**PULP PLATFORM**
Open Source Hardware, the way it should be!

# *Deployment of DNN on Extreme Edge Devices (1)*

**Alessio Burrello <alessio.burrello@unibo.it>**
Francesco Conti <f.conti@unibo.it>

ETH*zürich*

# Bringing DNN Inference to the Edge

**ImageNet Top-1 Accuracy**
vs **Memory Footprint**

- Most entries > 10 MB

- Pareto Frontier Acc vs Memory
  (from 50% @ 0.5Mparam
  to 85% @ 445 Mparam)

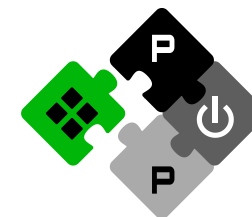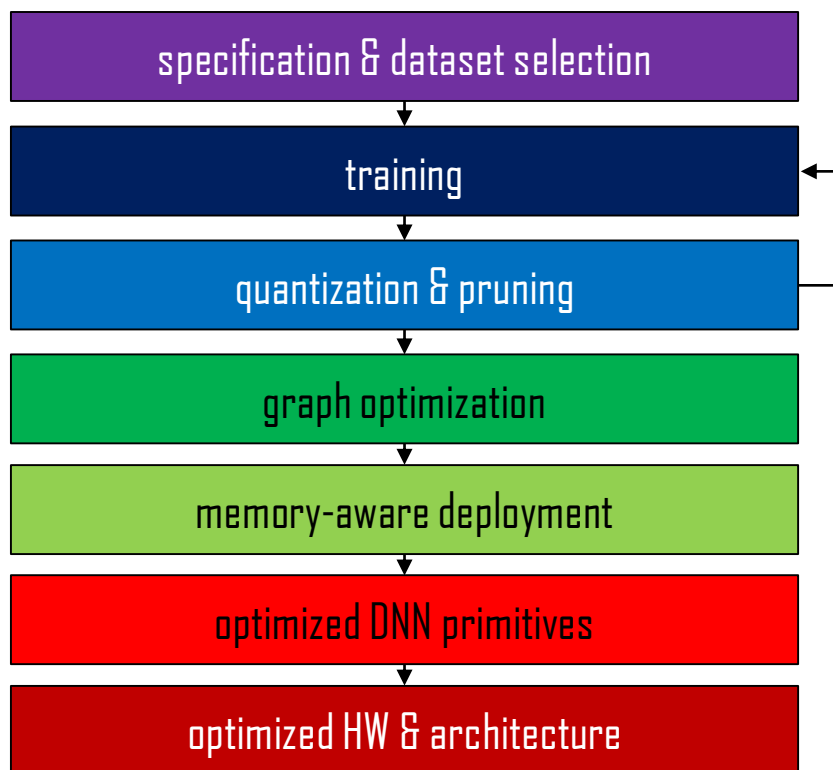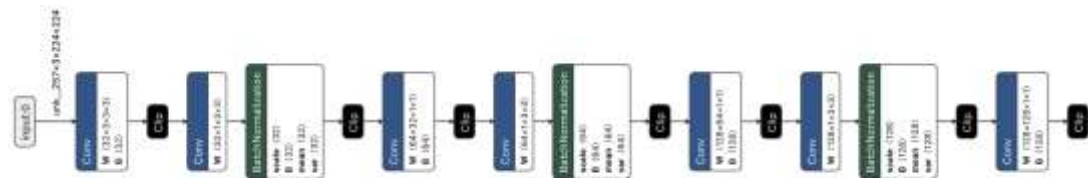- Almost always require off-chip DRAM
  even for **ULP**!

**1.0-MobileNetV1-224**

# Unibo Flow

- **Actually enabling execution of real-world sized DNNs at extreme edge is still a challenge**
  - most state-of-the-art (e.g. CMSIS-NN) shown on very small DNNs & datasets, e.g. CIFAR10
  - challenge #1: small and manually managed on-chip memory (512 kB L2, 64 kB fast L1 on most PULP-based chips)
  - challenge #2: better support for efficient integer computation, not floating point

- **We show the Unibo Flow, a vertically integrated framework for deployment of DNNs on PULP-based extreme edge platforms**
  - from algorithm definition (PyTorch) to running the DNN on the embedded platform (e.g., on GreenWaves GAP8, Mr. Wolf, PULP simulators)
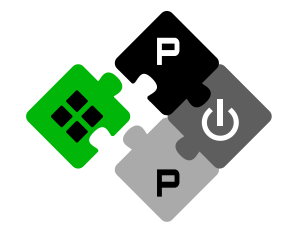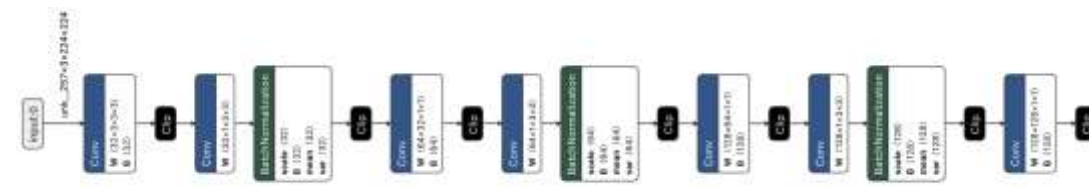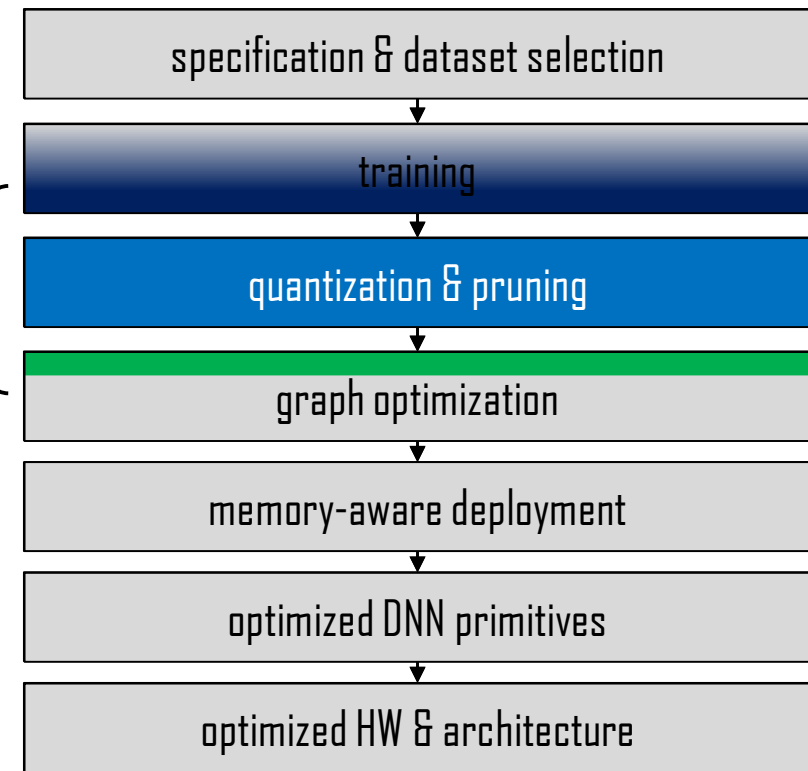
specification & dataset selection

training

quantization & pruning

graph optimization

memory-aware deployment

optimized DNN primitives

optimized HW & architecture

# Outline

1. **Intro on the UNIBO Flow**

2. **NEMO (*NE*ural *M*inimization for pyt*O*rch)**

    1. Topological Contraints

3. **DORY (*D*eployment *O*riented to memo*RY*)**

    1. Graph and Node reading

    2. Tiling

        • L3-L2 movement

        • L2-L1 movement

        • Data movement

    3. Template writing

4. **PULP-NN**

    1. Optimized backend

    2. Supported Layers

5. **How to Generate a Network**
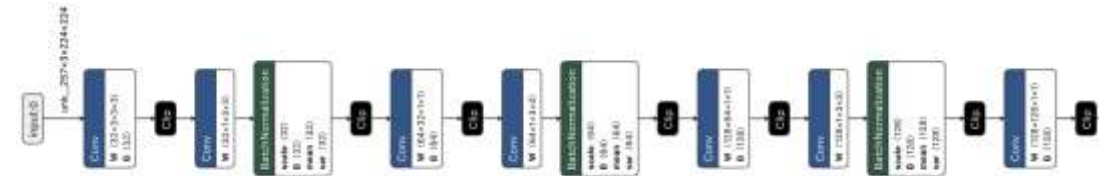
6. **Examples**
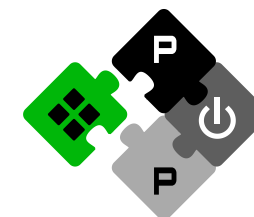
# Unibo Flow

**NEMO**
*NEural Minimization for pytOrch*

specification & dataset selection

training

quantization & pruning

graph optimization

memory-aware deployment

optimized DNN primitives

optimized HW & architecture

# Unibo Flow

**NEMO**
*NEural Minimization for pytOrch*

**DORY**
*Deployment Oriented to memoRY*

specification & dataset selection

training

quantization & pruning

graph optimization

memory-aware deployment

optimized DNN primitives

optimized HW & architecture

# Unibo Flow

**NEMO**
*NEural Minimization for pytOrch*

**DORY**
*Deployment Oriented to memoRY*

**PULP-NN**
*PULP Neural Network backend*

specification & dataset selection

training

quantization & pruning

graph optimization

memory-aware deployment

optimized DNN primitives

optimized HW & architecture

# Contributors

**NEMO**
*NEural Minimization for pytOrch*

**Francesco Conti**
**Marcello Zanghieri**
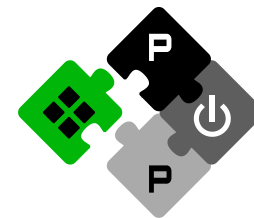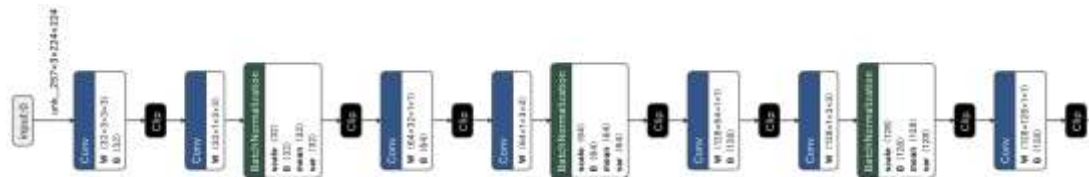**Leonardo Ravaglia**
**Lorenzo Lamberti**

