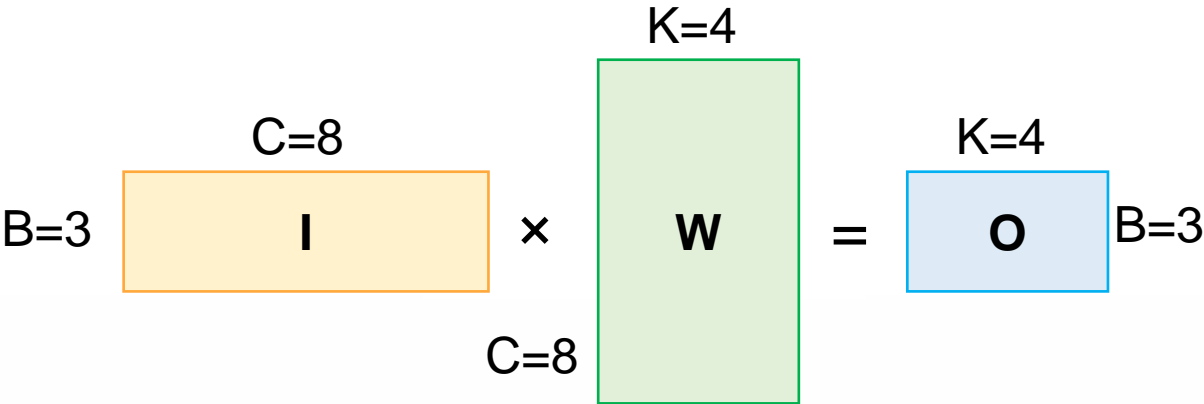


## Large Design Degrees of Freedom!

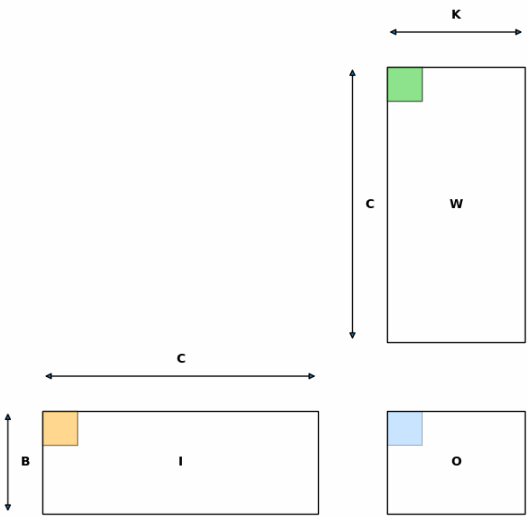




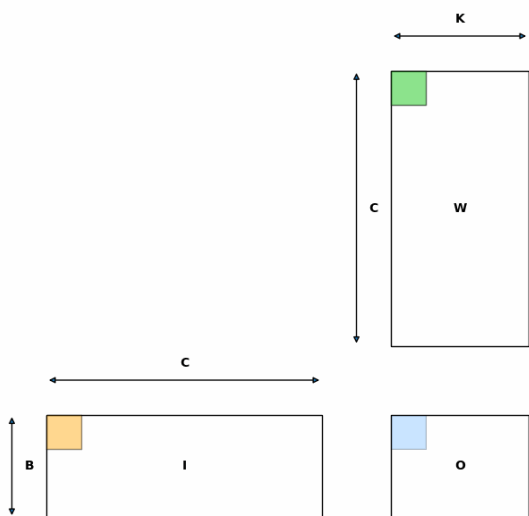




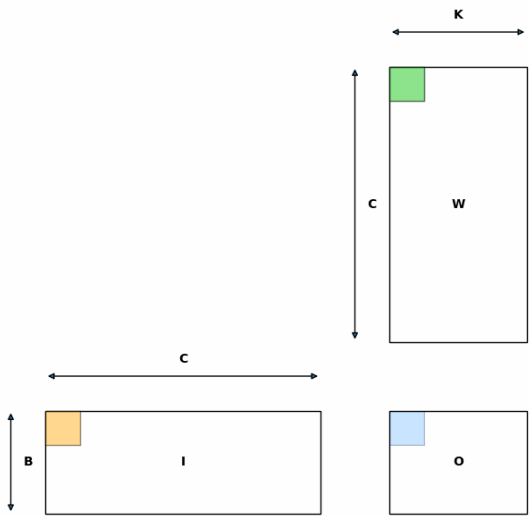
```
for b in [0: 3):  
  for k in [0: 4):  
    for c in [0: 8):  
      O[b][k] = I[b][c] × W[k][c]
```



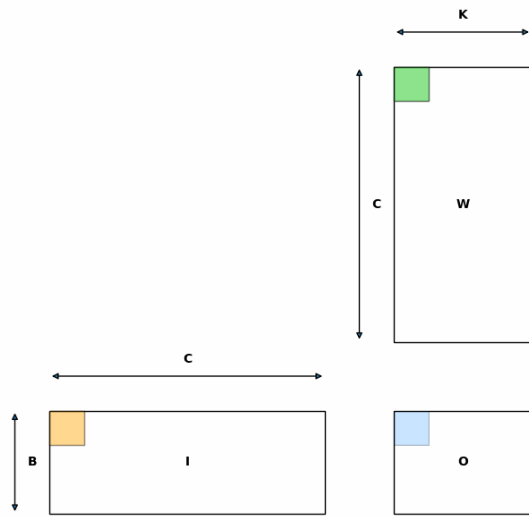
```
for b in [0: 3):  
  for k in [0: 4):  
    for c in [0: 8):  
      O[b][k] = I[b][c] × W[k][c]
```



```
for b in [0: 3):  
  for c in [0: 8):  
    for k in [0: 4):  
      O[b][k] = I[b][c] × W[k][c]
```



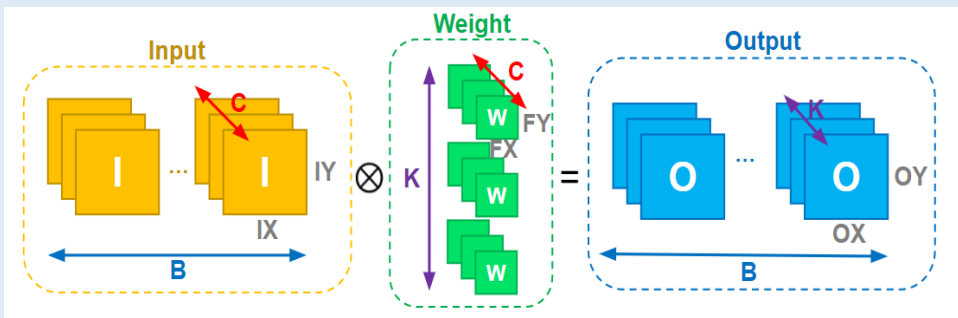
```
for k in [0: 4):  
  for c in [0: 8):  
    for b in [0: 3):  
      O[b][k] = I[b][c] × W[k][c]
```



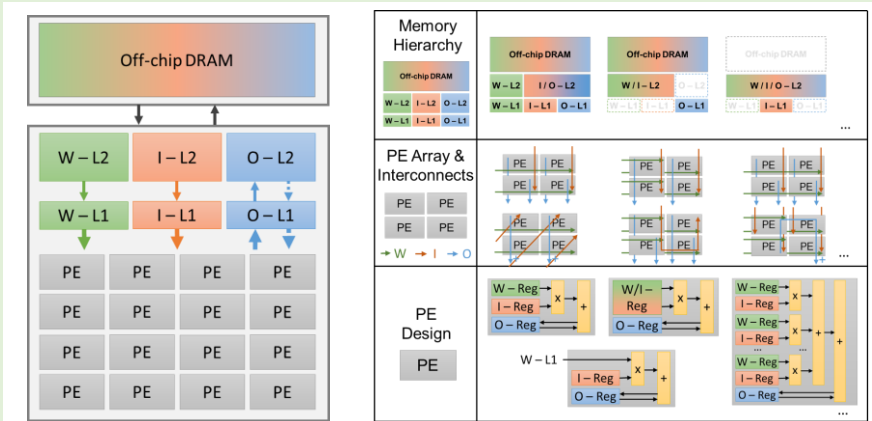
```
for k2 in [0: 2):  
  for c in [0: 8):  
    for b in [0: 3):  
      for k1 in [0: 2):  
        O[b][2k2+k1] = I[b][c] × W[2k2+k1][c]
```

# Co-Exploration

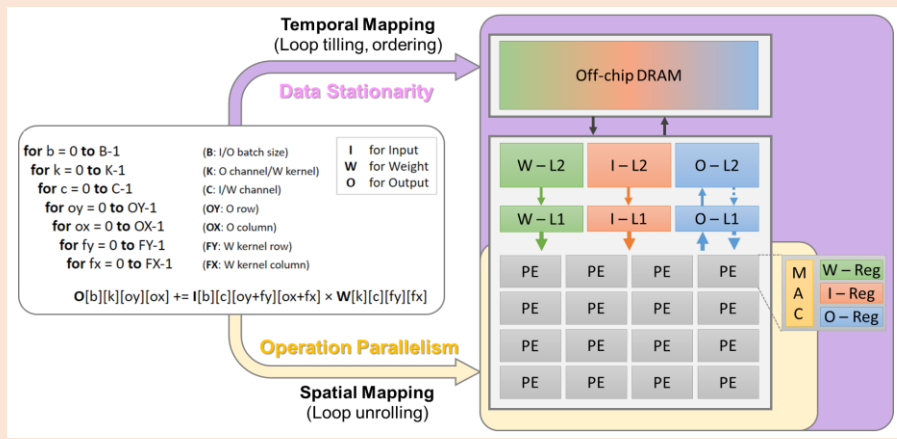
## Algorithm



## Hardware



## Mapping



## Technology and Others

**Technology:** 65nm/40nm/28nm/...,  
NVM, CIM, 3D IC, etc.

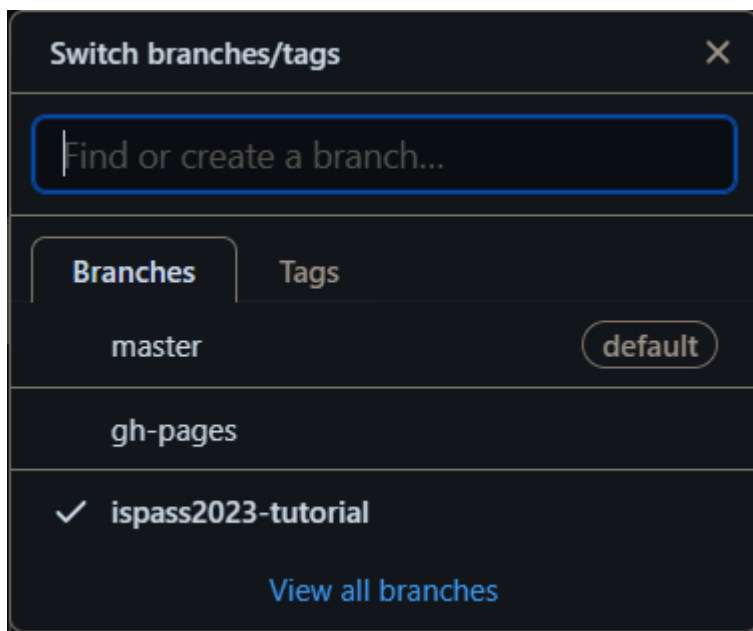
**Others:** Sparsity, various precisions,  
cross-layer execution, etc.

**HUGE** design space at each level & at combined levels.

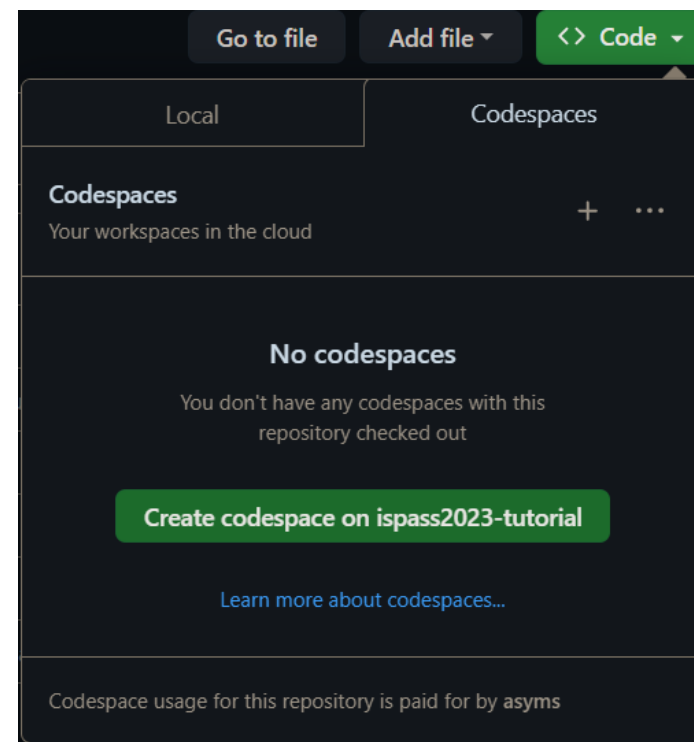
**Regular** workload & **Deterministic** processing flow & **Well-defined** HW components.

