



# WIP: Automatic DNN Deployment on Heterogeneous Platforms: the GAP9 Case Study

Energy-Efficient Embedded Systems Laboratory (UNIBO)

**Luka Macan** <[luka.macan@unibo.it](mailto:luka.macan@unibo.it)>

**Alessio Burrello** <[alessio.burrello@polito.it](mailto:alessio.burrello@polito.it)>

**Luca Benini** <[luca.benini@unibo.it](mailto:luca.benini@unibo.it)>

**Francesco Conti** <[f.conti@unibo.it](mailto:f.conti@unibo.it)>

**PULP Platform**

Open Source Hardware, the way it should be!



@pulp\_platform

[pulp-platform.org](http://pulp-platform.org)

[youtube.com/pulp\\_platform](https://youtube.com/pulp_platform)

# Goal: NN Deployment on the Edge



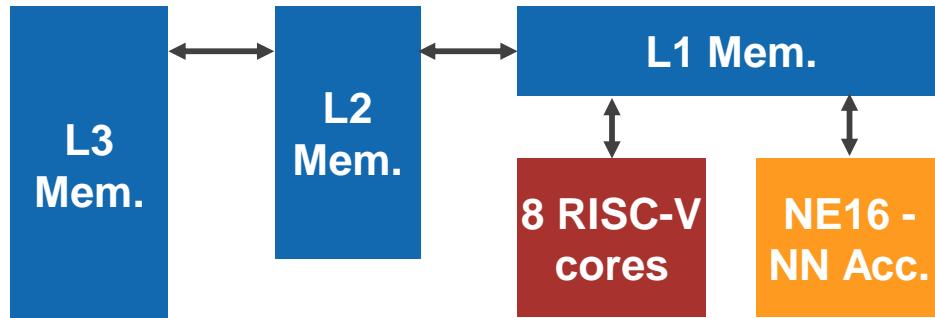
## Our target: NN accelerated embedded edge devices

- Heterogenous
- Ultra-low-power
- Software managed data "caches"

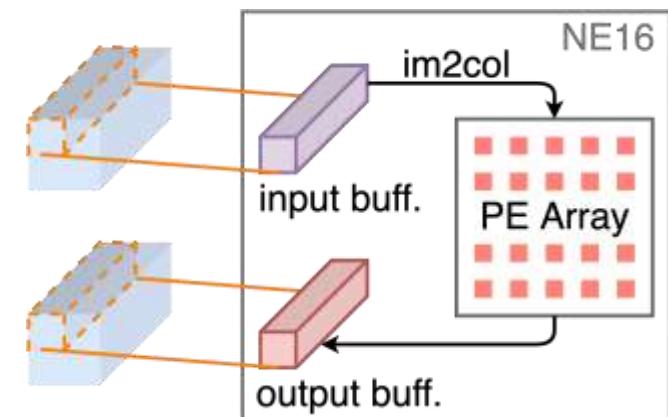


### Case study: GAP9 SoC

- Compute: Cluster of 8 RISC-V cores + NN Accelerator
- 0.33 mW/GOP<sup>1</sup>
- 3 levels of memory + DMAs