**DORY**
*Deployment Oriented to memoRY*

**Alessio Burrello**
**Francesco Conti**
**Thorir Ingolfsson**

**PULP-NN**
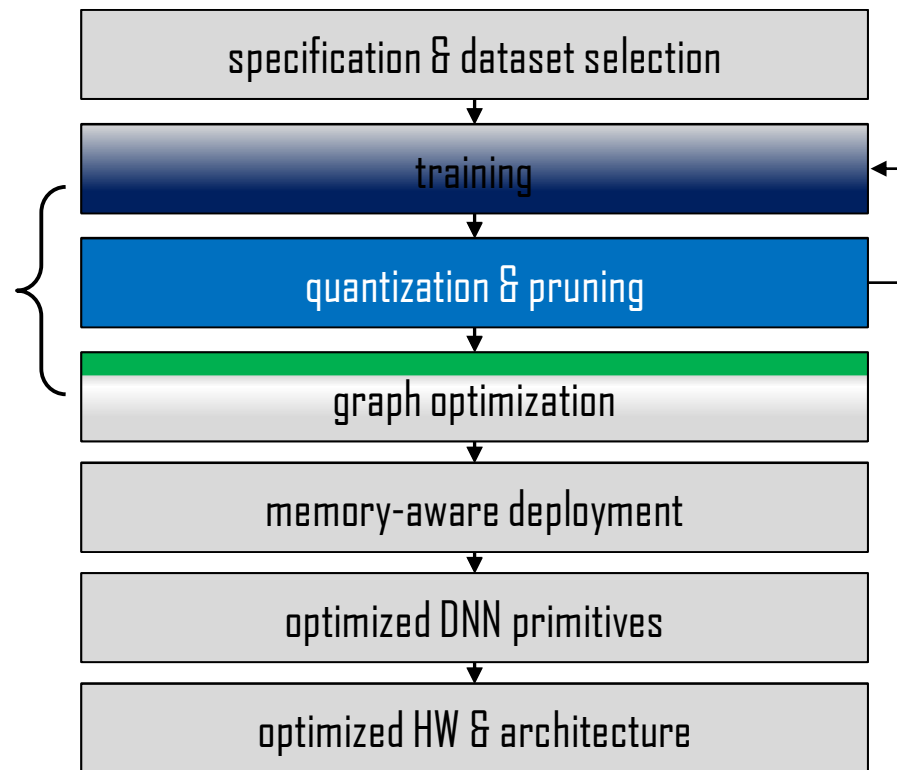*PULP Neural Network backend*

**Angelo Garofalo**
**Nazareno Bruschi**
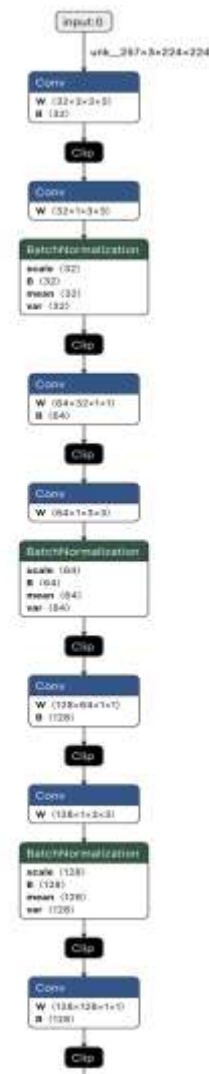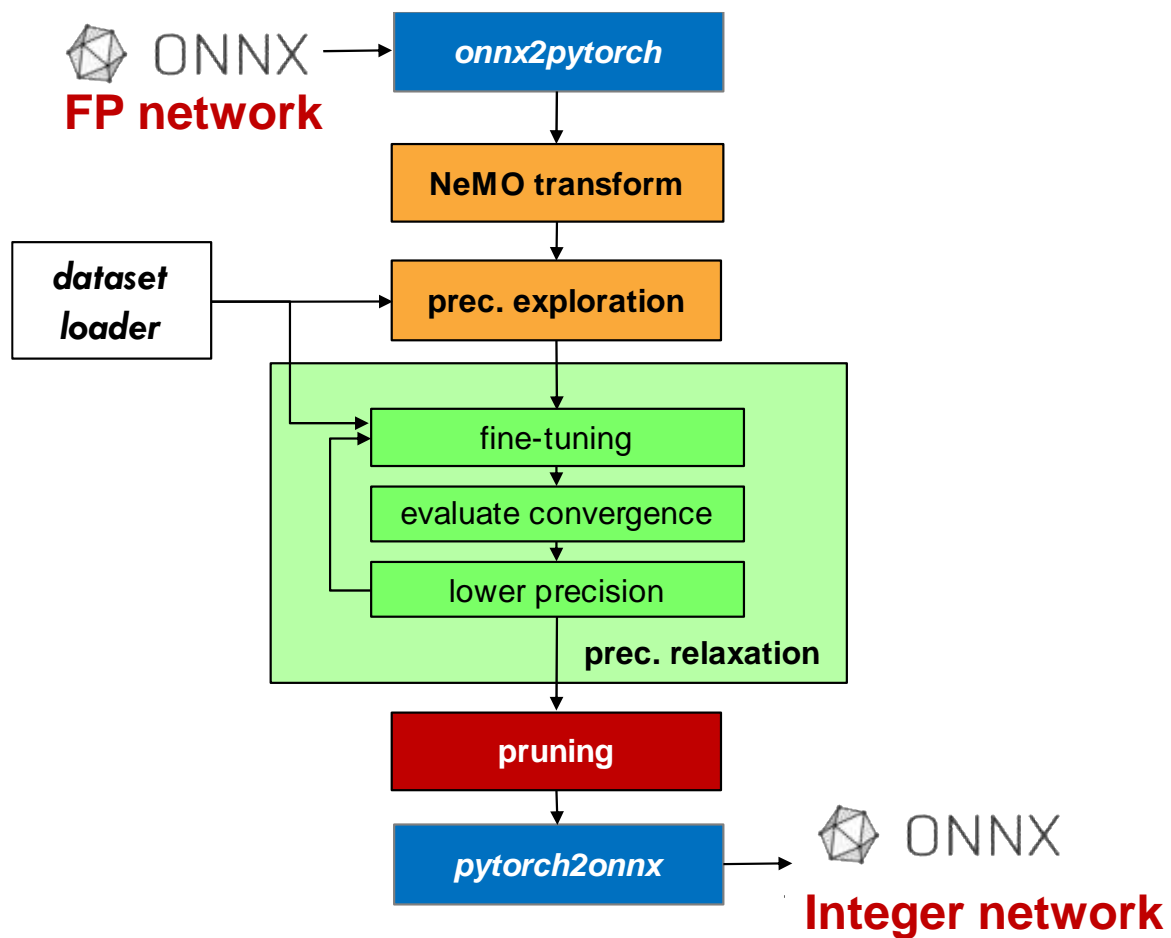
# NEMO: fp32 to full-integer networks

**NEMO**
**NE**ural **M**inimization for pyt**O**rch

From a full-precision representation to a fully integer (**not fixed-point**) HW-deployable one

| specification & dataset selection |
|:---:|
| training |
| quantization & pruning |
| graph optimization |
| memory-aware deployment |
| optimized DNN primitives |
| optimized HW & architecture |

# NEMO: quantization-aware retraining

# NEMO: topological constraints



**Integer BN**

**Quant**

1.  Recognize *super-layers* in the network
    *   typically, Conv+BN+Clip (quantization is implicit in QF format)

2.  Represent all tensors in the quantized form
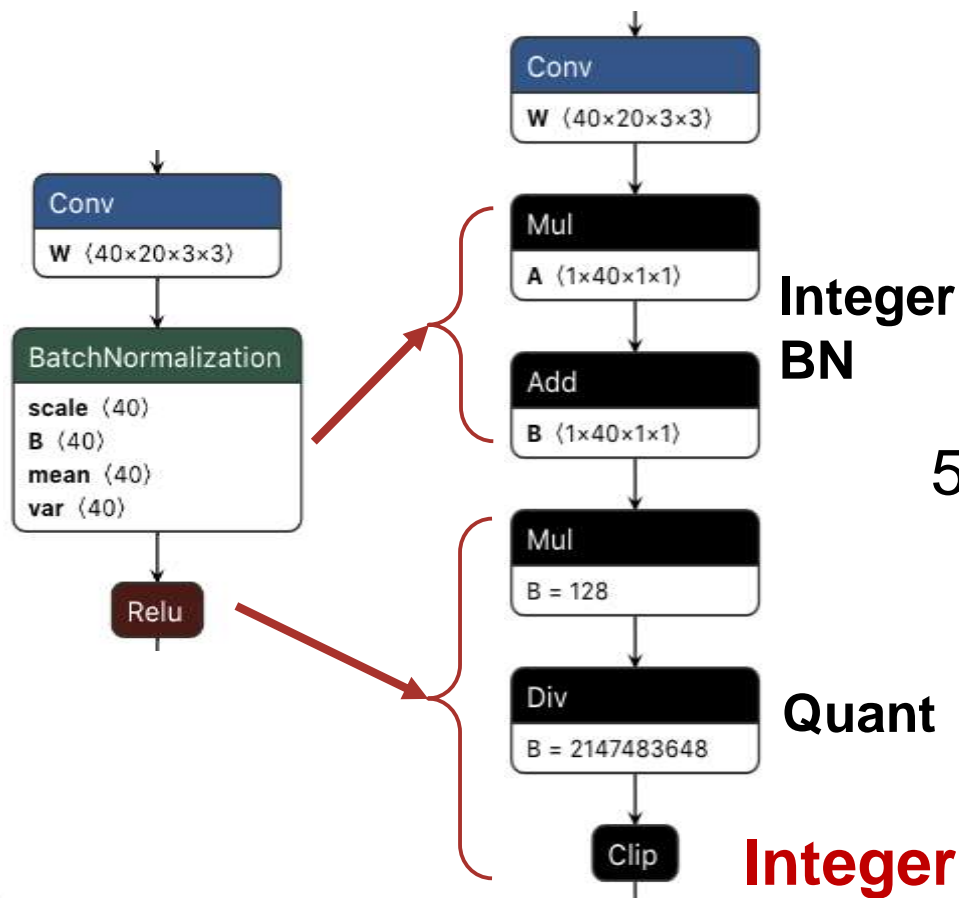
$$T = T_{int} \cdot \varepsilon_T$$

integer tensor
(**integer image**)          real-valued scalar
(**quantum**)