<https://github.com/ZigZag-Project/zigzag>

➤ Select ispass2023-tutorial



➤ Create codespace



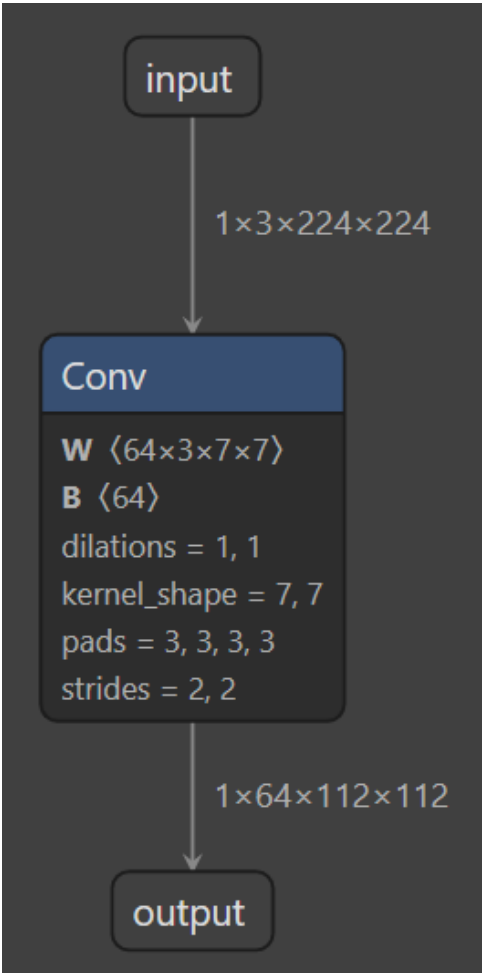
<https://github.com/ZigZag-Project/zigzag>

```
$ git clone git@github.com:ZigZag-Project/zigzag.git
$ cd zigzag
$ conda create --name my-zigzag-env python=3.10
$ conda activate my-zigzag-env
$ pip install -r requirements.txt
$ git checkout ispass2023-tutorial
$ code .
```

- Open lab1/main.py
- Expects three **arguments**:
  - accelerator
  - model
  - mapping
- Extracts names from the given arguments and sets inputs
- Defines the sequence of **stages** to be executed
- Runs the sequence of stages with inputs
- Plots the returned **CostModelEvaluation** (CME)

## Model (workload)

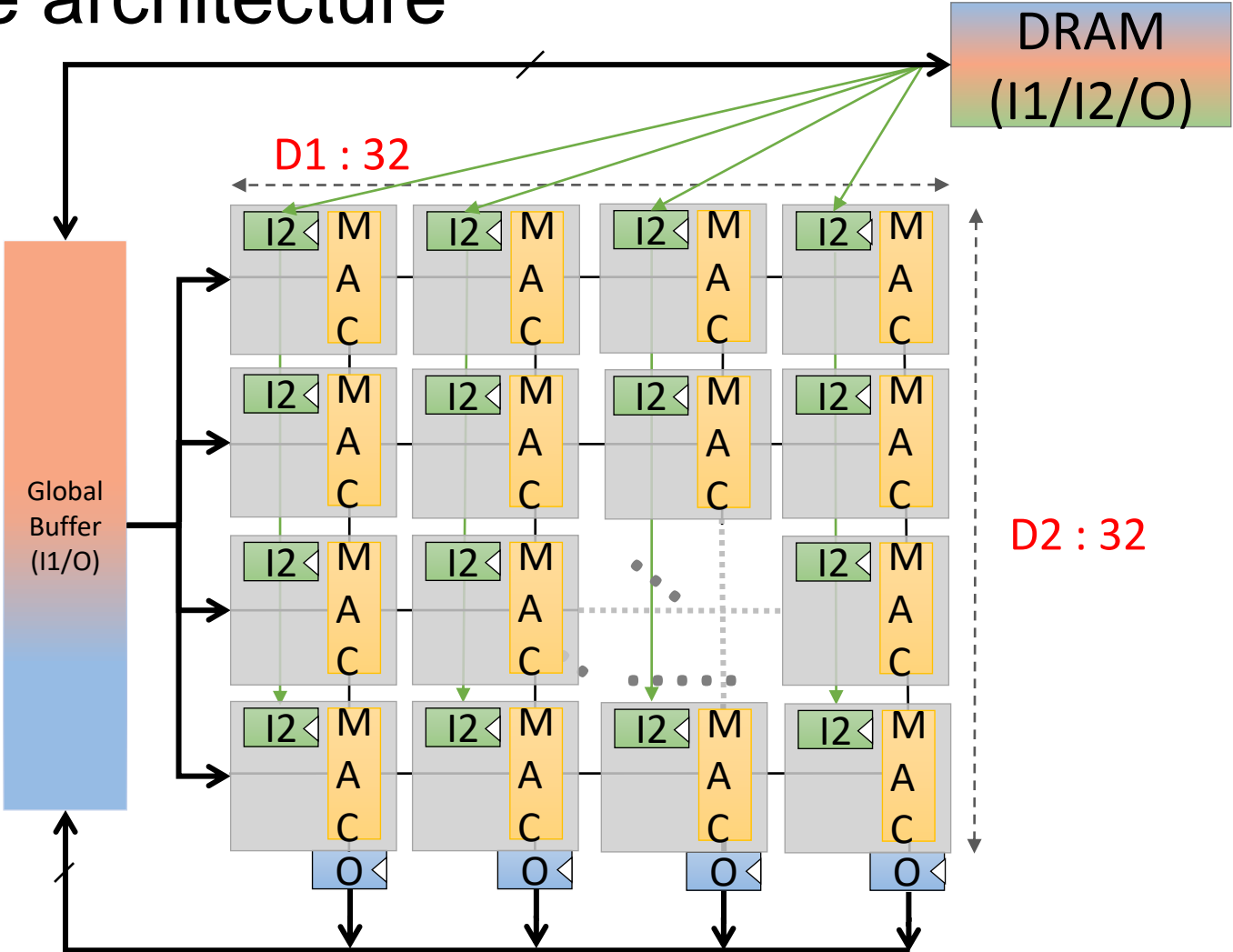
- First layer of ResNet18 (ONNX format)





## Accelerator

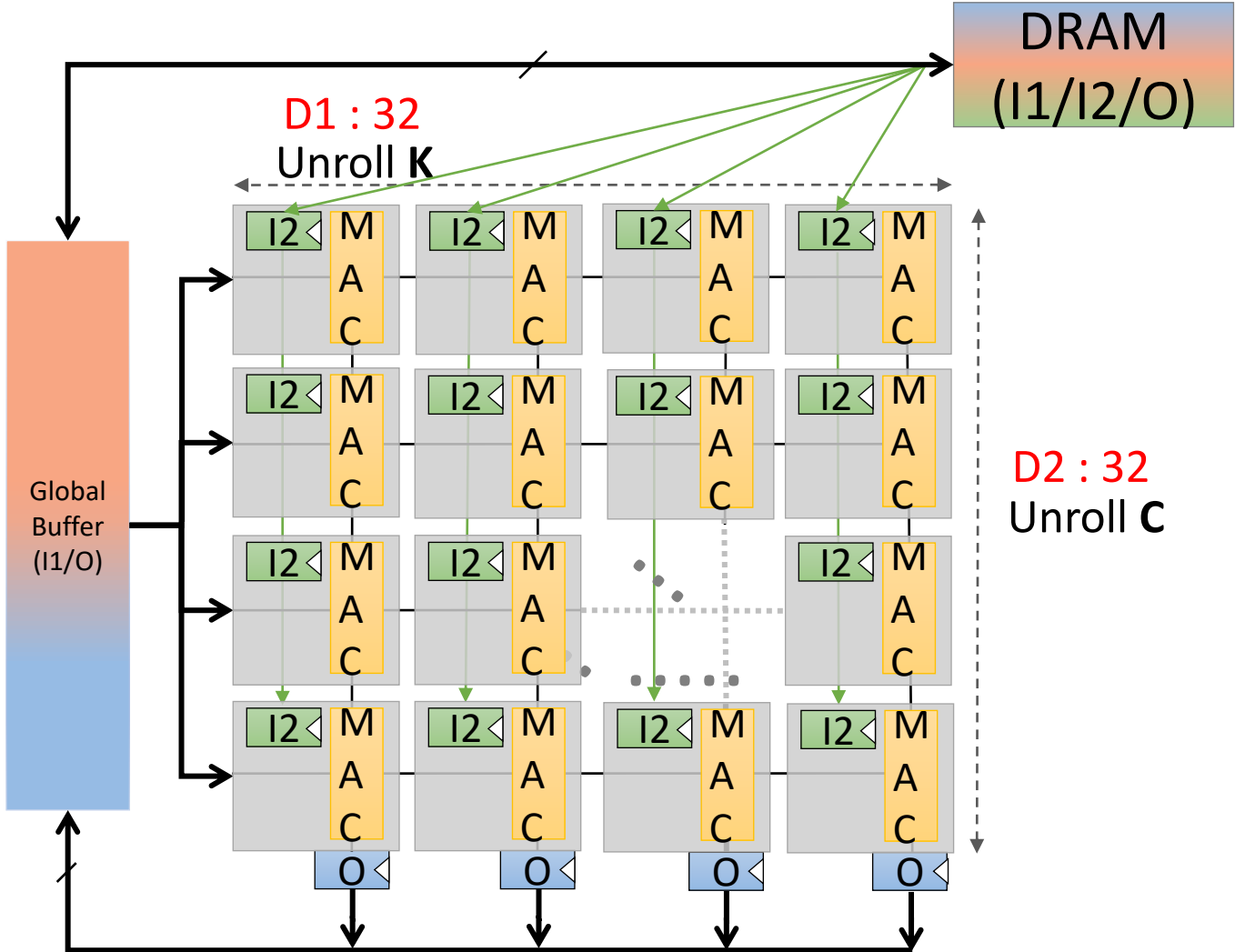
- TPU-like architecture



## Mapping

- Defines mapping of layers onto accelerator

```
mapping = {  
    "/conv1/Conv": { # first ResNet18 layer name in onnx model  
        "spatial_mapping": {"D1": ("K", 32), "D2": ("C", 32)},  
        "temporal_ordering": [  
            # Innermost loop  
            ("OX", 112),  
            ("OY", 112),  
            ("FX", 7),  
            ("FY", 7),  
            ("K", 2),  
            # Outermost loop  
        ],  
        "core_allocation": 1,  
        "memory_operand_links": {  
            "O": "O",  
            "W": "I2",  
            "I": "I1",  
        },  
    },  
}
```

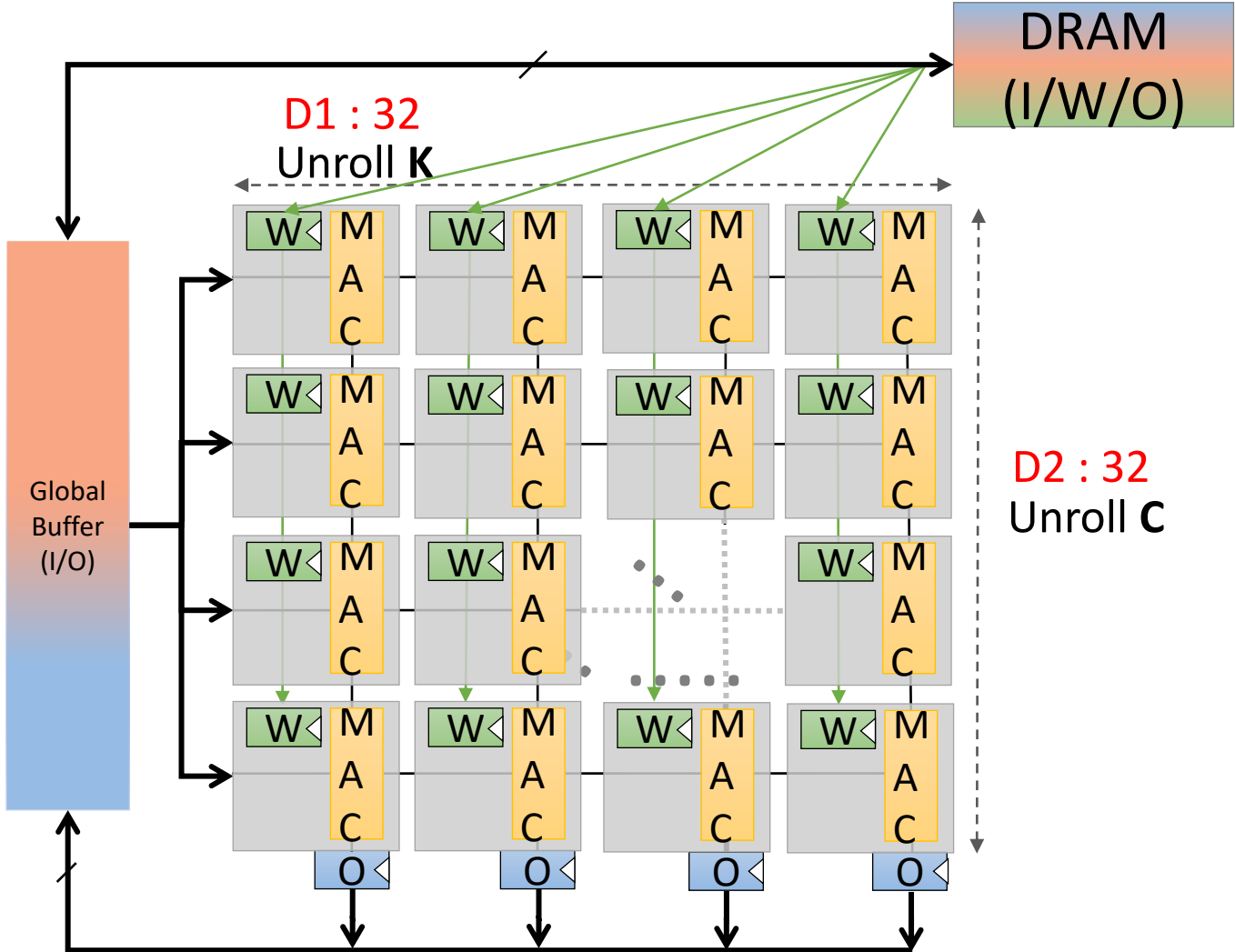


## Mapping

- Defines mapping of layers onto accelerator

```
mapping = {  
    "/conv1/Conv": { # first ResNet18 layer name in onnx model  
        "spatial_mapping": {"D1": ("K", 32), "D2": ("C", 32)},  
        "temporal_ordering": [  
            # Innermost loop  
            ("OX", 112),  
            ("OY", 112),  
            ("FX", 7),  
            ("FY", 7),  
            ("K", 2),  
            # Outermost loop  
        ],  
        "core_allocation": 1,  
        "memory_operand_links": {  
            "O": "O",  
            "W": "I2",  
            "I": "I1",  
        },  
    },  
}
```

