



# ExecuTorch on Small Bare-Metal Microcontrollers

Shane Slattery  
Pietra Ferreira  
William Jones  
Jeremy Bennett

Copyright © 2026 Embecosm. Freely available under a  
Creative Commons Attribution-ShareAlike license.



# The Cost of AI

50

# The Cost of AI

50

0.03

# The Cost of AI

50

0.03

10-100+

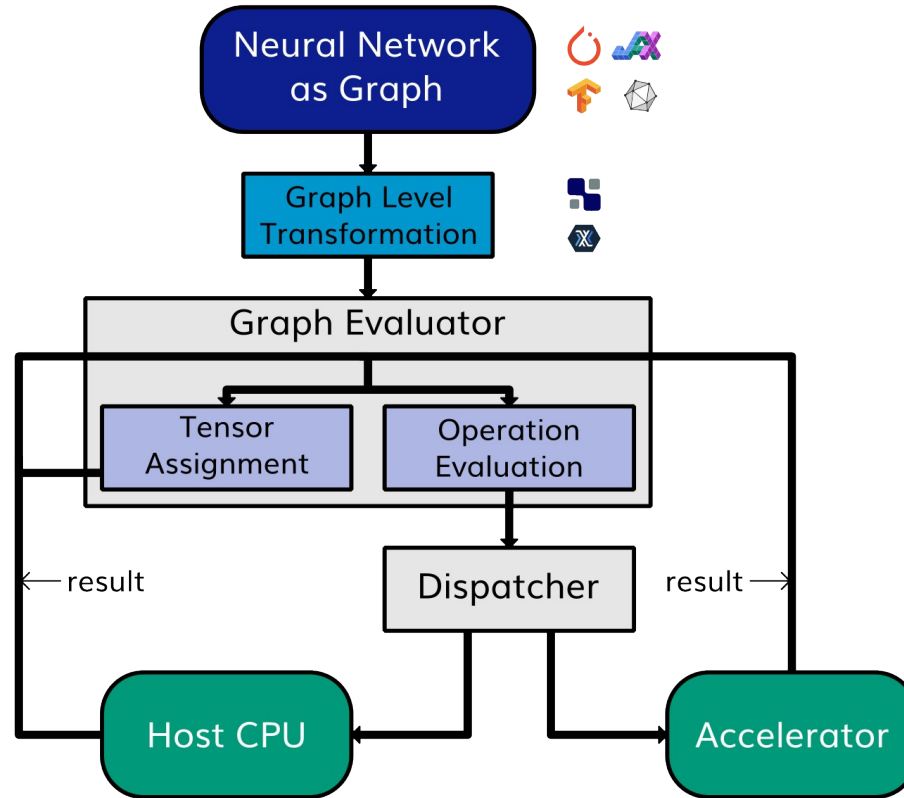


# Building AI at the Edge

Copyright © 2026 Embecosm. Freely available under a  
Creative Commons Attribution-ShareAlike license.