3.  Replace BN and Clip/Quant operations with equivalent working on quantized form and producing quantized tensors

# NEMO: topological constraints



**Integer BN**

**Quant**

**Integer-Deployable Network**

4. Keep track of $\varepsilon_T$ quanta along the network
   - linear operations produce outputs with smaller quantum (more bits)

   - non-linear activation produced outputs with quantum "collapsed" to a new value (usually requiring less bits) with **requantization**
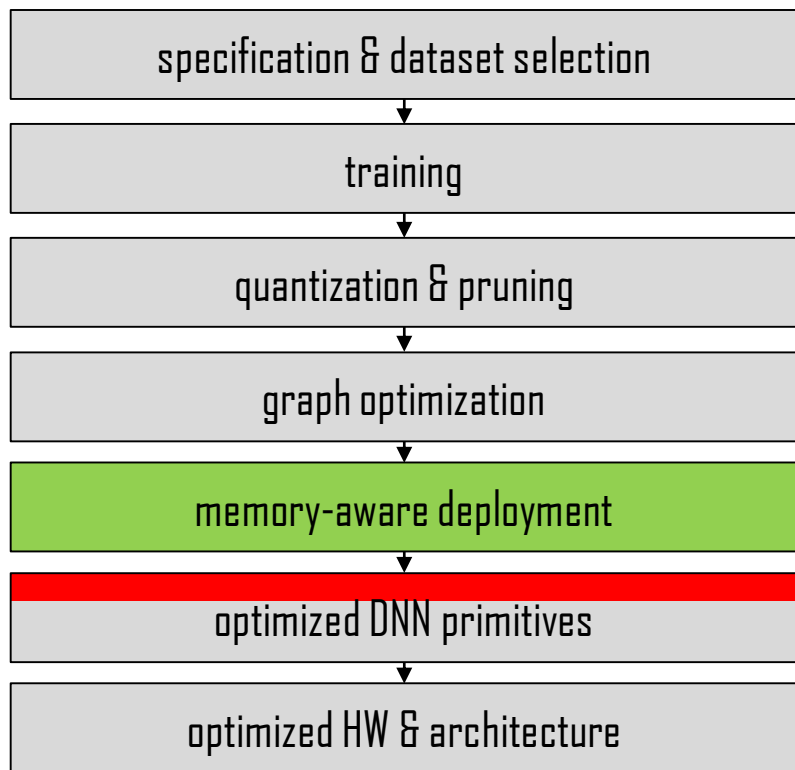
5. Replace all tensors by their integer image

$$T \rightarrow T_{int}$$

# DORY: Tiling & Code Generation

specification & dataset selection

training

quantization & pruning

graph optimization

memory-aware deployment

optimized DNN primitives

optimized HW & architecture

**DORY**
*Deployment Oriented to memoRY*

From an int8 quantized onnx network to a C compilable and runnable network

# DORY: Tiling & Code Generation

**DORY**
**D**eployment **O**riented to memo**RY**

1. Reading of the ONNX output
   1. Recognize backend implemented nodes
   2. Reconstruct the graph with backend nodes input-output dimensions
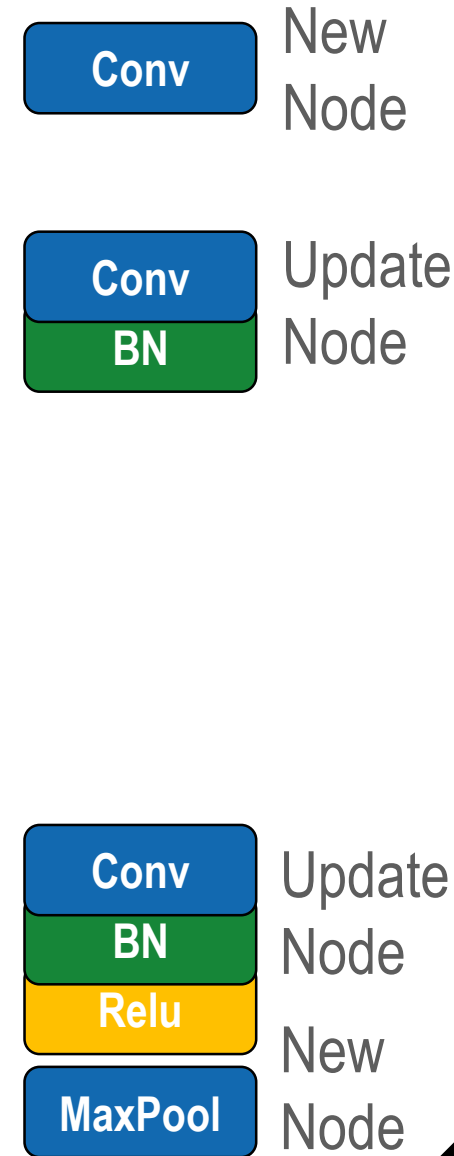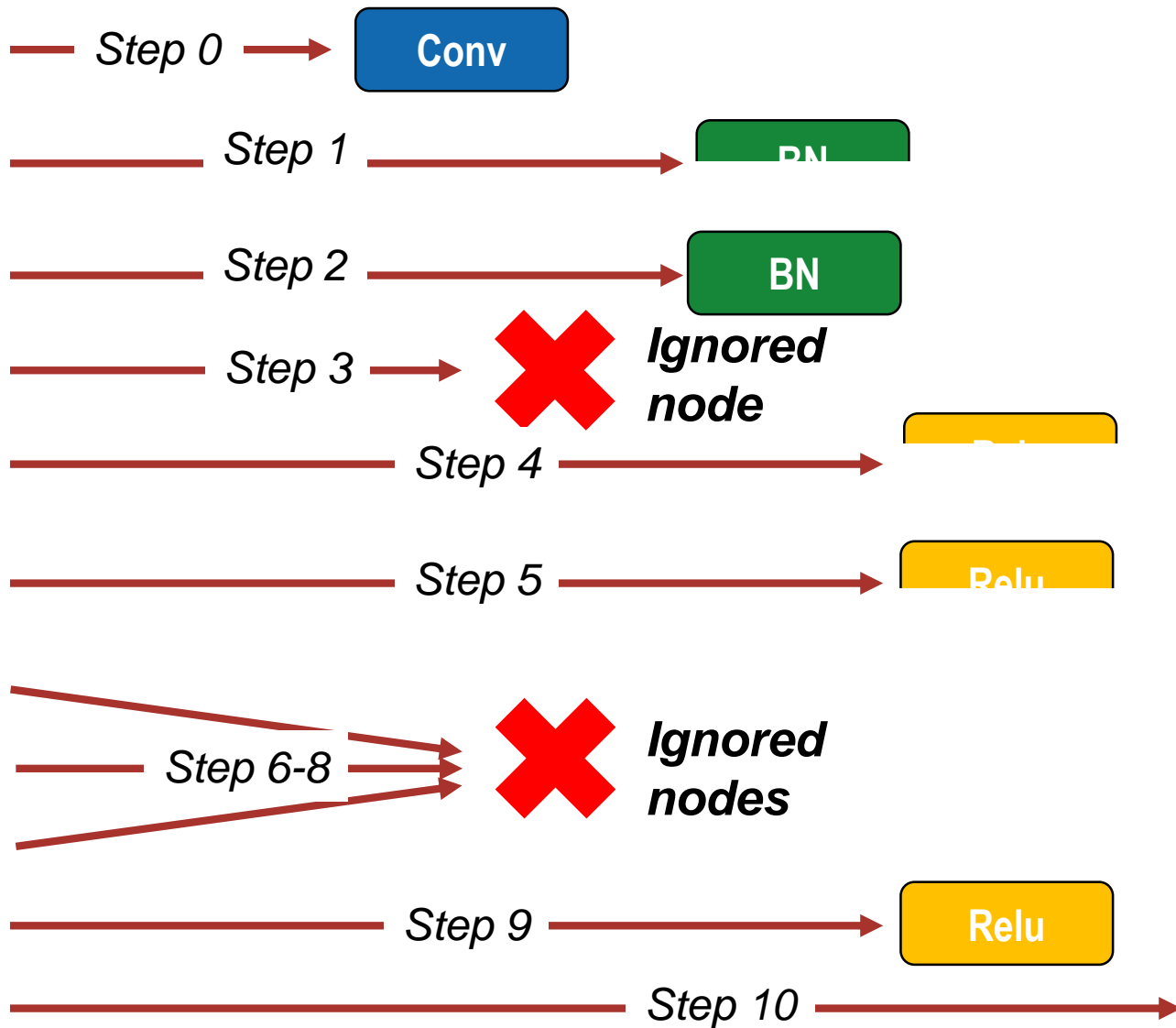2. Layer-by-Layer tiling
   1. L3-L2 tiling
   2. L2-L1 tiling
   3. Memory allocation in L2
3. Layer template compilation
4. Network compilation

# DORY: Tiling & Code Generation

***DORY***
***D**eployment **O**riented to memo**RY***

## 1. Reading of the ONNX output
1. **Recognize backend implemented nodes**
2. **Reconstruct the graph with backend nodes input-output dimensions**
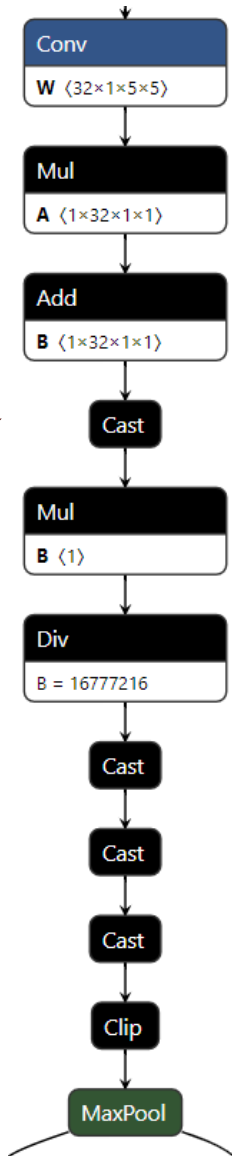
## 2. Layer-by-Layer tiling
1. L3-L2 tiling
2. L2-L1 tiling
3. Memory allocation in L2