$I_1 \rightarrow I$      $I_2 \rightarrow W$      $O \rightarrow O$



## First experiment:

- model = "lab1/resnet18\_first\_layer.onnx"
- accelerator = "zigzag.inputs.examples.hardware.TPU\_like"
- mapping = "mapping"
  
- Run lab1/main.py

```
(my-zigzag-env) asymons@micaszb03:~/zigzag$ python lab1/main.py --model lab1/resnet18_first_layer.onnx --accelerator zigzag.inputs.examples.hardware.TPU_like --mapping mapping
```



## Second experiment:

- Modify lab1/mapping.py:
  - Change temporal loop ordering
  - Run lab1/main.py

```
(my-zigzag-env) asymons@micaszb03:~/zigzag$ python lab1/main.py --model lab1/resnet18_first_layer.onnx --accelerator zigzag.inputs.examples.hardware.TPU_like --mapping mapping
```

- Copy lab1/main.py → lab2/main.py
- Replace TemporalOrderingConversionStage → LomaStage
- Change **dump\_filename\_pattern**
- Change plotting **save\_path**
  
- Copy lab1/mapping.py → lab2/mapping.py
- Remove **temporal\_ordering**

### First experiment:

- model = "lab2/resnet18\_first\_layer.onnx"
- accelerator = "zigzag.inputs.examples.hardware.TPU\_like"
- mapping = "mapping"
  
- Run lab2/main.py

```
(my-zigzag-env) asymons@micaszb03:~/zigzag$ python lab2/main.py --model lab2/resnet18_first_layer.onnx --accelerator zigzag.inputs.examples.hardware.TPU_like --mapping mapping
```

- Copy lab2/main.py → lab2/main\_user\_defined.py
- Replace ONNXModelParserStage → WorkloadParserStage
- Change **dump\_filename\_pattern**
- Change plotting **save\_path**

## Second experiment:

- model = “resnet18\_first\_layer”
- accelerator = “zigzag.inputs.examples.hardware.TPU\_like”
- mapping = “mapping”
  
- Run lab2/main\_user\_defined.py

```
(my-zigzag-env) asymons@micaszb03:~/zigzag$ python lab2/main_user_defined.py --model resnet18_first_layer --accelerat  
or zigzag.inputs.examples.hardware.TPU_like --mapping mapping
```

- ZigZag is also distributed on PyPI

```
~/zigzag$ pip install zigzag-dse
```

- API call for common use-case

```
from zigzag.api import get_hardware_performance_zigzag
def get_hardware_performance_zigzag(
    workload,
    accelerator,
    mapping,
    opt="latency",
    dump_filename_pattern="outputs/{datetime}.json",
    pickle_filename="outputs/list_of_cmes.pickle",
):
```

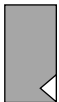


- Lab 3 & 4 after the break
- Start at 10.30 AM

- Open lab3/inputs/hardware/c\_k.py
  - Definition of multiplier array
  - Definition of memory hierarchy
  - Definition of core