*Simplified system view of the GAP9 SoC*

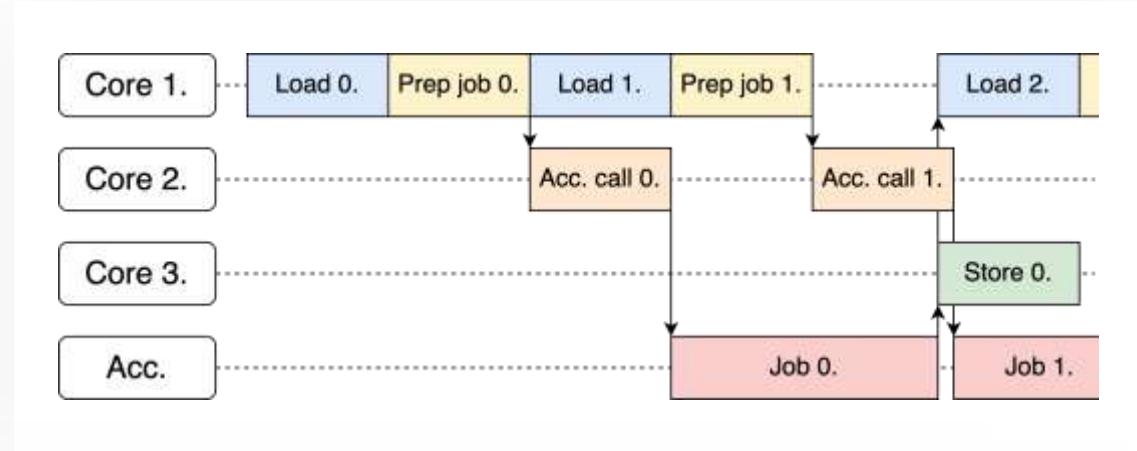


*NE16 Accelerator high-level view*

# Inter-layer optimizations

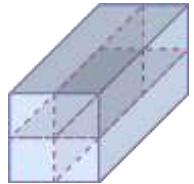


- **Software pipelining & Task parallelization**

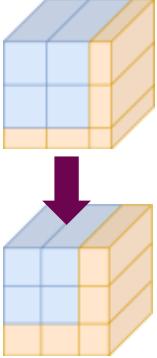


- **Tiling heuristics**

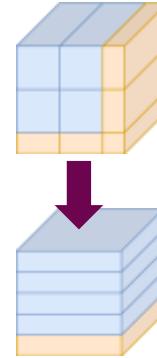
1. Tiles divisible  
with input buffer



2. Balanced tiles



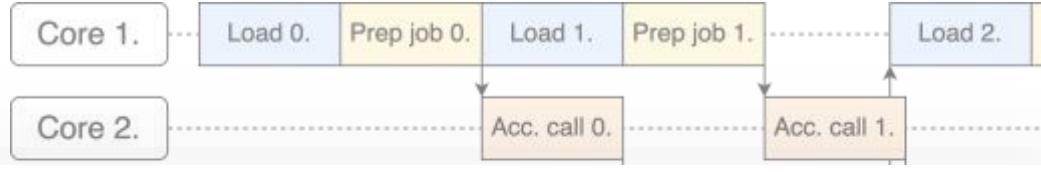
3. Full dimensions



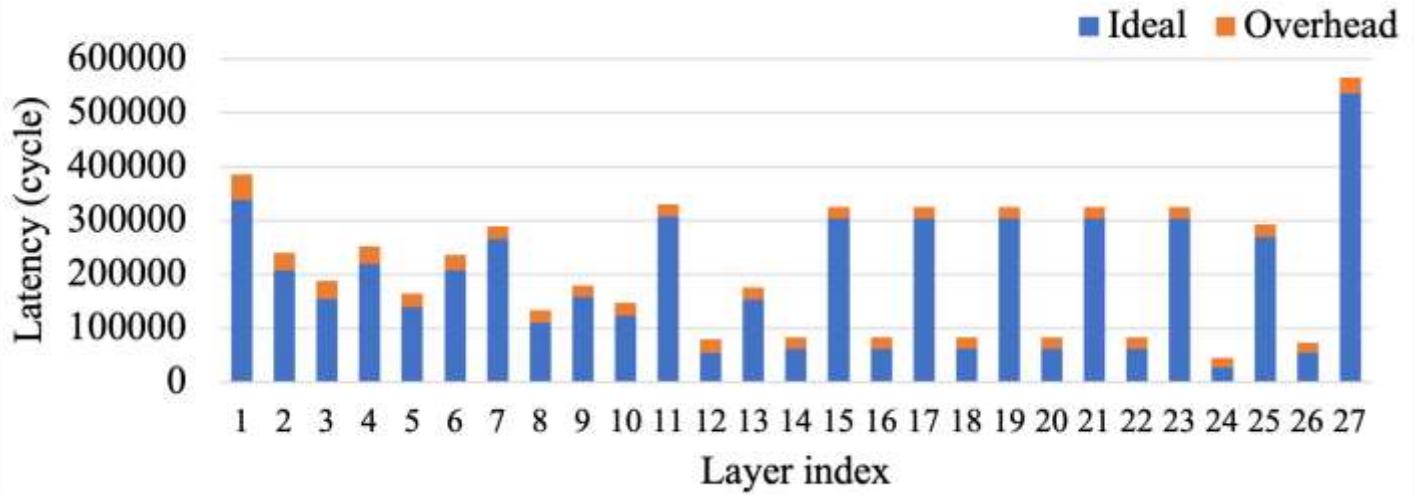
# Inter-layer optimizations



- Software pipelining & Task parallel



- Tiling heuris

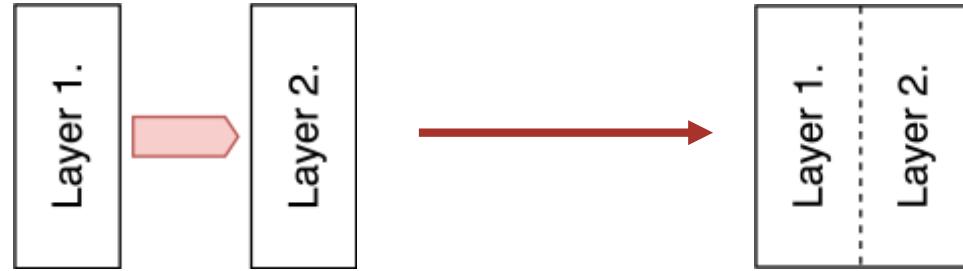