## 3. Layer template compilation

## 4. Network compilation

# DORY: ONNX Decoding

**Graph Parsing**

| Conv | | |
|------|------|------|
| W ⟨32×1×5×5⟩ | | |

Mul
A ⟨1×32×1×1⟩

Add
B ⟨1×32×1×1⟩

Cast

Mul
B ⟨1⟩

Div
B = 16777216

Cast

Cast

Cast

Clip

MaxPool

— Step 0 → **Conv**

Step 1 → **BN**

Step 2 → **BN**

Step 3 → ❌ *Ignored node*

Step 4 → **Relu**

Step 5 → **Relu**

Step 6-8 → ❌ *Ignored nodes*

Step 9 → **Relu**

Step 10 → **MaxPool**

**Conv** — New Node

**Conv** / **BN** — Update Node

**Conv** / **BN** / **Relu** — Update Node

**MaxPool** — New Node

# DORY: ONNX Decoding

**Conv**
**W** ⟨40×20×3×3⟩

**Mul**
**A** ⟨1×40×1×1⟩

**Add**
**B** ⟨1×40×1×1⟩

**Mul**
B = 128

**Div**
B = 2147483648

**Clip**

***ONNX***
***READER***

*New node_iterating:*
**ConvBNRelu** ──────→ ***Layer name***
**Filter Dimension**
**Stride**
**Padding**
**Groups**
**MACs**
**In-Out dimensions**
*k: present*
*lambd: present*
*outmul: present*
*outshift: present*
*Input branch: No*
*Output branch: No*
*Input: 93*
*Output: 105*

# DORY: ONNX Decoding



**Conv**
W ⟨40×20×3×3⟩

**Mul**
A ⟨1×40×1×1⟩

**Add**
B ⟨1×40×1×1⟩

**Mul**
B = 128

**Div**
B = 2147483648

**Clip**

**ONNX READER**

New node_iterating:
**ConvBNRelu** ⟶ **Layer name**
**Filter Dimension**
**Stride**
**Padding**
**Groups**
**MACs**
**In-Out dimensions**

**Conv/Linear Parameters**

k: present
lambd: present
outmul: present
outshift: present
Input branch: No
Output branch: No
Input: 93
Output: 105

# DORY: ONNX Decoding

Conv
W ⟨40×20×3×3⟩

Mul
A ⟨1×40×1×1⟩

Add
B ⟨1×40×1×1⟩

Mul
B = 128

Div
B = 2147483648

Clip

**ONNX READER**

New node_iterating:
**ConvBNRelu** ⟶ **Layer name**
**Filter Dimension**
**Stride**
**Padding**
**Groups**
**MACs**
**In-Out dimensions**

**Conv/Linear Parameters**

k: present
lambd: present

**Batchnorm: in x k + λ**

outmul: present
outshift: present
Input branch: No
Output branch: No
Input: 93
Output: 105

# DORY: ONNX Decoding



**Conv**
W ⟨40×20×3×3⟩

**Mul**
A ⟨1×40×1×1⟩

**Add**
B ⟨1×40×1×1⟩

**Mul**
B = 128

**Div**
B = 2147483648

**Clip**

*ONNX READER*

New node_iterating:
**ConvBNRelu** ⟶ *Layer name*
**Filter Dimension**
**Stride**
**Padding**
**Groups**
**MACs**
**In-Out dimensions**

*Conv/Linear Parameters*

k: present
lambd: present

**Batchnorm: in x k + λ**

outmul: present
outshift: present

**Relu: clip8(in x mul >> shift)**

Input branch: No
Output branch: No
Input: 93
Output: 105

# DORY: ONNX Decoding



**Conv**
W ⟨40×20×3×3⟩

**Mul**
A ⟨1×40×1×1⟩

**Add**
B ⟨1×40×1×1⟩

**Mul**
B = 128

**Div**
B = 2147483648

**Clip**

*ONNX READER* →

*New node_iterating:*

**ConvBNRelu** ⟶ *Layer name*

**Filter Dimension**
**Stride**
**Padding**
**Groups**
**MACs**
**In-Out dimensions**

*Conv/Linear Parameters*

*k: present*
*lambd: present*

**Batchnorm: in x k + λ**

*outmul: present*
*outshift: present*

**Relu: clip8(in x mul >> shift)**

*Input branch: No*
*Output branch: No*
*Input: 93*
*Output: 105*

*Network topology parameters*

# DORY: Tiling & Code Generation

**DORY**
***D**eployment **O**riented to memo**RY***

1. Reading of the ONNX output
    1. Recognize backend implemented nodes
    2. Reconstruct the graph with backend nodes input-output dimensions
2. **Layer-by-Layer tiling**
    1. **L3-L2 tiling**
    2. **L2-L1 tiling**
    3. **Memory allocation in L2**
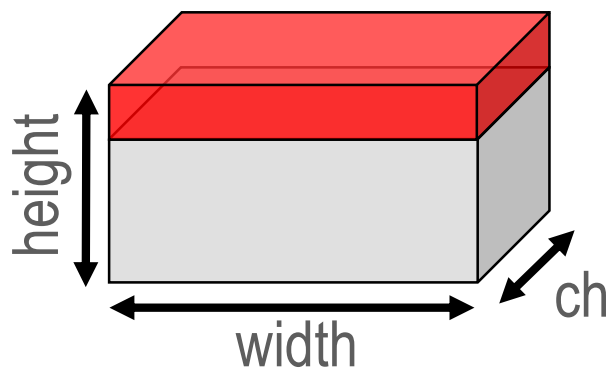3. Layer template compilation
4. Network compilation
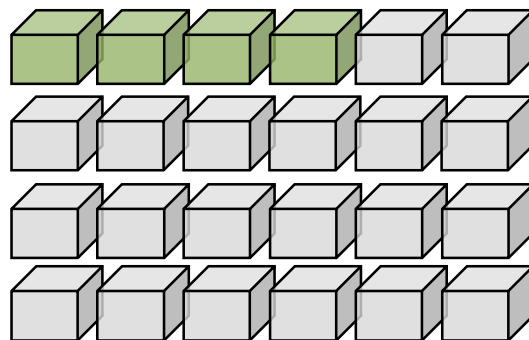
# DORY: Tiler

**L3** / **L2** tiling
64 MB / 512 kB



small memory

big memory

layer 1

layer 2

# DORY: Tiler – L3/L2

**L3/L2 Tiling:**

- Large L3 Memory → Enable Big Networks 🙂
- Small Memory Bandwidth → Slow Down Execution 🙁

**L3/L2 Tiling *steps*:**

1. Input tiling 🙂

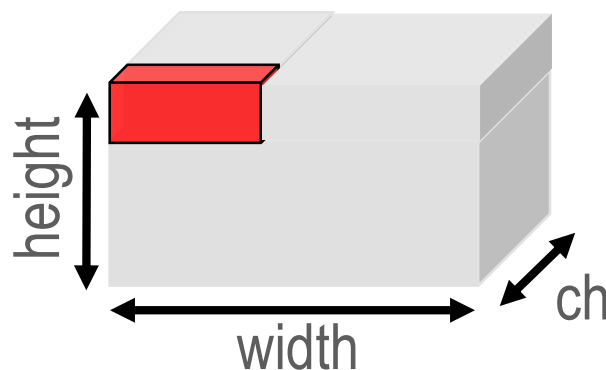2. Weights tiling 🙂

3. Output tiling 🙁

Output act.          Input act.



All tiles from L3 to L2 are 1D. Only uDMA linear transfers are required.

# DORY: Tiler – L2/L1

**L3** / **L2** tiling
64 MB / 512 kB

small memory

**L2** / **L1** tiling
512 kB / 64 kB

big memory



layer 1

layer 2

# DORY: Tiler – L2/L1

**L2/L1 Tiling:**

- Relatively low L2 Memory 🙁
- Large Memory Bandwidth 🙂

All tiles from L3 to L2 are 3D

height
width
ch

L2/L1 tiling is formalized as an **optimization problem.**

We use **Constraint Programming** to formalize the problem and find a feasible solution