memory\_instances

rf\_1B



l1\_w



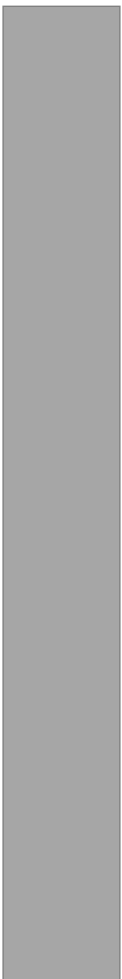
l1\_io



l2\_io



dram



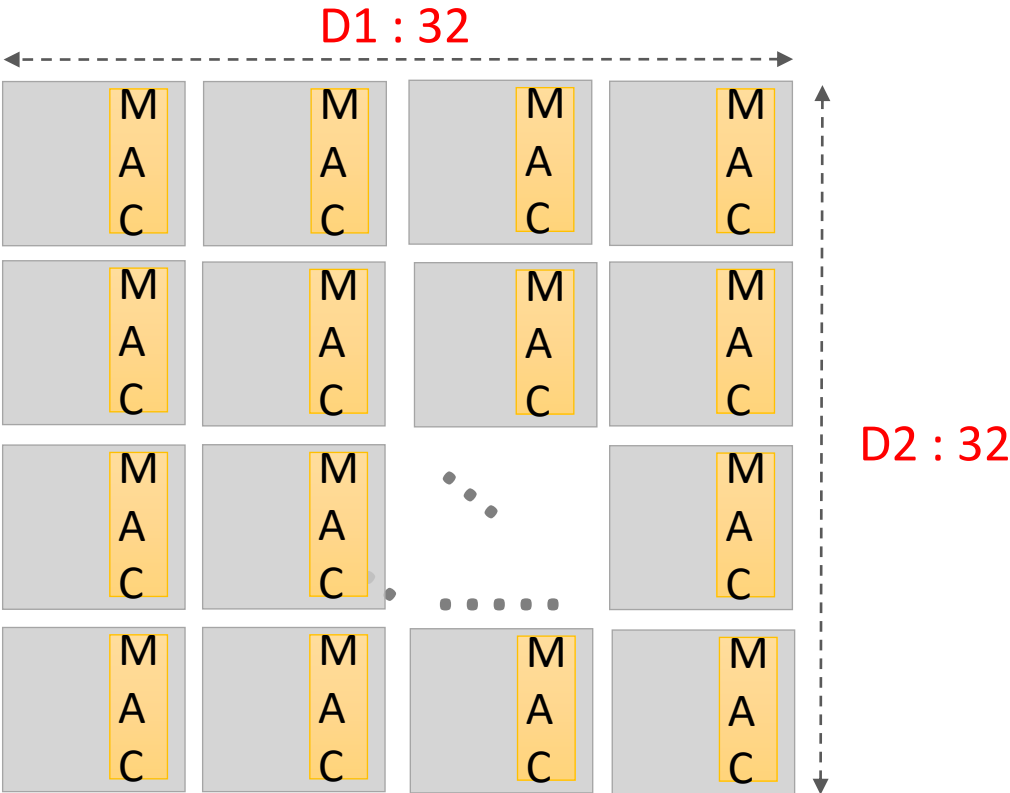
rf\_4B



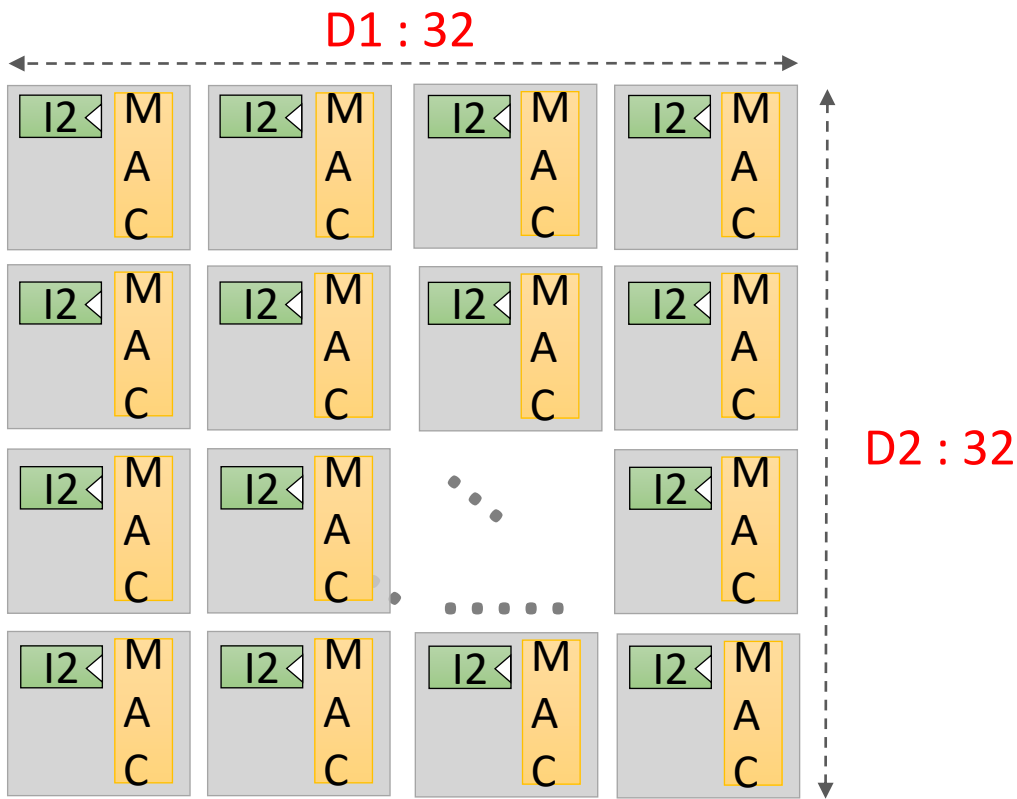
l2\_w



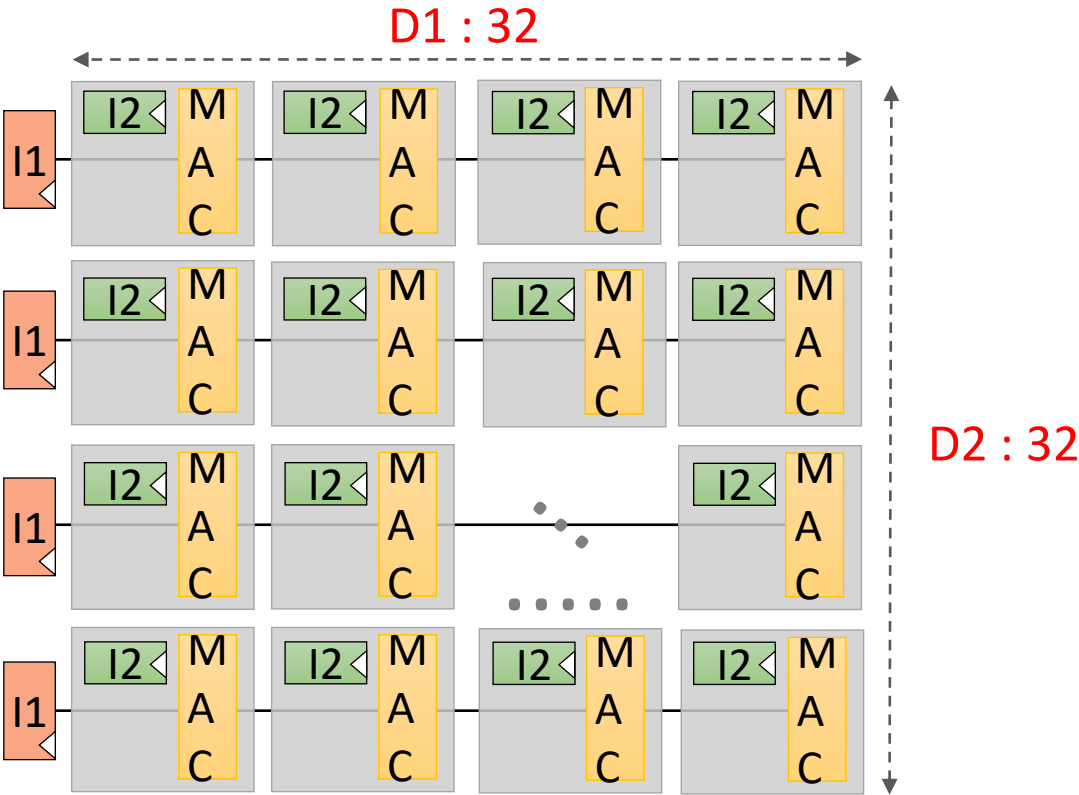
multiplier\_array



multiplier\_array  
rf\_1B

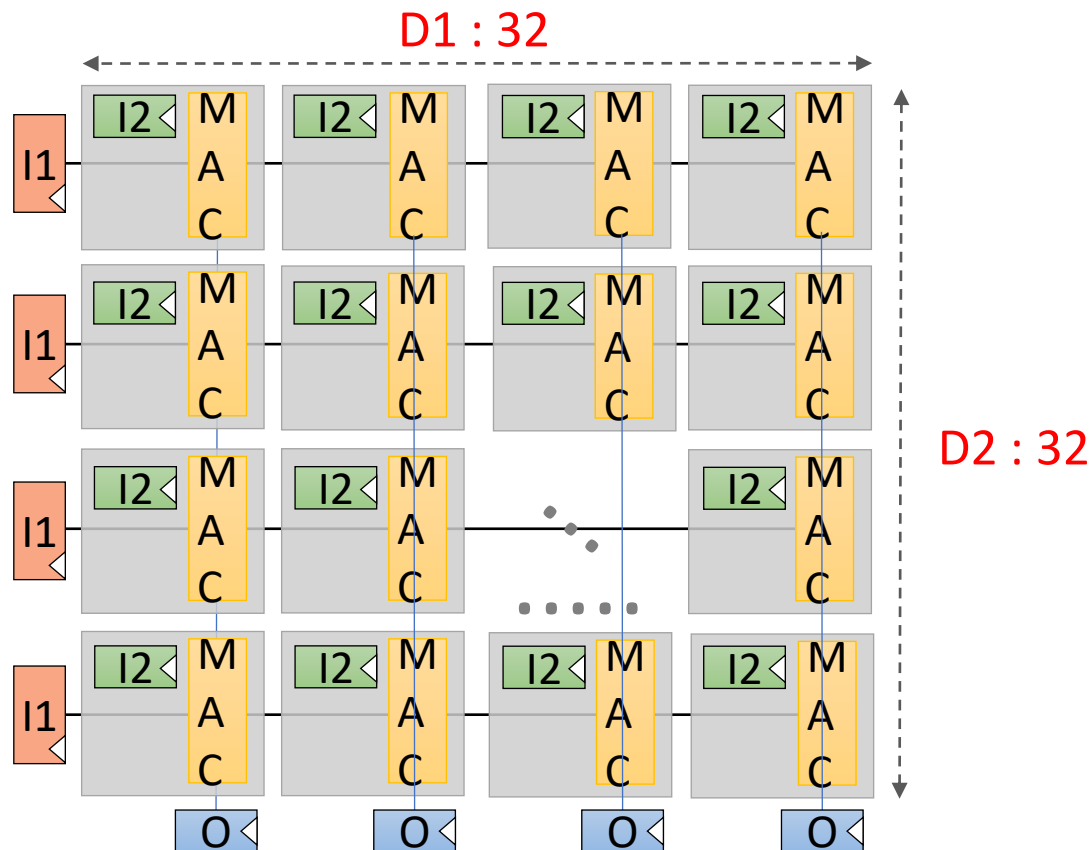


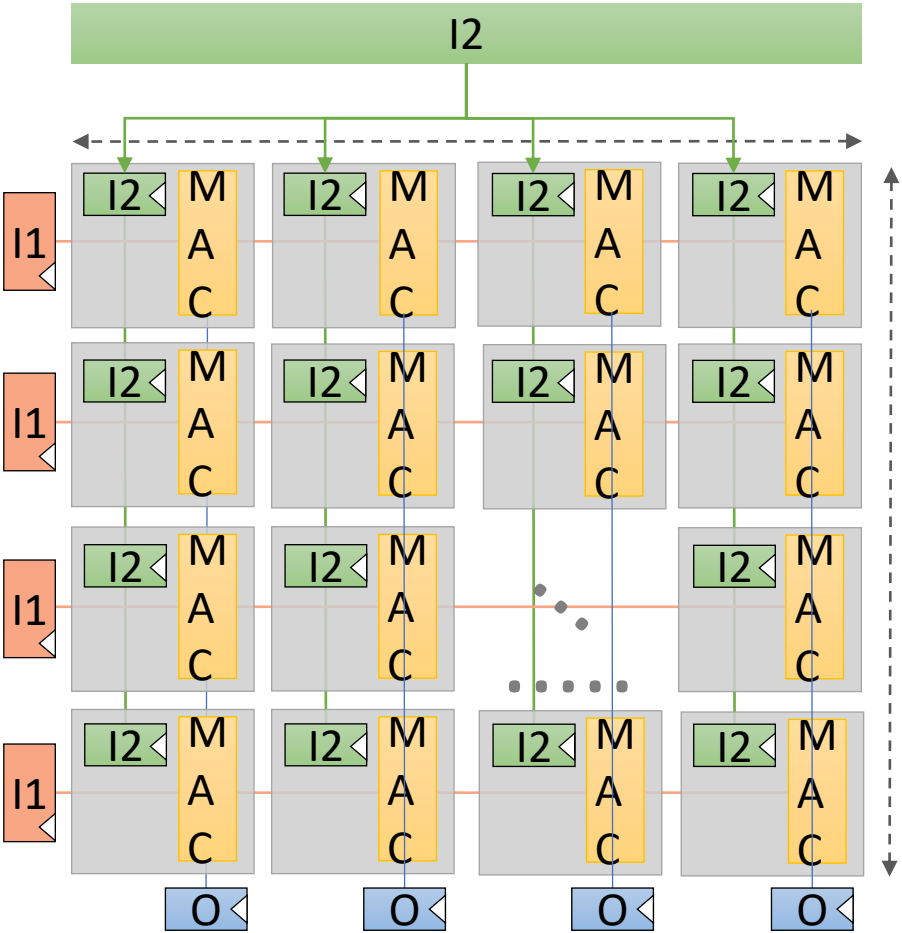
multiplier\_array  
rf\_1B  
rf\_1B





multiplier\_array  
rf\_1B  
rf\_1B  
rf\_4B





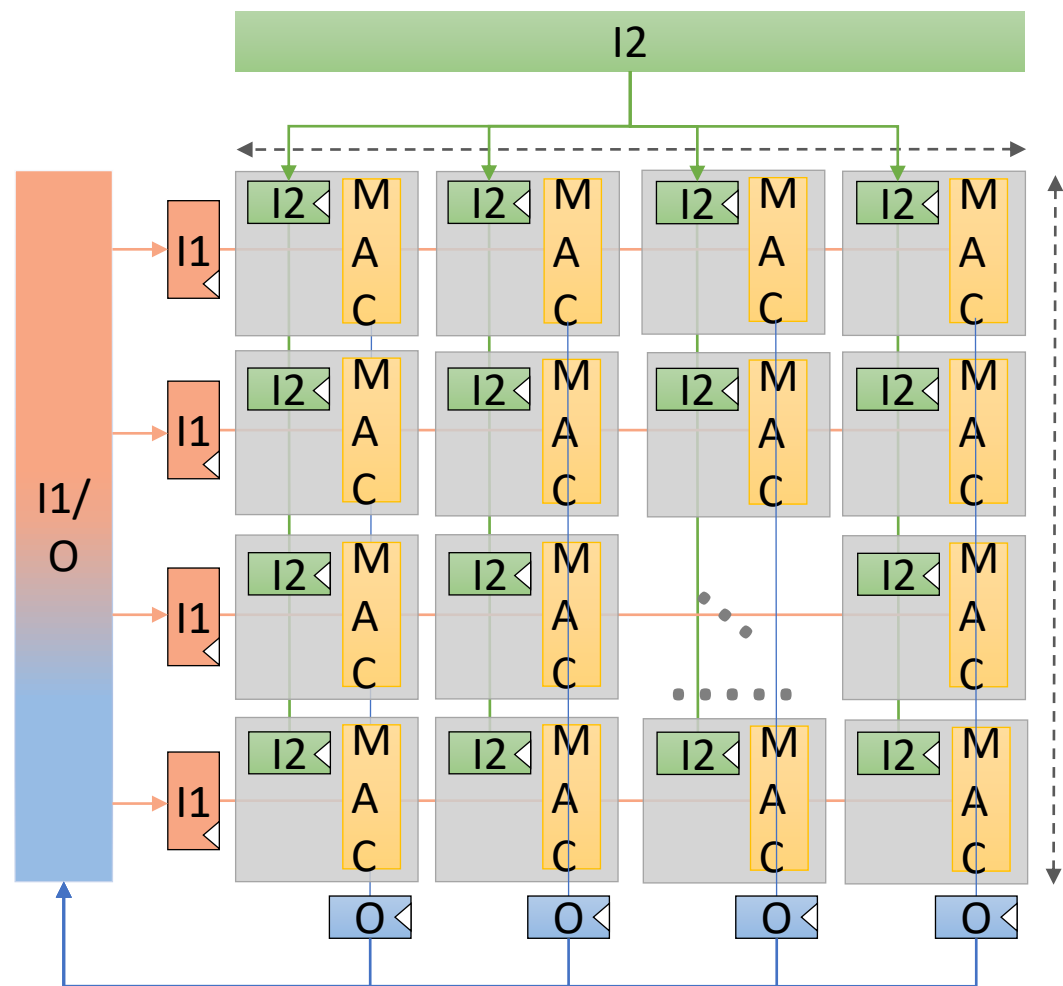
multiplier\_array

rf\_1B

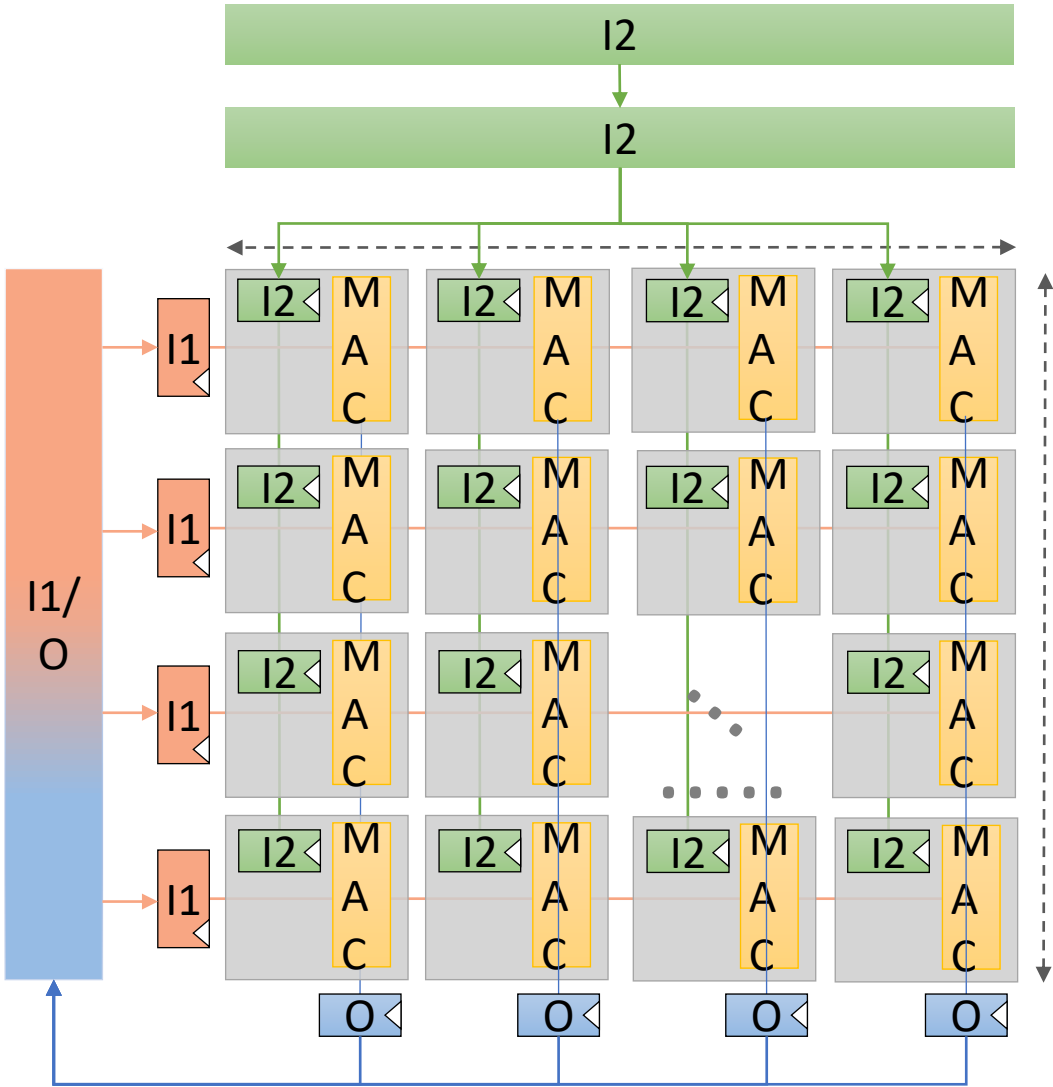
rf\_1B

rf\_4B

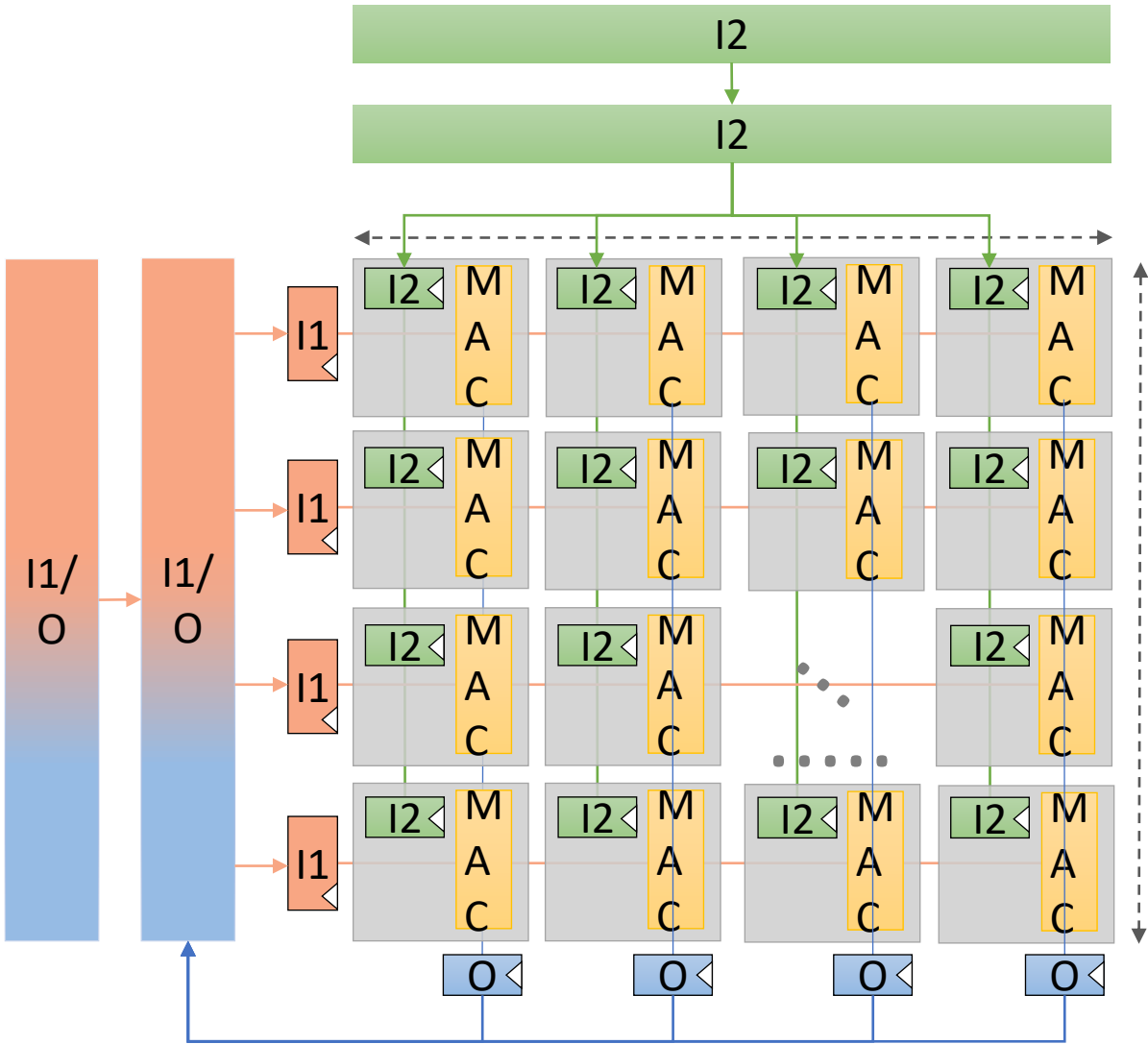
I1\_w



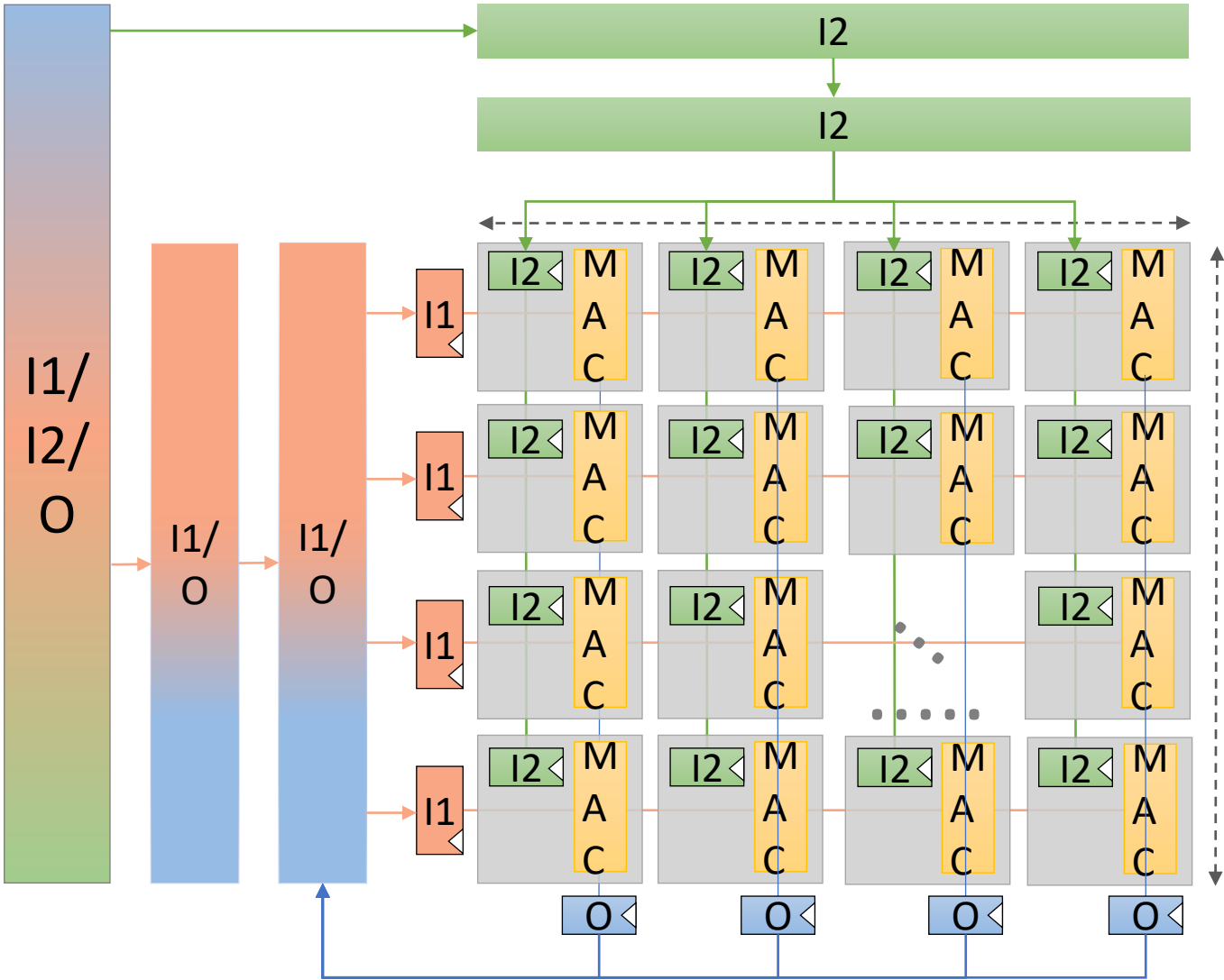
multiplier\_array  
rf\_1B  
rf\_1B  
rf\_4B  
l1\_w  
l1\_io



multiplier\_array  
rf\_1B  
rf\_1B  
rf\_4B  
I1\_w  
I1\_io  
I2\_w



multiplier\_array  
rf\_1B  
rf\_1B  
rf\_4B  
l1\_w  
l1\_io  
l2\_w  
l2\_io



multiplier\_array  
rf\_1B  
rf\_1B  
rf\_4B  
l1\_w  
l1\_io  
l2\_w  
l2\_io  
dram

- Open lab3/inputs/hardware/c\_k.py
  - Definition of multiplier array
  - Definition of memory hierarchy
  - Definition of core
- Open lab3/inputs/mapping/...
  - Definition of spatial mappings

- Open lab3/inputs/hardware/c\_k.py
  - Definition of multiplier array