# What are PyTorch and ExecuTorch?



# What are PyTorch and ExecuTorch?



# What are PyTorch and ExecuTorch?

## PyTorch

- Generic ML framework
- Python based
- Open source





# What are PyTorch and ExecuTorch?

## PyTorch

- Generic ML framework
- Python based
- Open source

Training **AND** Inference



# What are PyTorch and ExecuTorch?

## PyTorch

- Generic ML framework
- Python based
- Open source

Training **AND** Inference



## ExecuTorch

- PyTorch extension
- Resource limited

# What are PyTorch and ExecuTorch?

## PyTorch

- Generic ML framework
- Python based
- Open source

Training **AND** Inference



## ExecuTorch

- PyTorch extension
- Resource limited

Inference **ONLY**

# Two Steps to ExecuTorch Development

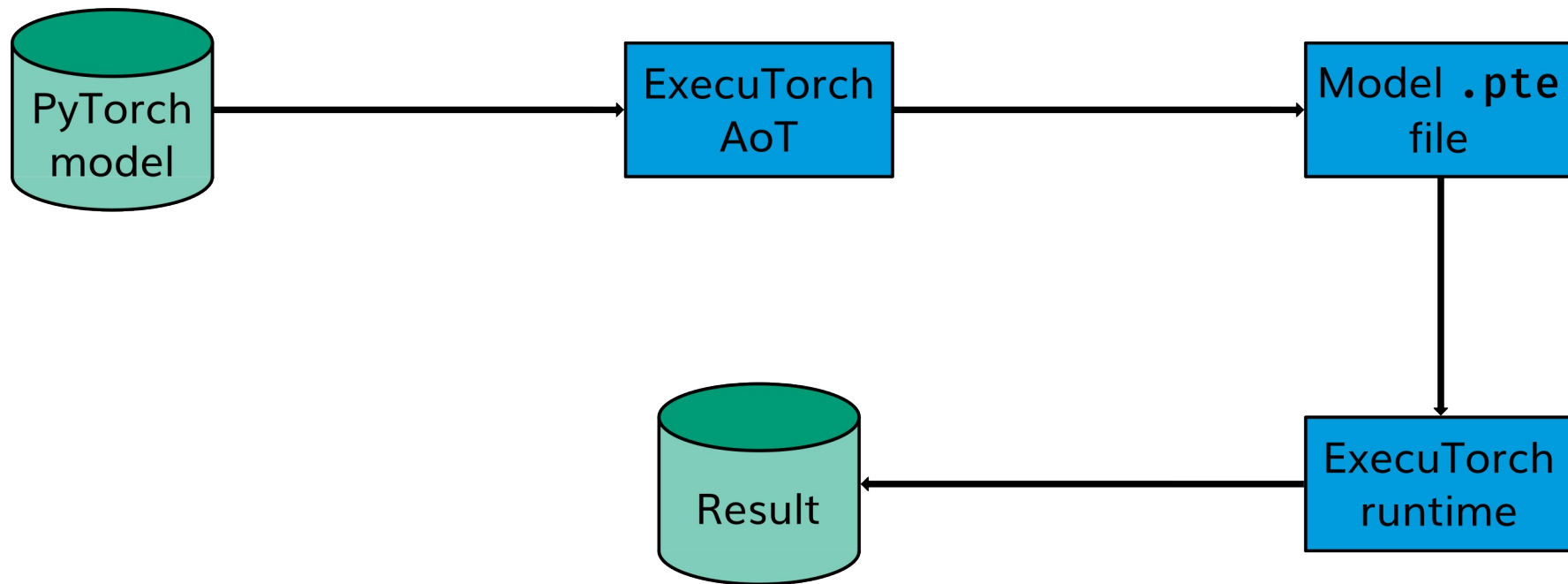
## 1. Build ExecuTorch

# Two Steps to ExecuTorch Development

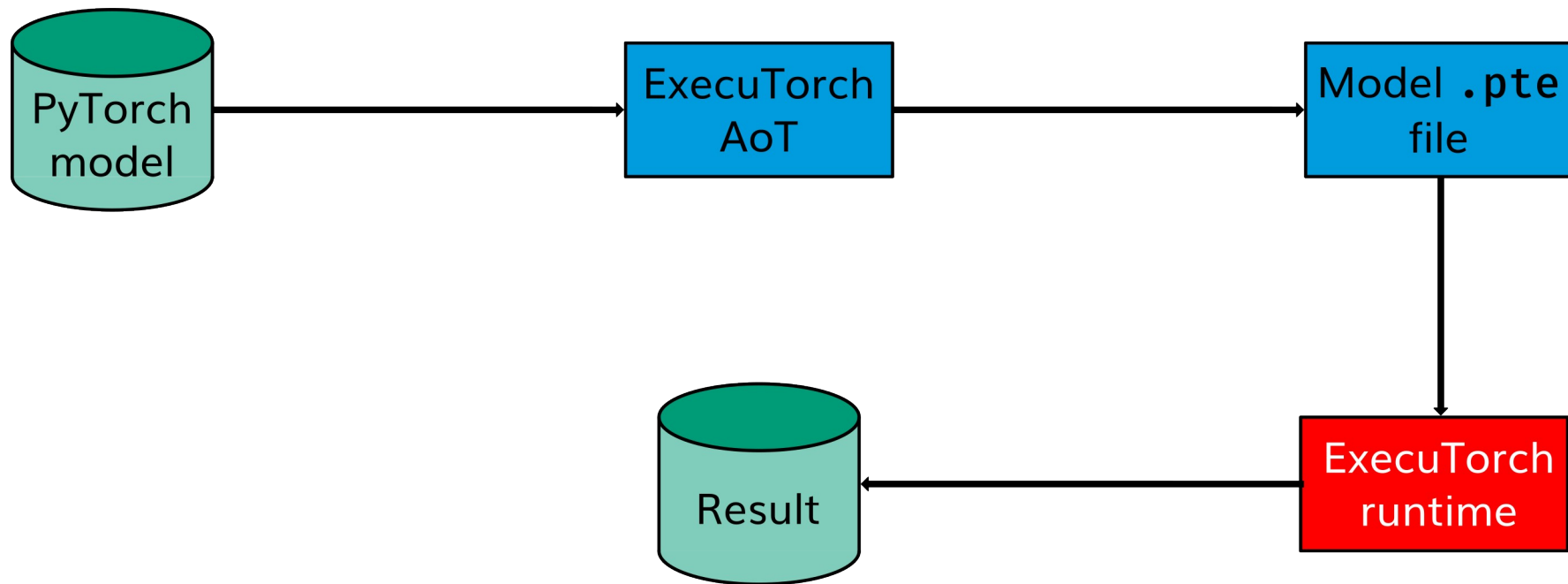
1.Build ExecuTorch