# DORY: Tiler – L2/L1

**Constraint Programming** problem → tiles size

$$\mathbf{cost} = \max Size(W_{tile}) + Size(x_{tile}) + Size(y_{tile})$$

**MEMORY** → s.t. $Size(W_{tile}) + Size(x_{tile}) + Size(y_{tile}) < L1size$

**GEOMETRY** → s.t. $\{y_{tile}[ch_{out}] = W_{tile}[ch_{out}], ...\}$

**EFF. HEURISTICS** → $\mathbf{cost}' = \mathbf{cost} + \{y_{tile}[ch_{out}] \text{ divisible by } 4, ...\}$

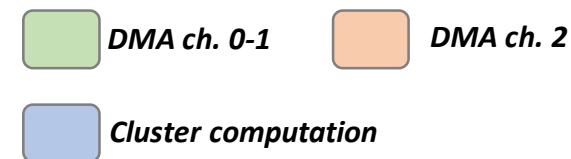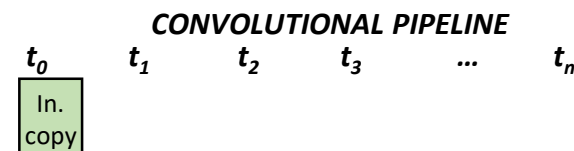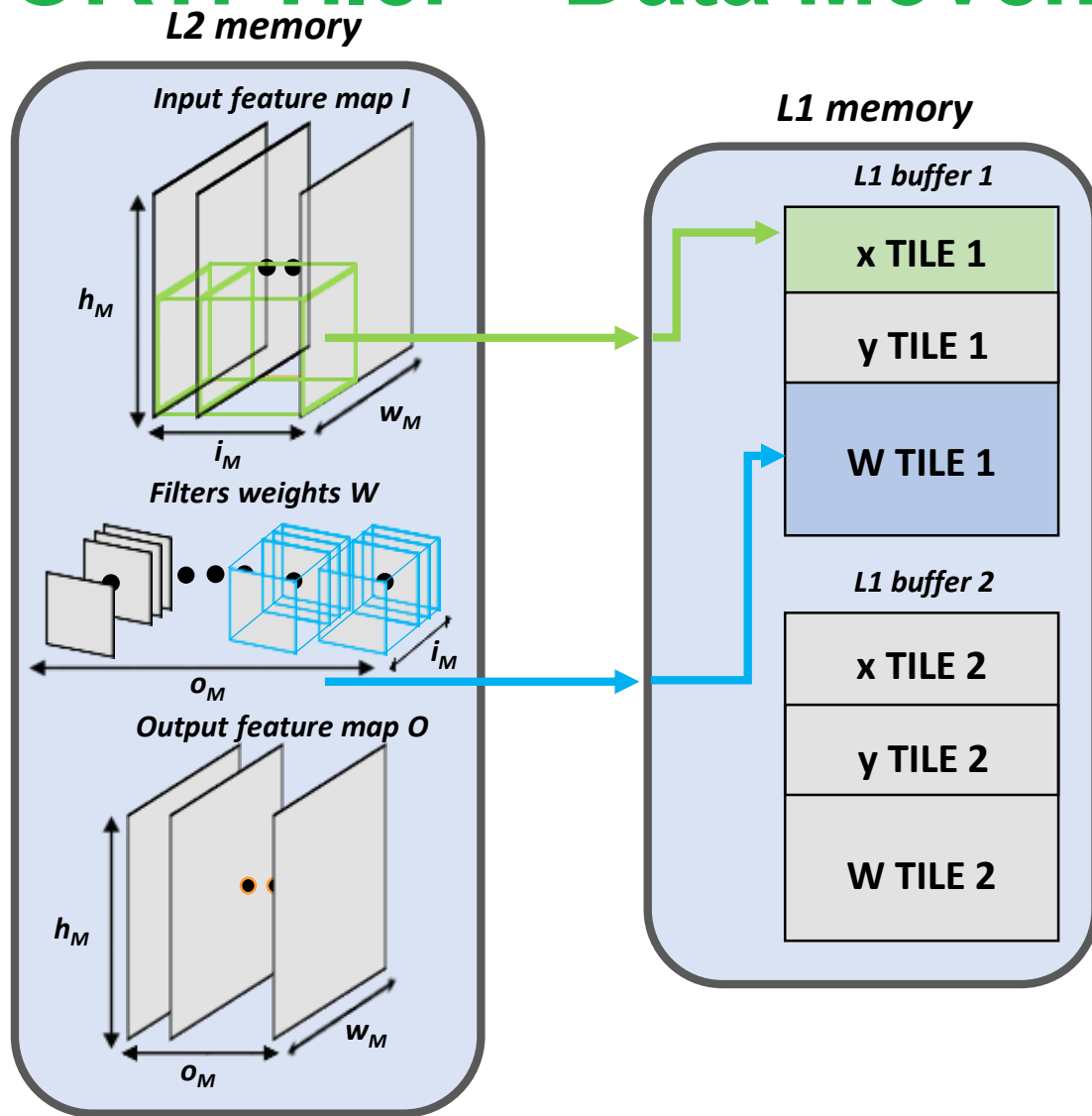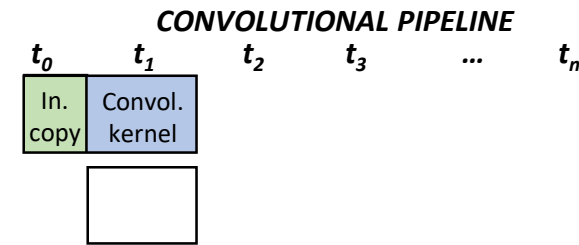*Performance is maximum for configurations that use PULP-NN primitives more efficiently (e.g., full parallelism)*
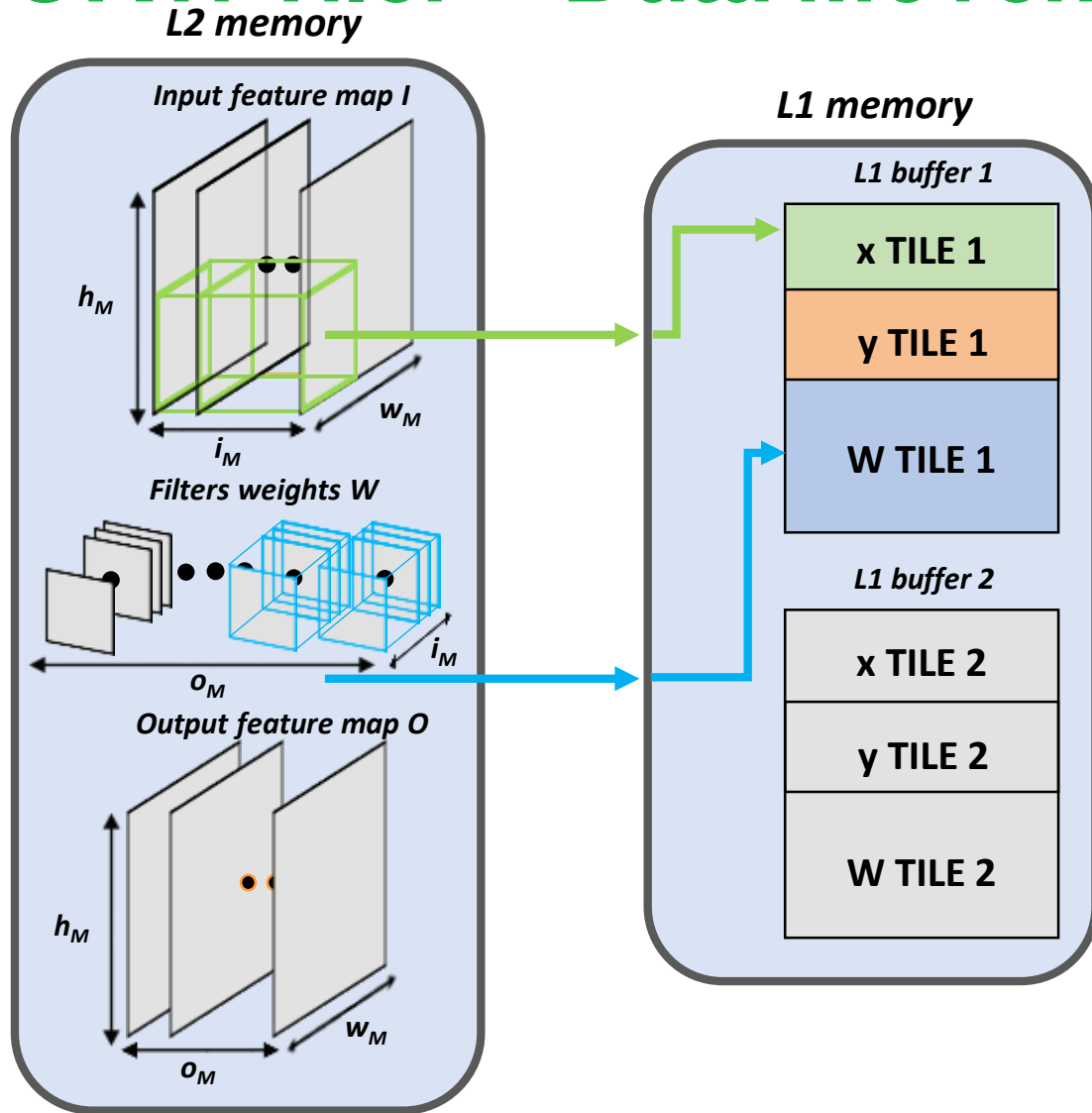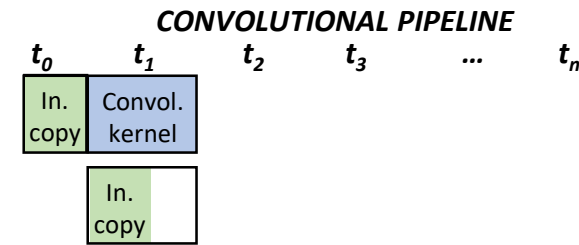
ONNX

**Integer DNN**

Google **ORTools**

**Integer DNN + tile sizes**

# DORY: Tiler – Data Movement

# DORY: Tiler – Data Movement

# DORY: Tiler – Data Movement



**L2 memory**

Input feature map I

$h_M$   $i_M$   $w_M$

Filters weights W

$o_M$   $i_M$

Output feature map O

$h_M$   $o_M$   $w_M$

**L1 memory**

L1 buffer 1

x TILE 1
y TILE 1
W TILE 1

L1 buffer 2

x TILE 2
y TILE 2
W TILE 2

**CONVOLUTIONAL PIPELINE**

$t_0$   $t_1$   $t_2$   $t_3$   ...   $t_n$

In. copy | Convol. kernel
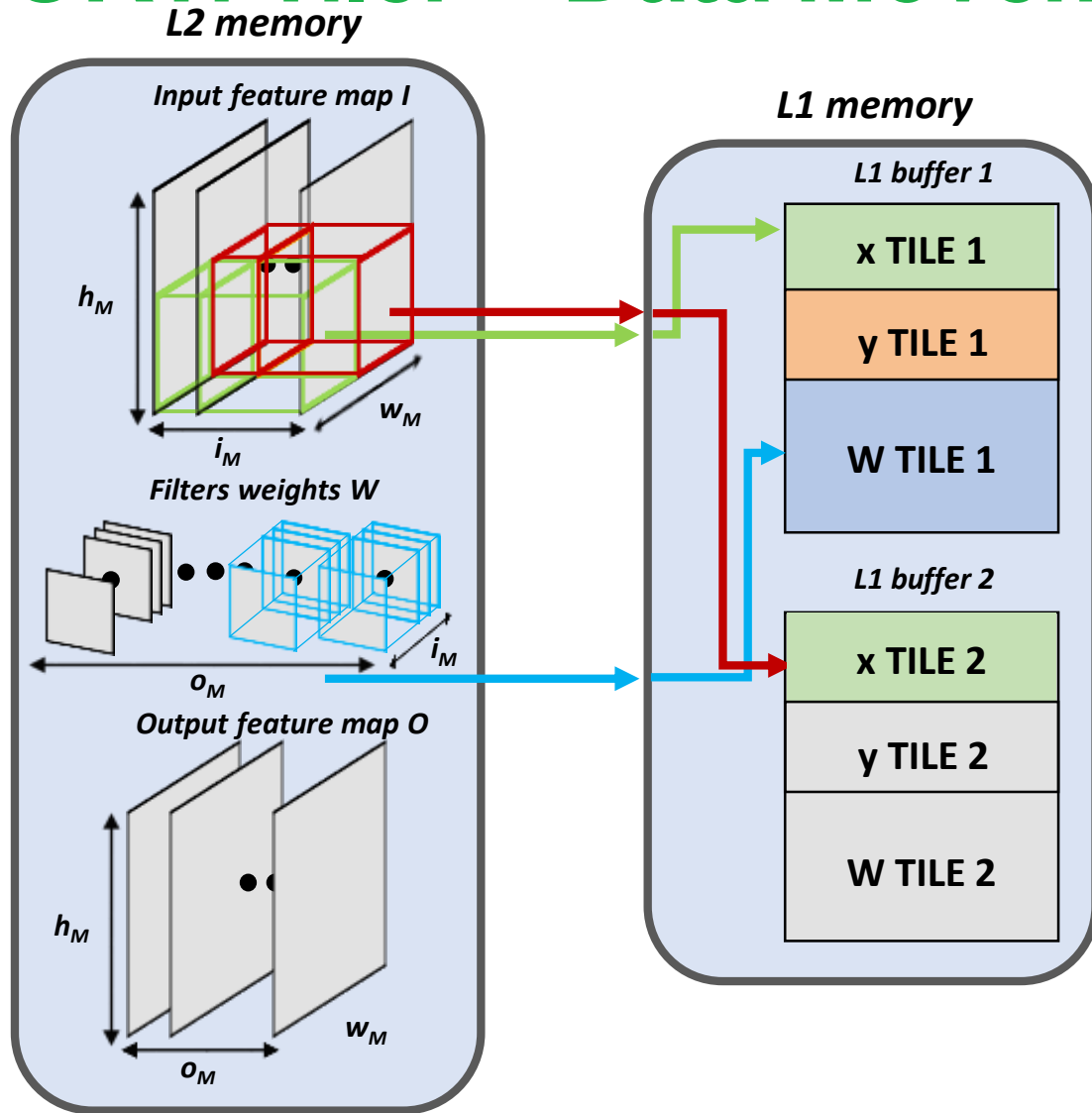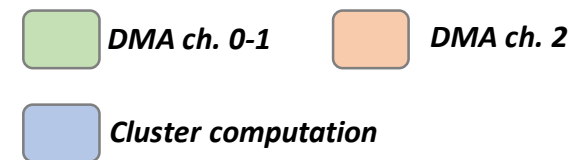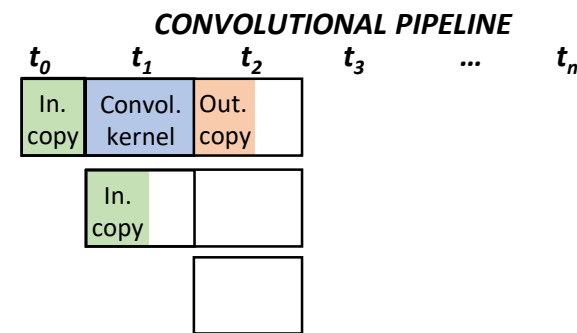
In. copy