  - Definition of memory hierarchy
  - Definition of core
  
- Open lab3/inputs/mapping/...
  - Definition of spatial mappings
  
- Open lab3/main.py
  - Uses API call for every core architecture mapping
  - Uses API call for architecture comparison plot



## First experiment:

- Run lab3/main.py

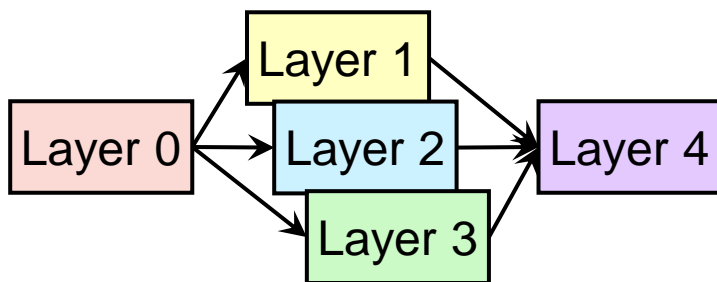
```
(my-zigzag-env) asymons@micaszb03:~/zigzag$ python lab3/main.py
```

- Hardware accelerator model based on array of multipliers and attached memory hierarchy
- Hardware performance estimation of DNN layer through analytical cost model
- Optimization of layer mapping through different stages
- Enables co-exploration of accelerator & mapping

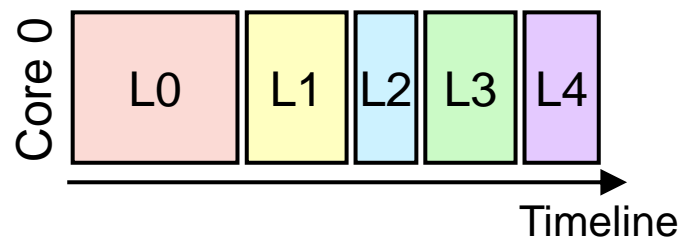
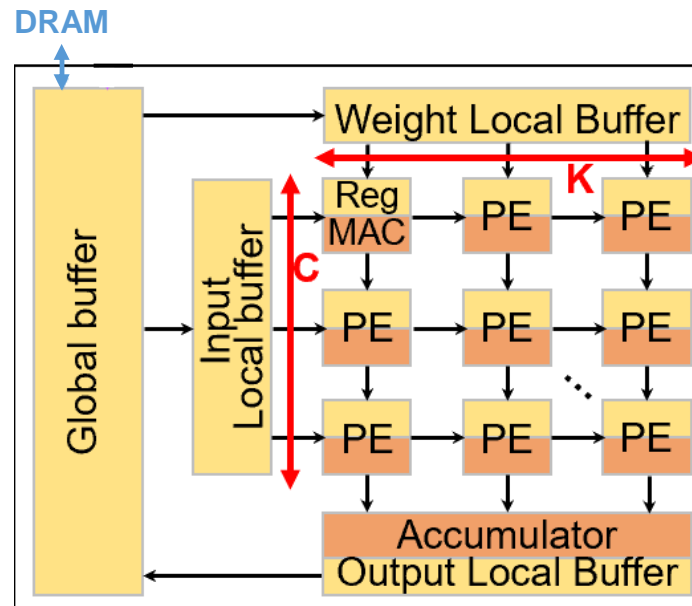
Frameworks	Workload	Hardware	Mapping
ZigZag	A NN layer	Single-core accelerator	Single-layer mapping
Stream	One/more NN(s)	Multi-core accelerator	Fine-grained layer-fused mapping

Focus on **Stream** for the rest of the session

## An example workload

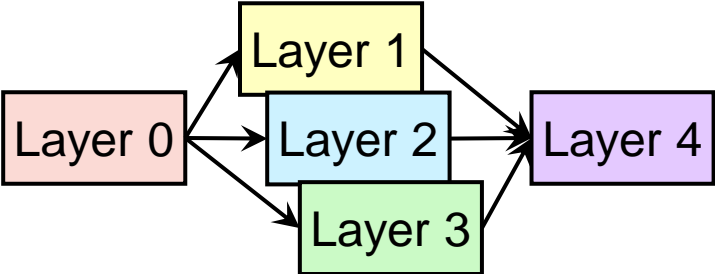


## An example accelerator

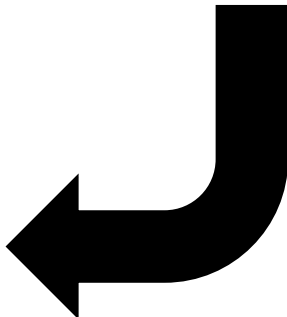
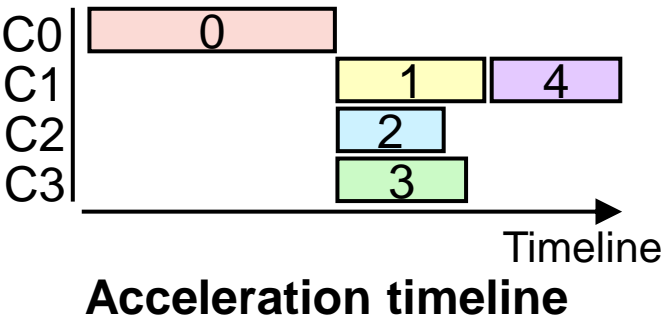
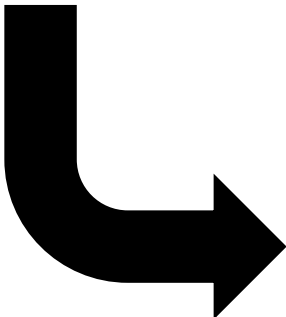
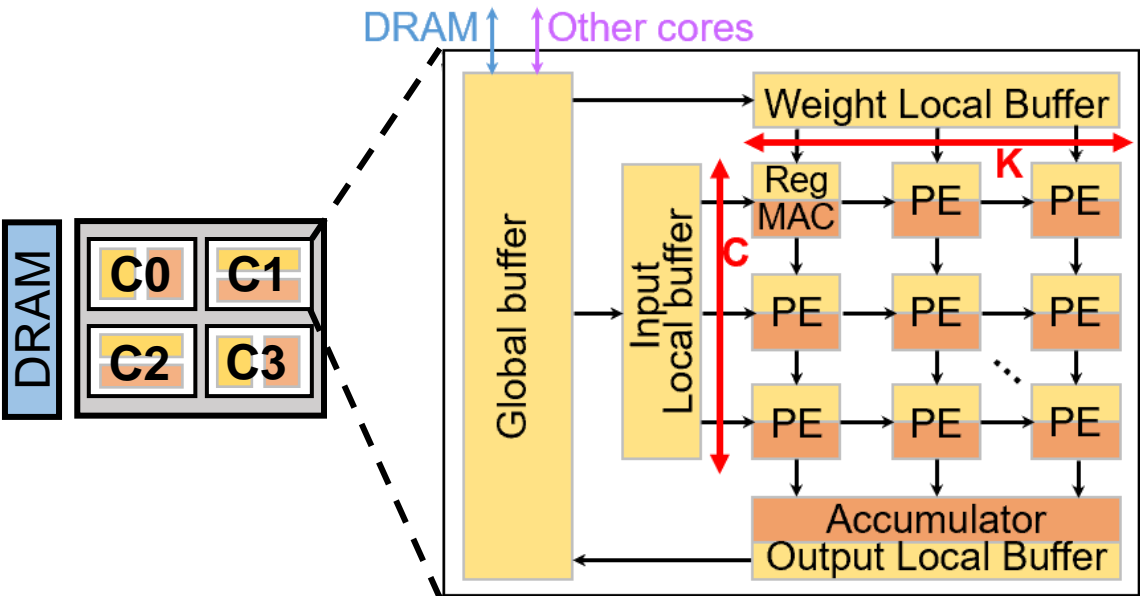