**12% overall latency overhead over ideal**

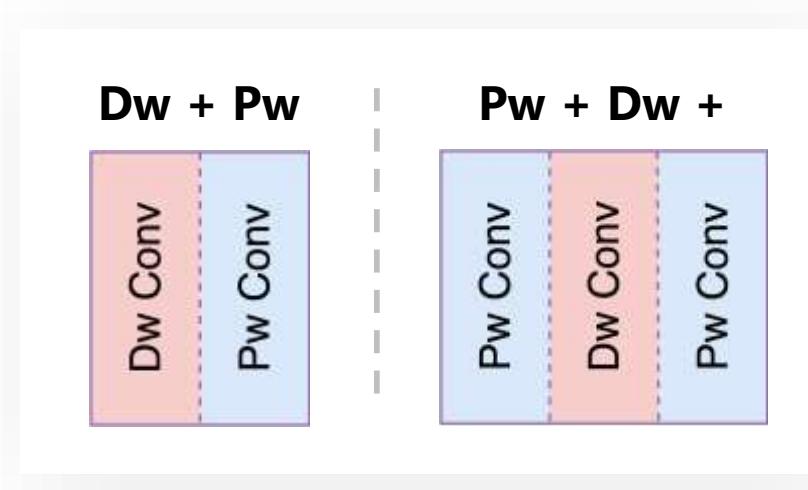
# Intra-layer optimizations



- Layer fusion



- Choice of layers to fuse



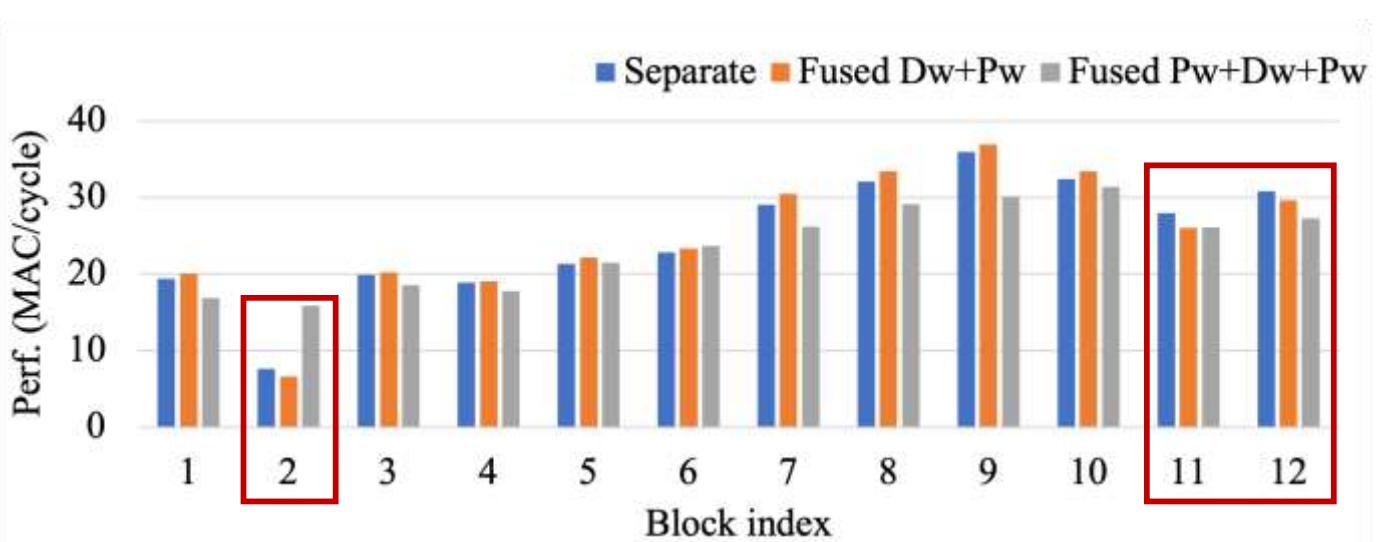
# Intra-layer optimizations



- Layer fusion



- Choice of layer



*MN-V2 unique Inverted Bottleneck blocks*

# Conclusion



- Achieved high accelerator utilization (~90%) and 3.44x performance improvement over the predecessor, GAP8
- Accelerators need several techniques to be applied for them to come close to ideal utilization
- Automatic deployment tools will play a more and more critical role as the hardware gets more complex