DMA ch. 0-1      DMA ch. 2

Cluster computation

# DORY: Tiler – Data Movement

# DORY: Tiler – Data Movement



**L2 memory**

Input feature map I

Filters weights W

Output feature map O
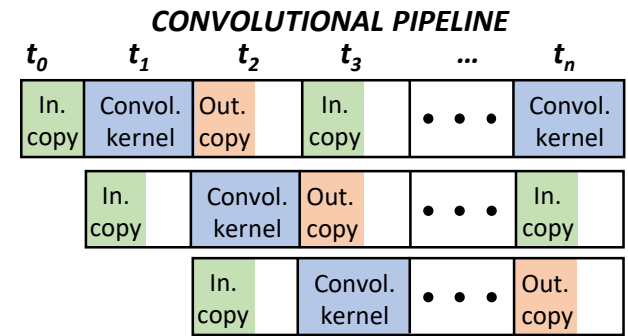
**L1 memory**

L1 buffer 1

x TILE 1

y TILE 1

W TILE 1

L1 buffer 2

x TILE 2

y TILE 2

W TILE 2

**CONVOLUTIONAL PIPELINE**

DMA ch. 0-1

DMA ch. 2

Cluster computation

# DORY: Tiling & Code Generation

**DORY**
**D**eployment **O**riented to memo**RY**

1. Reading of the ONNX output
    1. Recognize backend implemented nodes
    2. Reconstruct the graph with backend nodes input-output dimensions
2. Layer-by-Layer tiling
    1. L3-L2 tiling
    2. L2-L1 tiling
    3. Memory allocation in L2
3. **Layer template compilation**
4. Network compilation

# DORY: Template Writing

**Neural Network Layers generation**

mako.template → python compilation of c templates

```
dory_dma_memcpy_3d(input_0, ${args});
dory_dma_memcpy_3d(weights_0, ${args});
dory_dma_wait();
for (i=0; i<${tile_dim_nof * tile_dim_nif * tile_dim_h * tile_dim_w}; i++)
    dory_dma_memcpy_3d(input_i+1, ${args});
    dory_dma_memcpy_3d(weights_i+1, ${args});
    pulp_nn_conv(input_i, weights_i, output, ${args});
    dory_dma_wait();
    dory_dma_memcpy_3d(output, ${args});
```

# DORY: Template Writing

## Neural Network Layers generation

mako.template → python compilation of c templates

**Network exported parameters**

```
dory_dma_memcpy_3d(input_0, ${args});
dory_dma_memcpy_3d(weights_0, ${args});
dory_dma_wait();
for (i=0; i<${tile_dim_nof * tile_dim_nif * tile_dim_h * tile_dim_w}; i++)
    dory_dma_memcpy_3d(input_i+1, ${args});
    dory_dma_memcpy_3d(weights_i+1, ${args});
    pulp_nn_conv(input_i, weights_i, output, ${args});
    dory_dma_wait();
    dory_dma_memcpy_3d(output, ${args});
```

**pulp_nn kernel**

# DORY: Template Writing

## Neural Network Layers generation

mako.template → python compilation of c templates

**L2/L1 memory copies**

```
dory_dma_memcpy_3d(input_0, ${args});
dory_dma_memcpy_3d(weights_0, ${args});
dory_dma_wait();
for (i=0; i<${tile_dim_nof * tile_dim_nif * tile_dim_h * tile_dim_w}; i++)
    dory_dma_memcpy_3d(input_i+1, ${args});
    dory_dma_memcpy_3d(weights_i+1, ${args});
    pulp_nn_conv(input_i, weights_i, output, ${args});
    dory_dma_wait();
    dory_dma_memcpy_3d(output, ${args});
```

→ **First tile allocation**

# DORY: Template Writing

## Neural Network Layers generation

mako.template → python compilation of c templates

```
dory_dma_memcpy_3d(input_0, ${args});
dory_dma_memcpy_3d(weights_0, ${args});          → First tile allocation
dory_dma_wait();
for (i=0; i<${tile_dim_nof * tile_dim_nif * tile_dim_h * tile_dim_w}; i++)
    dory_dma_memcpy_3d(input_i+1, ${args});
    dory_dma_memcpy_3d(weights_i+1, ${args});
    pulp_nn_conv(input_i, weights_i, output, ${args});
    dory_dma_wait();                                 → Tile loop
    dory_dma_memcpy_3d(output, ${args});
```

# DORY: Template Writing

## Neural Network Layers generation

mako.template → python compilation of c templates

```
dory_dma_memcpy_3d(input_0, ${args});
dory_dma_memcpy_3d(weights_0, ${args});
dory_dma_wait();
for (i=0; i<${tile_dim_nof * tile_dim_nif * tile_dim_h * tile_dim_w}; i++)
    dory_dma_memcpy_3d(input_i+1, ${args});         → Async Data movement
    dory_dma_memcpy_3d(weights_i+1, ${args});
    pulp_nn_conv(input_i, weights_i, output, ${args});   → Kernel Computation
    dory_dma_wait();
    dory_dma_memcpy_3d(output, ${args});            → Async Data movement
```

# DORY: Tiling & Code Generation

**DORY**
**D**eployment **O**riented to memo**RY**

1. Reading of the ONNX output
   1. Recognize backend implemented nodes
   2. Reconstruct the graph with backend nodes input-output dimensions
2. Layer-by-Layer tiling
   1. L3-L2 tiling
   2. L2-L1 tiling
   3. Memory allocation in L2
3. Layer template compilation
4. **Network compilation**

# DORY: Network Generation

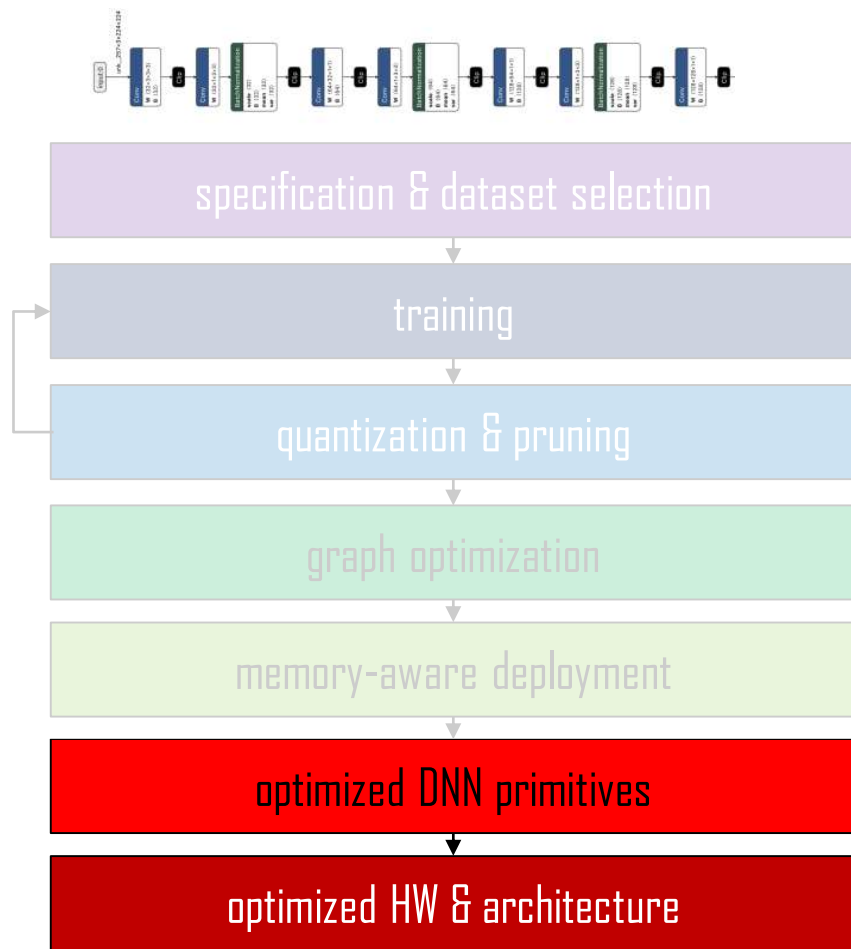**Neural Network generation** ➔ mako.template

```
for (int i = 0; i < ${len(PULP_Nodes_Graph)}; i++)
    pi_cl_ram_read_wait(&buff_req1);
    pi_cl_ram_read(&ram, transfer_weights, ${args}, &buff_req1);
    switch (i)
    {
    % for i in range(len(PULP_Nodes_Graph)):
        case ${i}:
            ${func_name[i]}(args);
            break;
    % endfor
    }
    dory_L2_memory_management();
```