## Acceleration timeline

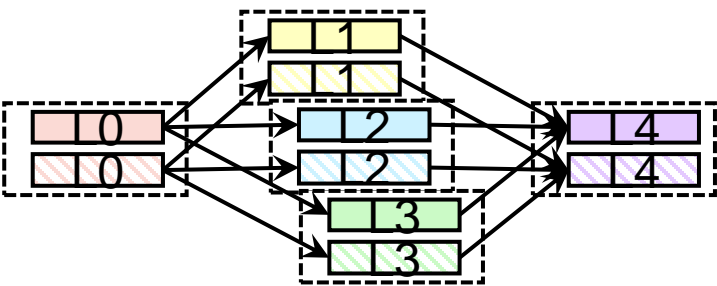
An example workload



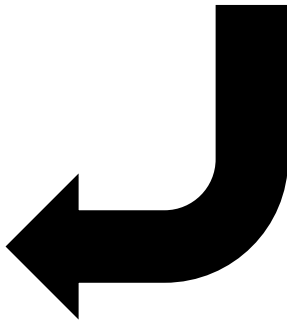
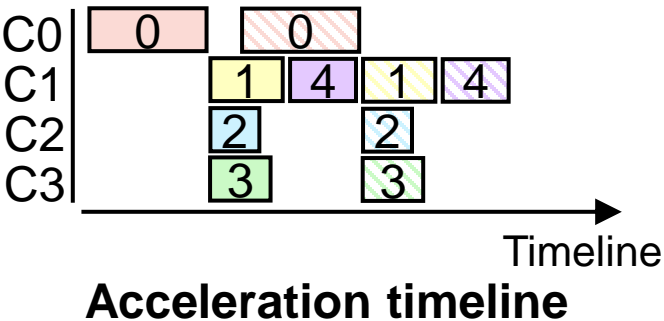
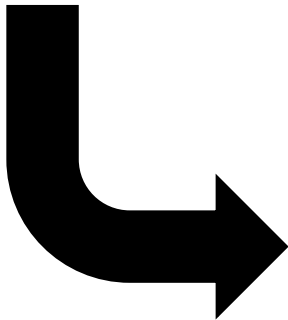
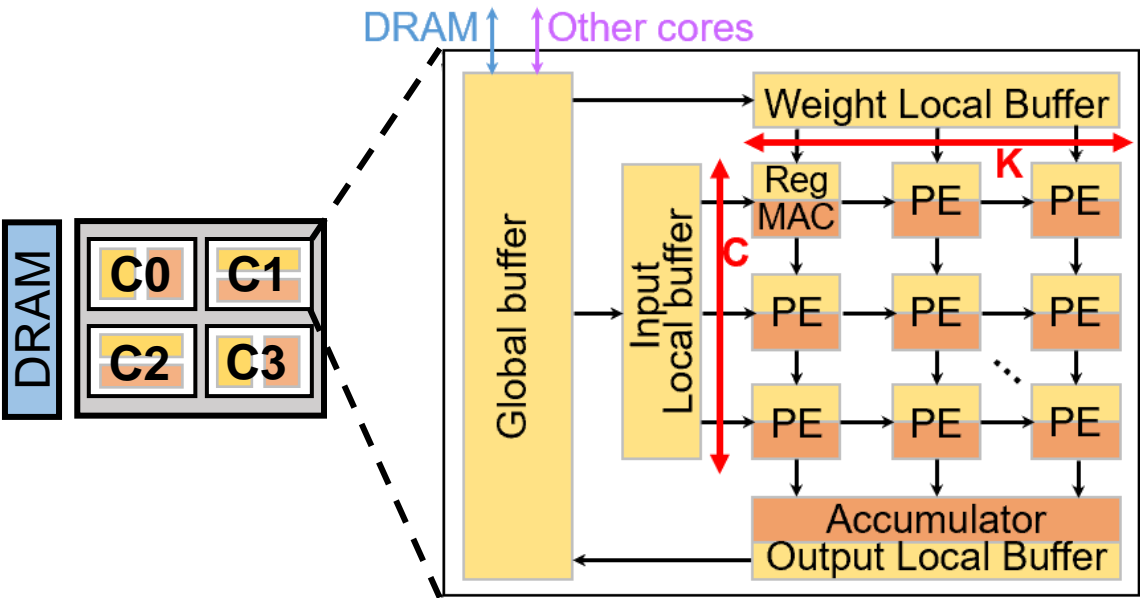
A multi-core accelerator



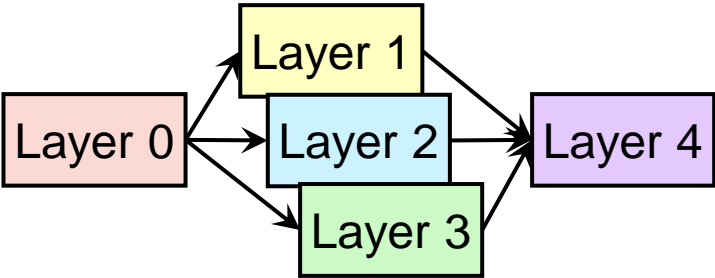
Tiled for layer-fused processing



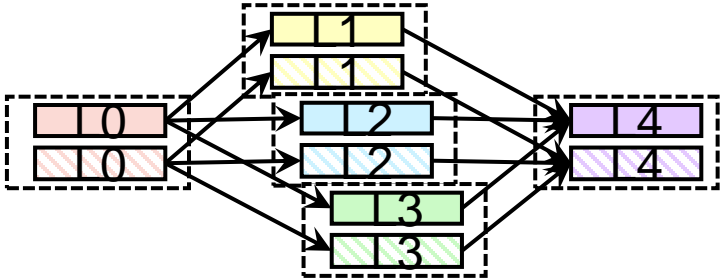
A multi-core accelerator



An example workload



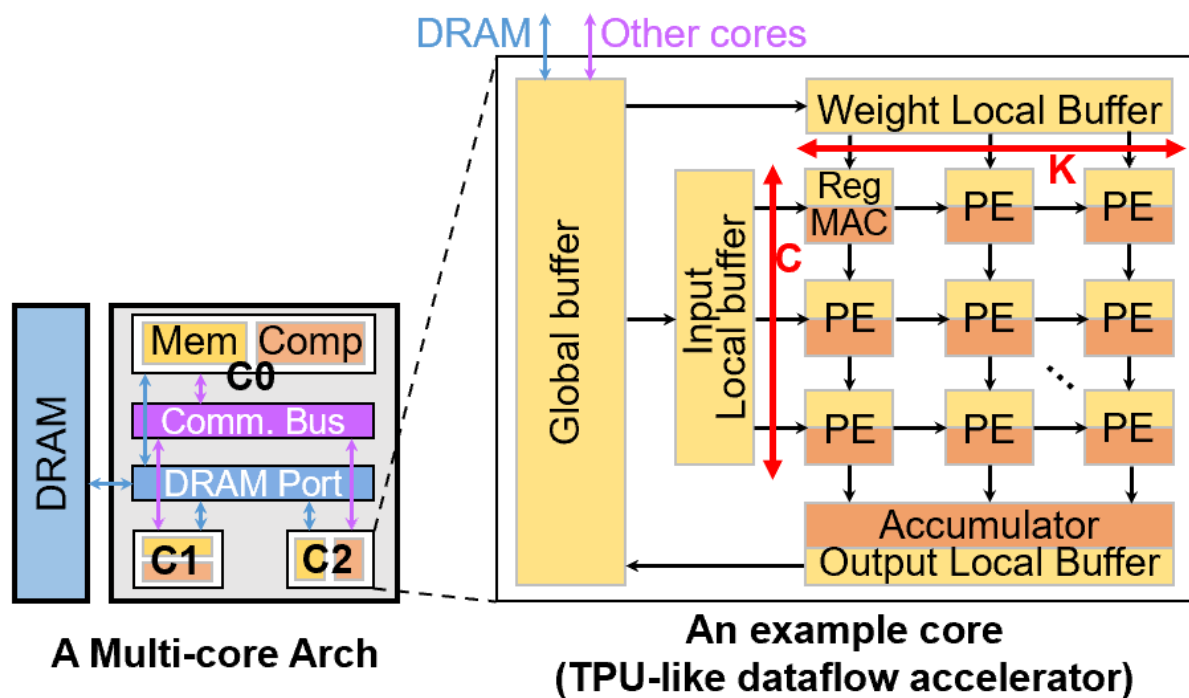
Tiled for layer-fused scheduling



Schedule the workload to hardware accelerators

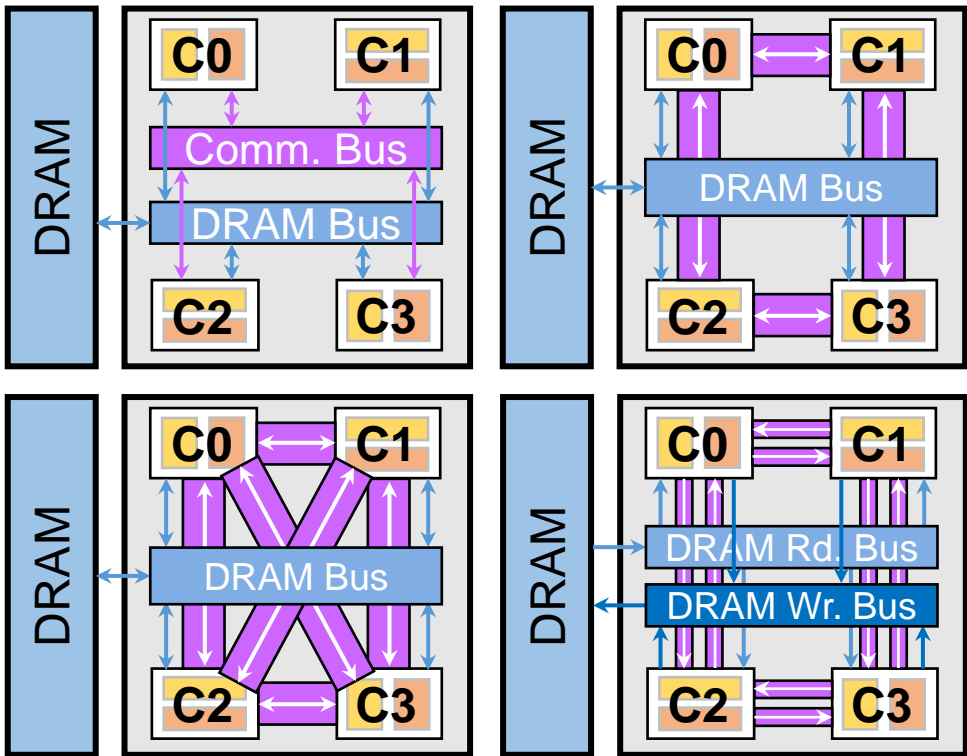
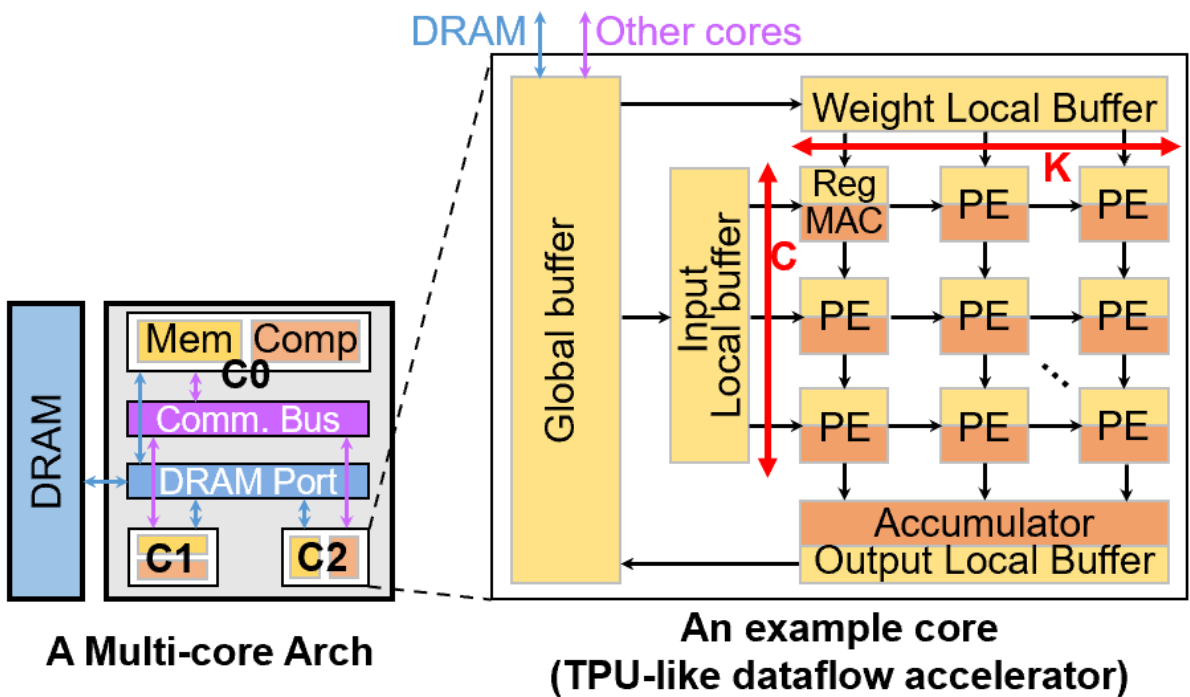
Hardware Schedule	Single-core	Multi-core
Layer-by-layer	(1) Core 0: L0 → L1 → L2 → L3 → L4	(3) Timeline for C0, C1, C2, C3: C0: L0 C1: L1 → L4 C2: L2 C3: L3
SotA frameworks	ZigZag, Timeloop, ..., Stream	Kwon et al. 2021, Stream
Layer-fused	(2) Core 0: 0 1 2 3 4 0 1 2 3 4	(4) Timeline for C0, C1, C2, C3: C0: 0 0 C1: 1 4 1 4 C2: 2 2 C3: 3 3
SotA frameworks	DeFiNES, TVM-Cascade, Stream	Stream