# DORY: Network Generation

**Neural Network generation** → mako.template

**Loop over layers**

```
for (int i = 0; i < ${len(PULP_Nodes_Graph)}; i++)
        pi_cl_ram_read_wait(&buff_req1);
        pi_cl_ram_read(&ram, transfer_weights, ${args}, &buff_req1);
        switch (i)
        {
        % for i in range(len(PULP_Nodes_Graph)):
            case ${i}:
                ${func_name[i]}(args);
                break;
        % endfor
        }
        dory_L2_memory_management();
```

# DORY: Network Generation

**Neural Network generation** → mako.template

**L3 DMA weights memory copy**

```
for (int i = 0; i < ${len(PULP_Nodes_Graph)}; i++)
        pi_cl_ram_read_wait(&buff_req1);
        pi_cl_ram_read(&ram, transfer_weights, ${args}, &buff_req1);
        switch (i)
        {
        % for i in range(len(PULP_Nodes_Graph)):
            case ${i}:
                ${func_name[i]}(args);
                break;
        % endfor
        }
        dory_L2_memory_management();
```

# DORY: Network Generation

**Neural Network generation** → mako.template

```
for (int i = 0; i < ${len(PULP_Nodes_Graph)}; i++)
    pi_cl_ram_read_wait(&buff_req1);
    pi_cl_ram_read(&ram, transfer_weights, ${args}, &buff_req1);
    switch (i)
    {
    % for i in range(len(PULP_Nodes_Graph)):
        case ${i}:
            ${func_name[i]}(args);
            break;
    % endfor
    }
    dory_L2_memory_management();
```

**Convolutional layers**

# DORY: Network Generation

**Neural Network generation → mako.template**

```
for (int i = 0; i < ${len(PULP_Nodes_Graph)}; i++)
    pi_cl_ram_read_wait(&buff_req1);
    pi_cl_ram_read(&ram, transfer_weights, ${args}, &buff_req1);
    switch (i)
    {
    % for i in range(len(PULP_Nodes_Graph)):
        case ${i}:
            ${func_name[i]}(args);
            break;
    % endfor
    }
    dory_L2_memory_management();
```

**L2 memory allocation/deallocation**

# PULP-NN: Optimized Back-End



specification & dataset selection

training

quantization & pruning

graph optimization

memory-aware deployment

optimized DNN primitives

optimized HW & architecture

***PULP-NN***
*Parallel **U**LP*
***N**eural **N**etwork library*

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

↓

## HWC format



Channels

height

width

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

47

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

HWC format

```
lp.setup
  p.lw  w0, 4(W0!)
  p.lw  w1, 4(W1!)
  p.lw  w2, 4(W2!)
  p.lw  w3, 4(W3!)
  p.lw  x1, 4(X0!)
  p.lw  x2, 4(X1!)
  pv.sdotsp.b    acc1, w0, x0
  pv.sdotsp.b    acc2, w0, x1
  pv.sdotsp.b    acc3, w1, x0
  pv.sdotsp.b    acc4, w1, x1
  pv.sdotsp.b    acc5, w2, x0
  pv.sdotsp.b    acc6, w2, x1
  pv.sdotsp.b    acc7, w3, x0
  pv.sdotsp.b    acc8, w3, x1
end
```
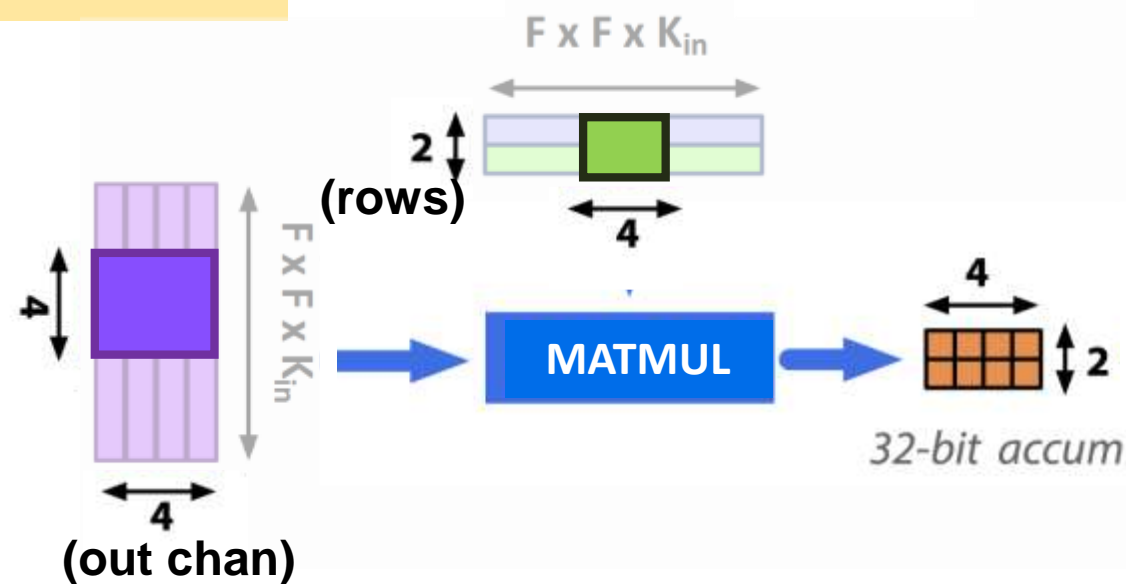
**Load 16 weights (8-bit)**
4 out chan, 4 in chan
address post-increment

$F \times F \times K_{in}$

MATMUL

32-bit accum

**(out chan)**

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

```
lp.setup
    p.lw   w0, 4(W0!)
    p.lw   w1, 4(W1!)
    p.lw   w2, 4(W2!)
    p.lw   w3, 4(W3!)
    p.lw   x1, 4(X0!)
    p.lw   x2, 4(X1!)
    pv.sdotsp.b      acc1, w0, x0
    pv.sdotsp.b      acc2, w0, x1
    pv.sdotsp.b      acc3, w1, x0
    pv.sdotsp.b      acc4, w1, x1
    pv.sdotsp.b      acc5, w2, x0
    pv.sdotsp.b      acc6, w2, x1
    pv.sdotsp.b      acc7, w3, x0
    pv.sdotsp.b      acc8, w3, x1
end
```
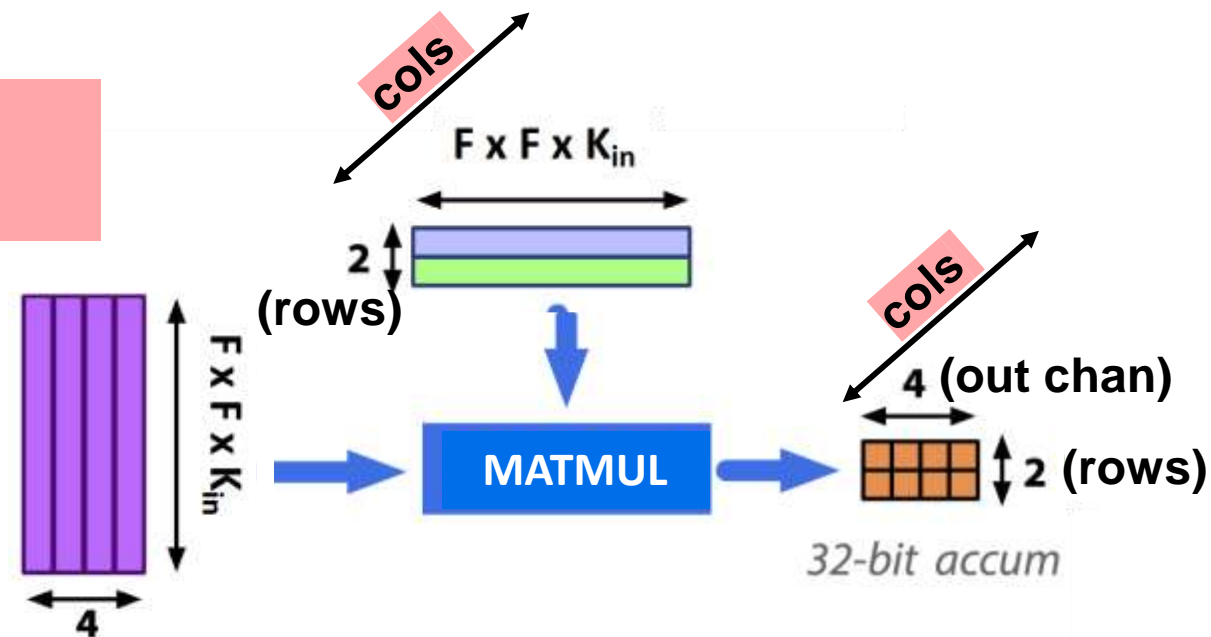
**Load 8 pixels**
2 rows, 4 in chan
address post-increment



$F \times F \times K_{in}$

**2 (rows)**

**4**

$F \times F \times K_{in}$

**4 (out chan)**

**MATMUL**

**4**

**2**

*32-bit accum*

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

```
lp.setup
    p.lw  w0, 4(W0!)
    p.lw  w1, 4(W1!)
    p.lw  w2, 4(W2!)
    p.lw  w3, 4(W3!)
    p.lw  x1, 4(X0!)
    p.lw  x2, 4(X1!)
    pv.sdotsp.b     acc1, w0, x0
    pv.sdotsp.b     acc2, w0, x1
    pv.sdotsp.b     acc3, w1, x0
    pv.sdotsp.b     acc4, w1, x1
    pv.sdotsp.b     acc5, w2, x0
    pv.sdotsp.b     acc6, w2, x1
    pv.sdotsp.b     acc7, w3, x0
    pv.sdotsp.b     acc8, w3, x1
end
```

**Load 8 pixels**
2 rows, 4 in chan
address post-increment



$F \times F \times K_{in}$

(rows)

MATMUL

32-bit accum

(out chan)

**Compute 32 MAC over 8 accumulators**
dot-product instructions

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

Loop over in chan, filter size

```
lp.setup
    p.lw   w0, 4(W0!)
    p.lw   w1, 4(W1!)
    p.lw   w2, 4(W2!)
    p.lw   w3, 4(W3!)
    p.lw   x1, 4(X0!)
    p.lw   x2, 4(X1!)
    pv.sdotsp.b    acc1, w0, x0
    pv.sdotsp.b    acc2, w0, x1
    pv.sdotsp.b    acc3, w1, x0
    pv.sdotsp.b    acc4, w1, x1
    pv.sdotsp.b    acc5, w2, x0
    pv.sdotsp.b    acc6, w2, x1
    pv.sdotsp.b    acc7, w3, x0
    pv.sdotsp.b    acc8, w3, x1
end
```



$F \times F \times K_{in}$

2 (rows)

4

$F \times F \times K_{in}$

4

4

MATMUL

4

2

32-bit accum

(out chan)

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

# PULP-NN: Optimized Back-End

Target **int8** execution of CONV, FC, ... primitives
1) maximize **data reuse in register file** 2) improve **kernel regularity** 3) exploit **parallelism**

```
lp.setup
  p.lw  w0, 4(W0!)
  p.lw  w1, 4(W1!)
  p.lw  w2, 4(W2!)
  p.lw  w3, 4(W3!)
  p.lw  x1, 4(X0!)
  p.lw  x2, 4(X1!)
  pv.sdotsp.b    acc1, w0, x0
  pv.sdotsp.b    acc2, w0, x1
  pv.sdotsp.b    acc3, w1, x0
  pv.sdotsp.b    acc4, w1, x1
  pv.sdotsp.b    acc5, w2, x0
  pv.sdotsp.b    acc6, w2, x1
  pv.sdotsp.b    acc7, w3, x0
  pv.sdotsp.b    acc8, w3, x1
end
```

**Parallelize over 8 cores**
column dimension



cols

$F \times F \times K_{in}$

2
(rows)

$F \times F \times K_{in}$

4

MATMUL

cols

4 (out chan)

2 (rows)

*32-bit accum*

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

# PULP-NN: Layers Supported (@ 25-2-2021)

**Convolutions**
- Conv_Ho_parallel (+bn, +Relu)
- Conv_HoWo_parallel (+bn, +Relu)
- Conv_Co_parallel (+bn, +Relu)

**Point-wise Convolutions**
- Pointwise_Ho_parallel (+bn, +Relu)
- Pointwise_HoWo_parallel (+bn, +Relu)
- Pointwise_Co_parallel (+bn, +Relu)

**Depth-wise Convolutions**
- Depthwise_3x3s1 (+bn, +Relu)
- Depthwise_generic (+bn, +Relu)

**Linear Layers**
- Linear (+bn, +Relu)
- Linear_out_fp32

**Other Layers**
- Add (+bn, +Relu)
- Avgpool
- Maxpool

**PULP-NN** [Garofalo 19] https://arxiv.org/abs/1908.11263

**https://github.com/pulp-platform/pulp-nn**

# Requirements – DORY + PULP-NN

- **DORY is available at https://github.com/pulp-platform/dory**
- **On Ubuntu 18.04 you need the following packages and tools:**
  - **python>=3.6** or **python3.5 with future-fstrings** package
  - **pulp-sdk** available at **https://github.com/pulp-platform/pulp-sdk**
  - Python packages:
    - **onnx>=1.8.1**
    - **torch>=1.5.1**
    - **pandas>=0.24.2**
    - **ortools>=8.0.8283**

- No installation required for DORY and PULP-NN

**https://github.com/pulp-platform/pulp-nn**

# Network Generation



**Integer Network + tile sizes**

```
Code Generation
from templates
```

**Network-level C code**
- L3/L2 transfer boilerplate
- double buffering for weights
- calls to layer-level code

**Layer-level C code**
- L2/L1 transfer boilerplate
- calls to PULP-NN backend library

**NEMO**
**Post-training Tutorial:**
https://github.com/pulp-platform/nemo

**DORY**
**Tutorial:**
https://github.com/pulp-platform/dory_examples

**Full stack tutorial in the SDK documentation**
https://github.com/pulp-platform/pulp-sdk

# Generate a neural network with default settings

- **Generate the default network**

```
pulp-user@pulp-box /pulp/dory $ cd dory_examples/
pulp-user@pulp-box /pulp/dory/dory_examples $ python3 network_generate.py
Creating annotated graph in Network_annotated_graph.log
Creating tiling profiling in Tiling_profiling.log
pulp-user@pulp-box /pulp/dory/dory_examples $
```

# Generate a neural network with default settings

- **Generate the default network**

```
pulp-user@pulp-box /pulp/dory $ cd dory_examples/
pulp-user@pulp-box /pulp/dory/dory_examples $ python3 network_generate.py
Creating annotated graph in Network_annotated_graph.log
Creating tiling profiling in Tiling_profiling.log
pulp-user@pulp-box /pulp/dory/dory_examples $
```

- **Inspect the two output files**

Network_annotated_graph

Tiling profiling

# Generate a neural network with default settings

- **Generate the default network**

- **Inspect the two output files**

Network_annotated_graph

Tiling profiling

L2-L1 tiling

L3-L2 tiling +
L2-L1 tiling

# Generate a neural network with default settings

- ## Run the network on pulp gvsoc

```
pulp-user@pulp-box /pulp/dory/dory_examples $ cd application/
pulp-user@pulp-box /pulp/dory/dory_examples/application $ make clean all run CORE=8 platform=gvsoc
```

**Weights checksum**

```
L3 Buffer alloc initial @ 3388608:      Ok

L3 Buffer alloc initial @ 2388608:      Ok

L3 Buffer alloc initial @ 1388608:      Ok
Layer 0  : Checksum = 120996     , FLASH 120996     , Check OK
Layer 1  : Checksum = 53504      , FLASH 53504      , Check OK
Layer 2  : Checksum = 303730     , FLASH 303730     , Check OK
Layer 3  : Checksum = 117611     , FLASH 117611     , Check OK
Layer 4  : Checksum = 1103930    , FLASH 1103930    , Check OK
Layer 5  : Checksum = 234609     , FLASH 234609     , Check OK
Layer 6  : Checksum = 2236755    , FLASH 2236755    , Check OK
Layer 7  : Checksum = 204746     , FLASH 204746     , Check OK
Layer 8  : Checksum = 4266224    , FLASH 4266224    , Check OK
Layer 9  : Checksum = 524077     , FLASH 524077     , Check OK
Layer 10 : Checksum = 8552435    , FLASH 8552435    , Check OK
Layer 11 : Checksum = 451595     , FLASH 451595     , Check OK
Layer 12 : Checksum = 17038175   , FLASH 17038175   , Check OK
Layer 13 : Checksum = 1057921    , FLASH 1057921    , Check OK
Layer 14 : Checksum = 33824230   , FLASH 33824230   , Check OK
```

**Activations checksum**

```
Layer 3 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
Layer 4 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
Layer 5 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
Layer 6 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
Layer 7 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
Layer 8 ended
Checksum in/out Layer : Ok
```

**Performance**

```
[0] : num_cycles: 20541790
[0] : MACs: 186400768
[0] : MAC/cycle: 9.074223
[0] : n. of Cores: 8
```

# Change default settings

- **Set of arguments that you can pass to DORY**

```
parser = argparse.ArgumentParser(formatter_class=RawTextHelpFormatter)
parser.add_argument('--network_dir', default = "./examples/MobilenetV1/", help = 'directory of the onnx file of the network')
parser.add_argument('--l1_buffer_size', type=int, default = 38000, help = 'L1 buffer size. IT DOES NOT INCLUDE SPACE FOR STACKS.')
parser.add_argument('--l2_buffer_size', type=int, default = 380000, help = 'L2 buffer size.')
parser.add_argument('--master_stack', type=int, default = 4096, help = 'Cluster Core 0 stack')
parser.add_argument('--slave_stack', type=int, default = 3072, help = 'Cluster Core 1-7 stack')
parser.add_argument('--Bn_Relu_Bits', type=int, default = 32, help = 'Number of bits for Relu/BN')
parser.add_argument('--perf_layer', default = 'No', help = 'Yes: MAC/cycles per layer. No: No perf per layer.')
parser.add_argument('--verbose_level', default = 'Check_all+Perf_final', help = "None: No_printf.\nPerf_final: only total performance\nCheck
parser.add_argument('--chip', default = 'GAP8v3', help = 'GAP8v2 for fixing DMA issue. GAP8v3 otherise')
parser.add_argument('--sdk', default = 'pulp_sdk', help = 'gap_sdk or pulp_sdk')
parser.add_argument('--dma_parallelization', default = '8-cores', help = '8-cores or 1-core')
parser.add_argument('--fc_frequency', default = 100000000, help = 'frequency of fabric controller')
parser.add_argument('--cl_frequency', default = 100000000, help = 'frequency of cluster')
parser.add_argument('--optional', default = '8bit', help = '8bit, mixed, 1D_Conv')
args = parser.parse_args()
```

# Change default settings

- ## Enable layer performance verbose

```
pulp-user@pulp-box /pulp/dory/dory_examples $ python3 network_generate.py --perf_layer Yes
```

```
[0] Layer 3  : num_cycles: 814176      , MACs: 589824     , MAC/cycle: 0.724443, n. of Cores: 8
Layer 3 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
[0] Layer 4  : num_cycles: 757333      , MACs: 8388608    , MAC/cycle: 11.076512, n. of Cores: 8
Layer 4 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
[0] Layer 5  : num_cycles: 643655      , MACs: 1179648    , MAC/cycle: 1.832733, n. of Cores: 8
Layer 5 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
[0] Layer 6  : num_cycles: 1263171     , MACs: 16777216   , MAC/cycle: 13.281825, n. of Cores: 8
Layer 6 ended
Checksum in/out Layer : Ok
Checksum in/out Layer : Ok
[0] Layer 7  : num_cycles: 328293      , MACs: 294912     , MAC/cycle: 0.898319, n. of Cores: 8
Layer 7 ended
Checksum in/out Layer : Ok
```

- ## Change L1 maximum memory footprint

```
pulp-user@pulp-box /pulp/dory/dory_examples $ python3 network_generate.py --network_dir ./examples/MobilenetV1/ --l1_buffer_size 30000
```

```
[0] : num_cycles: 25922049
[0] : MACs: 186400768
[0] : MAC/cycle: 7.190819
[0] : n. of Cores: 8
```

```
pulp-user@pulp-box /pulp/dory/dory_examples $ python3 network_generate.py --network_dir ./examples/MobilenetV1/ --l1_buffer_size 10000
Creating annotated graph in Network_annotated_graph.log
Creating tiling profiling in Tiling_profiling.log
  Conv2d ERROR: no L2-L1 tiling found. Exiting...
```

- ## Generate a new network

```
pulp-user@pulp-box /pulp/dory/dory_examples $ python3 network_generate.py --network_dir ./examples/PenguiNet_32/
```

# Thanks for the attention