- ✓ Model single-core architecture (identical to ZigZag)

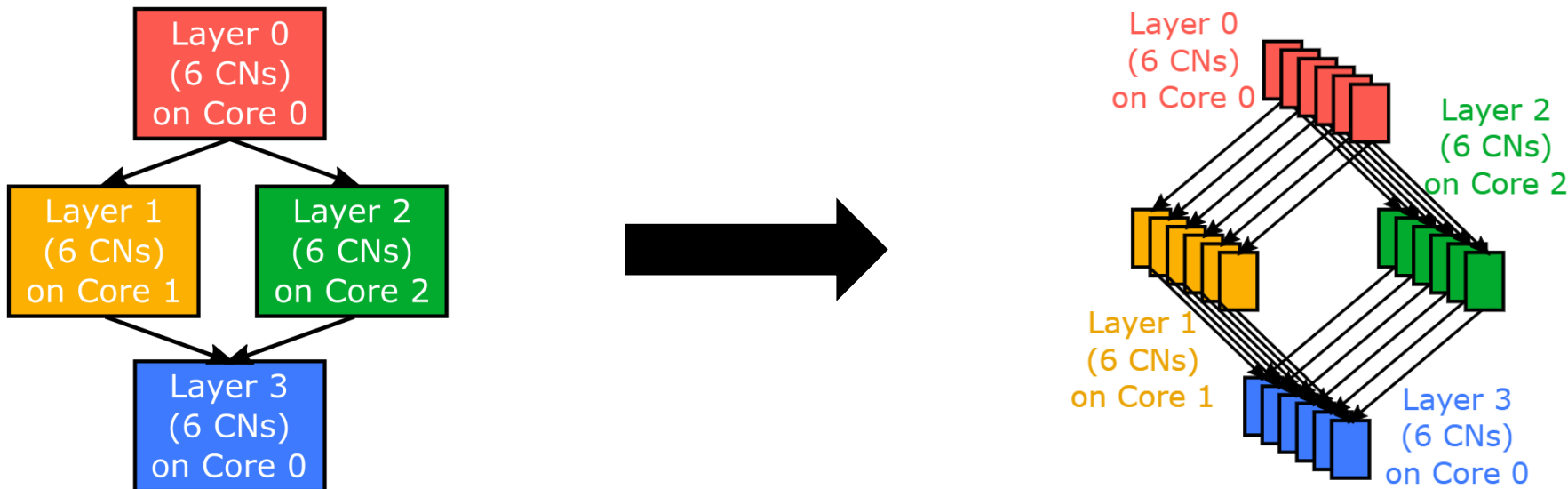




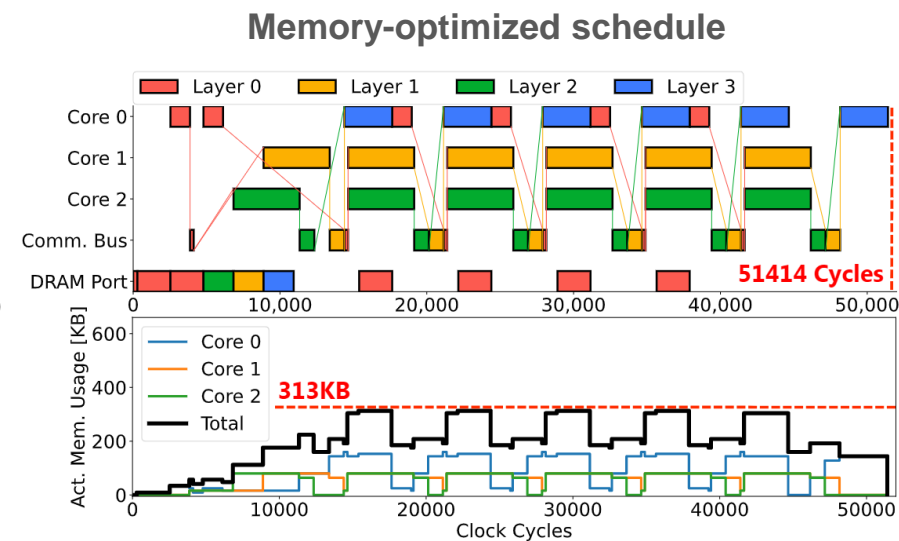
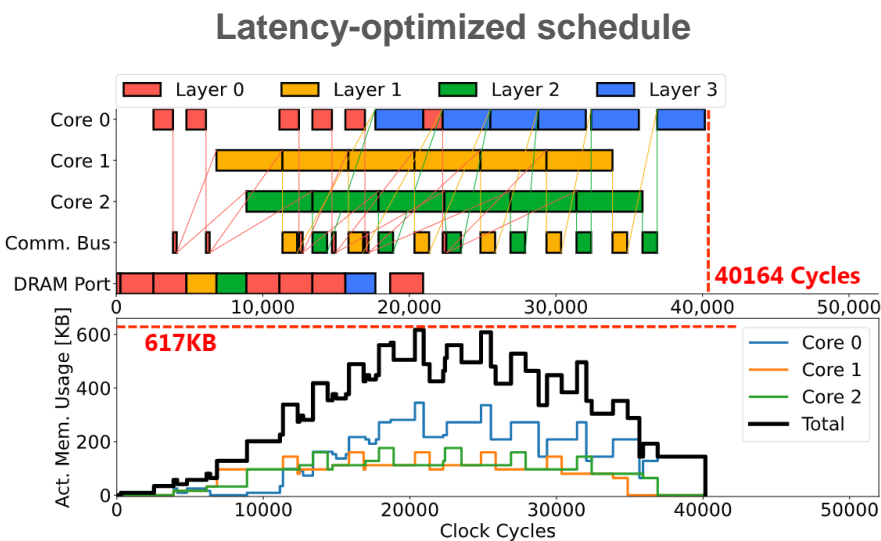
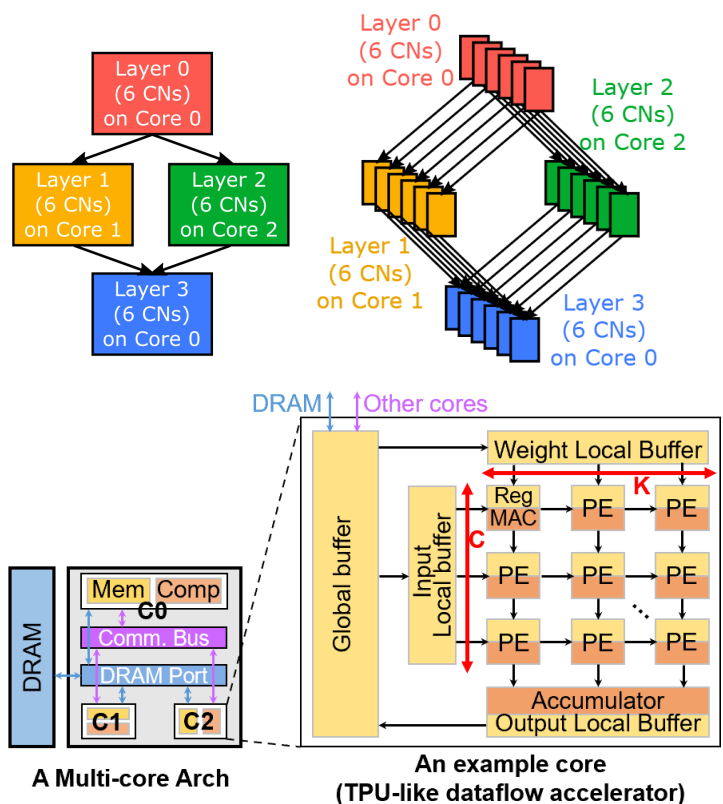
- ✓ Model single-core architecture (identical to ZigZag)
- ✓ Model different multi-core topologies



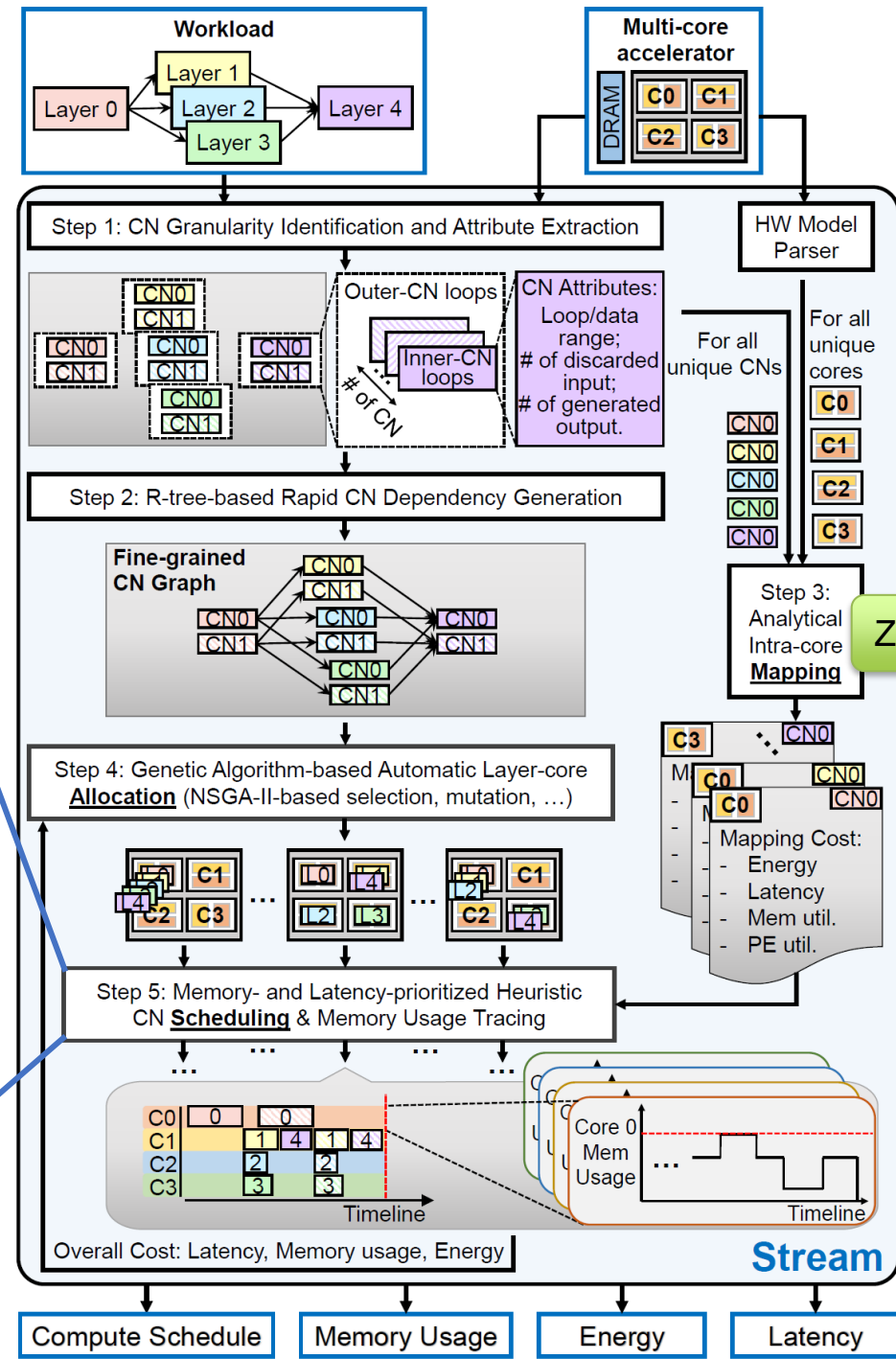
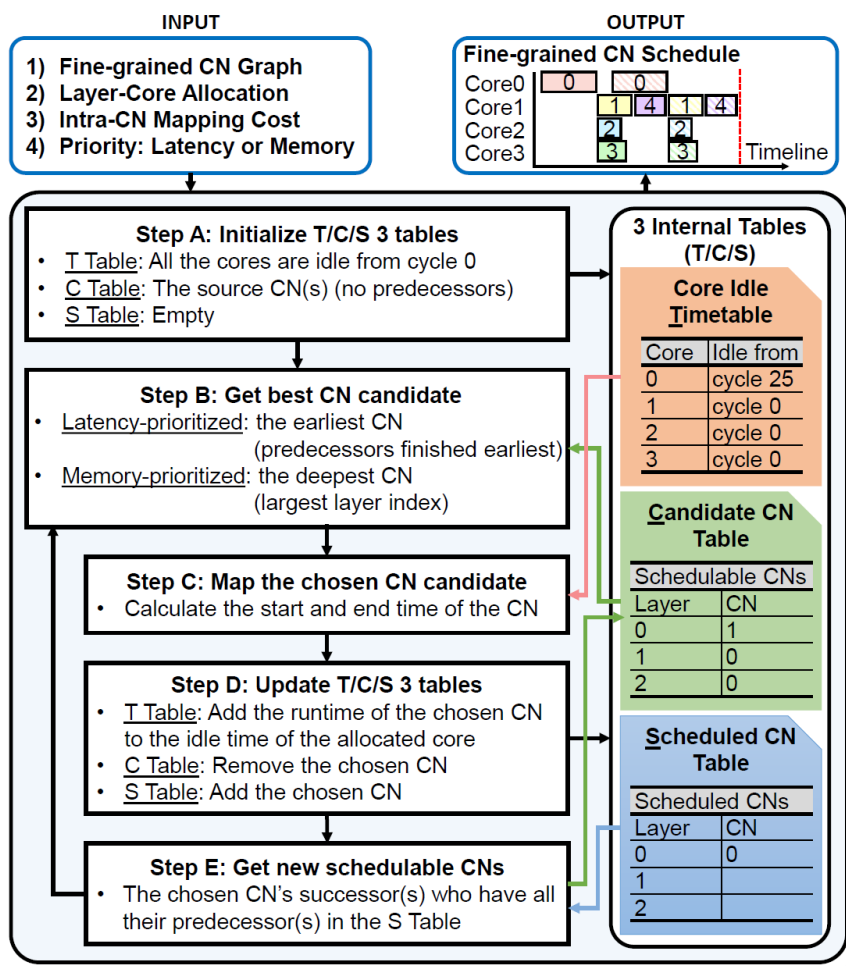
- ✓ Model single-core architecture (identical to ZigZag)
- ✓ Model different multi-core topologies
- ✓ Model different scheduling granularities



- ✓ Model single-core architecture (identical to ZigZag)
- ✓ Model different multi-core topologies
- ✓ Model different scheduling granularities
- ✓ Model different scheduling heuristics

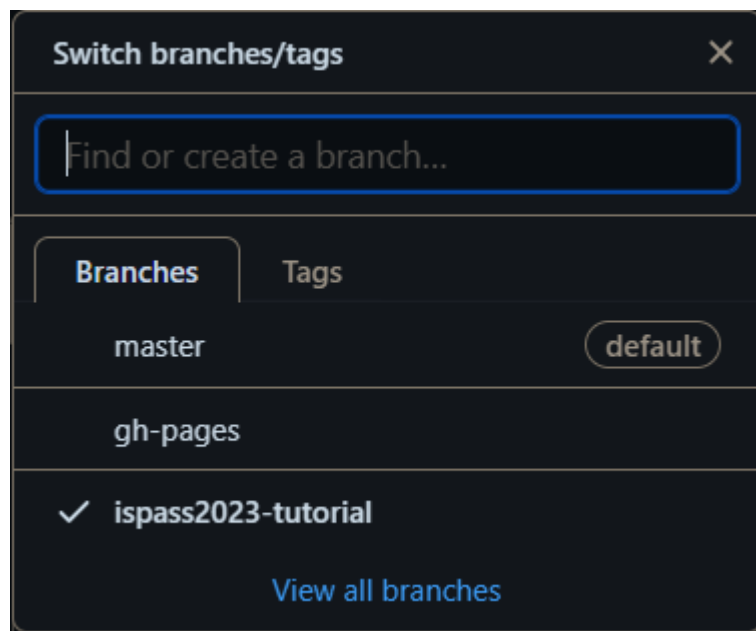


# Stream Overview

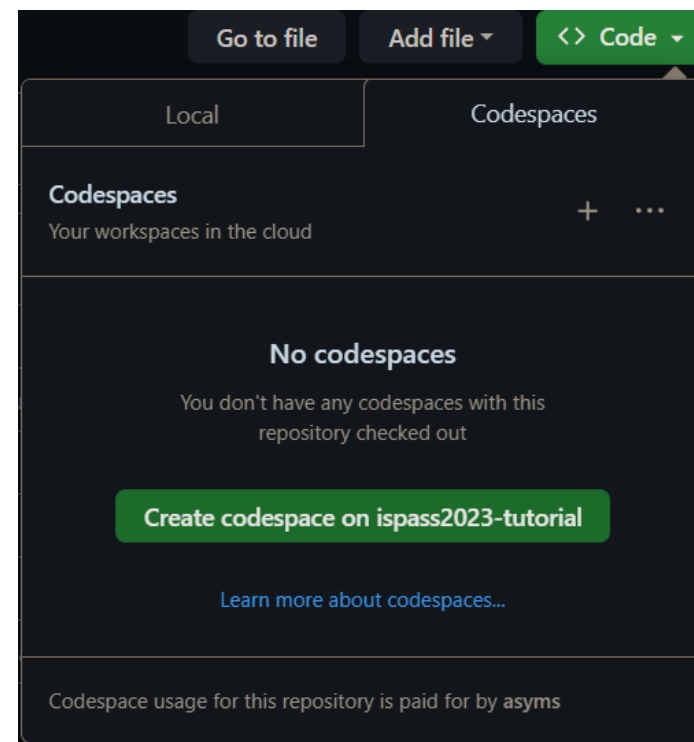


<https://github.com/ZigZag-Project/stream>

➤ Select ispass2023-tutorial



➤ Create codespace



<https://github.com/ZigZag-Project/stream>

```
$ git clone git@github.com:ZigZag-Project/stream.git
$ cd stream
$ conda create --name my-stream-env python=3.10
$ conda activate my-zigzag-env
$ pip install -r requirements.txt
$ git checkout ispass2023-tutorial
$ code .
```

- Open lab4/main\_fixed.py
- Defines inputs directly in file instead of arguments
- Extracts (from input names) and defines variables
- Sets up the sequence of stages
- Runs the stages
- Plots the results

## Workload

- Open  
lab4/inputs/workload/duplicated\_resnet18\_layer\_fixed.py
- First layer of ResNet18 duplicated 4 times
- Dependencies between the layers (**operand\_source**)
- **operator\_type** overloaded for fixed mapping



## Accelerator

- Open `lab4/inputs/hardware/heterogeneous_quadcore.py`
- Imports the “computational cores”
- Imports the pooling and simd cores
- Imports the offchip core
- Creates a 2D mesh of these cores
- Defines the multi-core Accelerator object

## Mapping

- Open lab4/inputs/mapping/mapping\_fixed.py
- Defines for each **operator\_type** the possible layer-core allocations

## First experiment:

- model = "...duplicated\_resnet18\_layer\_fixed"
- accelerator = "...heterogeneous\_quadcore"
- mapping = "...mapping\_fixed"
- Run lab4/main\_fixed.py

```
(my-stream-env) asymons@micaszb03:~/stream$ python lab4/main_fixed.py
```

## Second experiment:

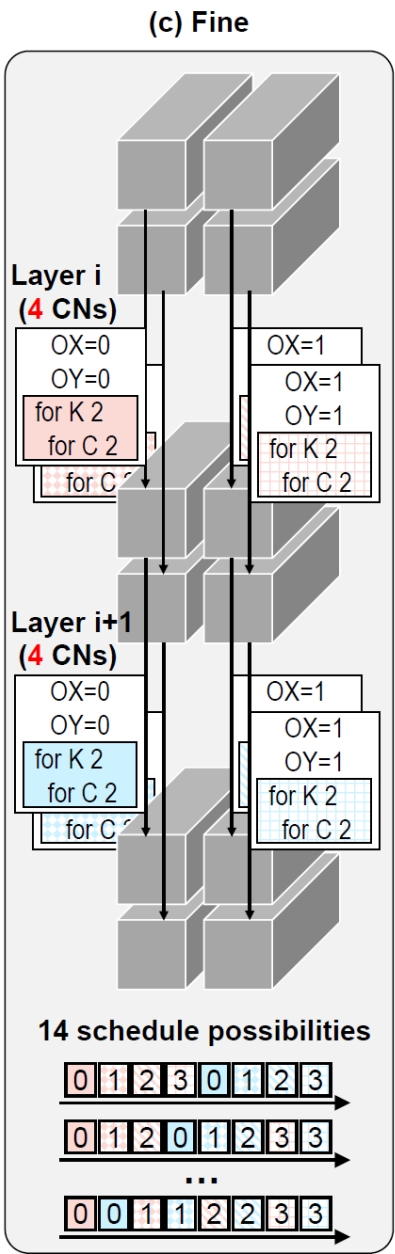
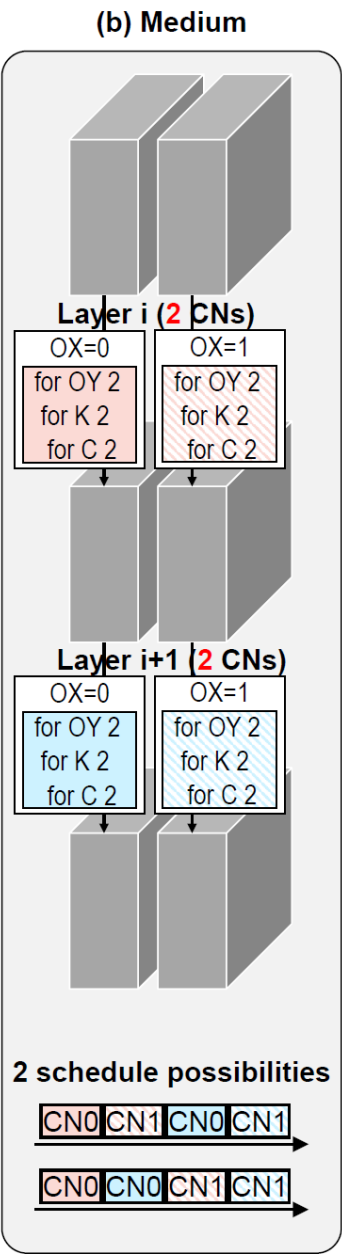
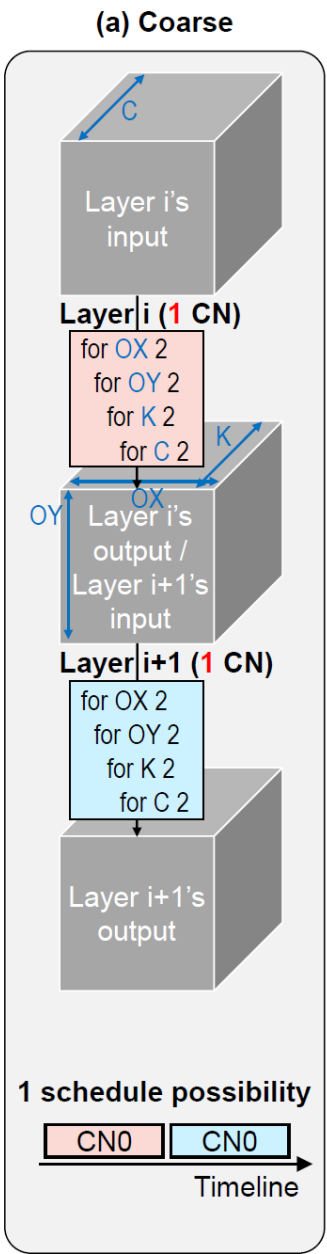
- What happens if we remove the layer dependencies?
- Remove the **operand\_source** and modify the **constant\_operands**
- Run lab4/main\_fixed.py

```
(my-stream-env) asymons@micaszb03:~/stream$ python lab4/main_fixed.py
```

### Third experiment:

- Allow genetic algorithm to find best layer-core allocation
- model = "...duplicated\_resnet18\_layer"
  - Modified **operator\_type**
- mapping = "...mapping"
  - Modified for flexible layer-core allocation
- Run lab4/main\_layer\_by\_layer.py

```
(my-stream-env) asymons@micaszb03:~/stream$ python lab4/main_layer_by_layer.py
```



**Computation node (CN)  
granularity impacts  
scheduling flexibility**

and others...

Core utilization  
Intra-CN data reuse  
Data loading overhead  
Control overhead

- Open lab4/main\_layer\_fused.py
- **hint\_loops** defines the outer-CN loops
  - hint\_loops = [("OY", 2)] means 2 CNs per layer
  - hint\_loops = [("OY", "all")] means OY CNs per layer

### Fourth experiment:

- Assess the impact of layer-fused processing
- Run lab4/main\_layer\_fused.py

```
(my-stream-env) asymons@micaszb03:~/stream$ python lab4/main_layer_fused.py
```



- Open lab4/main.py
- End-to-end ResNet18 onnx model
- Layer-by-layer

### Last experiment:

- End-to-end ResNet18 example
- Run lab4/main.py (layer-by-layer)
- If time permits:
  - Modify hint\_loops with (“OY”, “all”)
  - Re-run and analyze differences

- Automatically infer optimal CN granularity
- Integrate inter-core connect energy estimation framework
- Automatically search for optimal multi-core architectures
- Add optimization constraints (e.g. max latency, area, ...)
  
- Code generation for existing accelerators
- Automatic hardware generation from hardware templates

- ZigZag enables fast hardware performance estimation for specialized DNN accelerator architectures
- Mapping optimizations yield better performance
- Co-exploration of architecture with mapping
  
- Stream extends the capabilities to multi-core architectures employing layer-fused scheduling
- Unified hardware model for different architecture topologies
- Different scheduling granularities through computation node

- **Goal: Enabling Fast Architecture-Scheduling/Mapping DSE for Machine Learning Accelerators**
- Github: <https://github.com/ZigZag-Project>
  - ZigZag
  - ZigZag-Demo
  - DeFiNES
  - Stream
- ZigZag documentation: <https://zigzag-project.github.io/zigzag/>
- Stream documentation: Underway (end of May)
- ZigZag-related publications:  
<https://zigzag-project.github.io/zigzag/publications.html>