2.Customize Performance

# Building With ExecuTorch



# Building With ExecuTorch



# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory



# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory

## Compute

- few cores
- few accelerators

# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory

## Compute

- few cores
- few accelerators

No OS

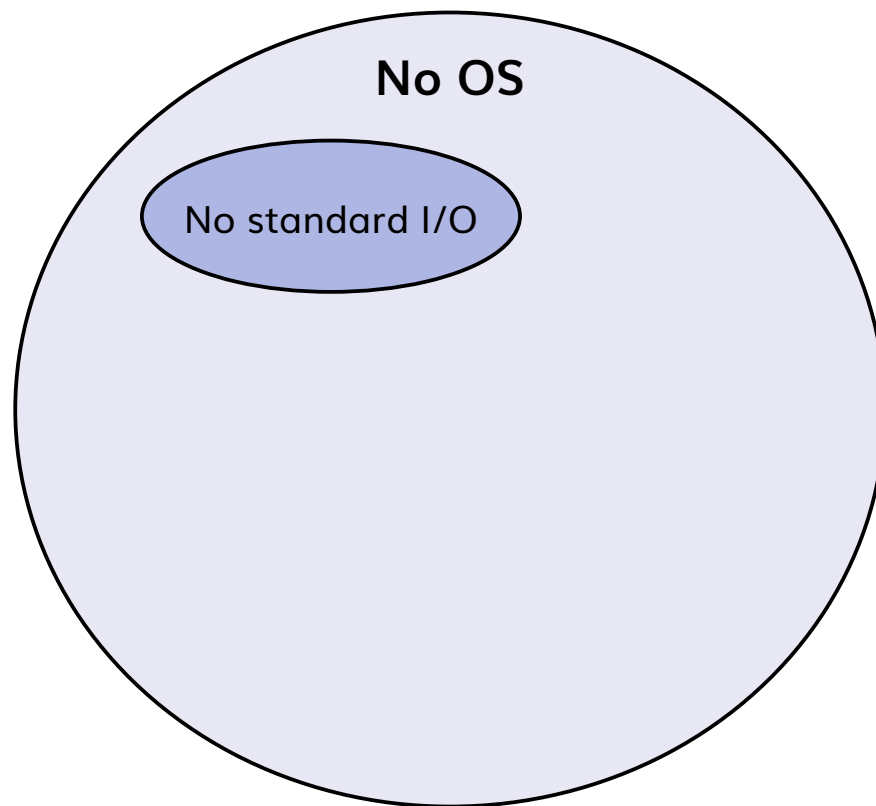
# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory

## Compute

- few cores
- few accelerators



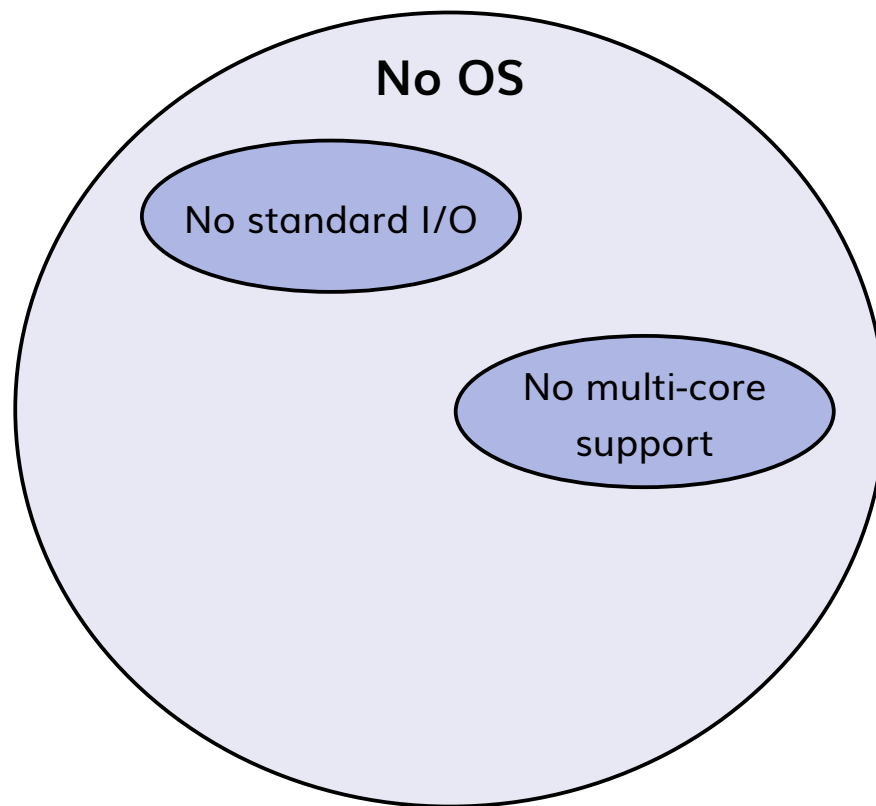
# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory

## Compute

- few cores
- few accelerators



# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory

## Compute

- few cores
- few accelerators

## No OS

No standard I/O

No multi-core  
support

No run-time  
code loading

# Our Experience: A RISC-V Platform

## Memory

- a few Megabytes
- little fast memory

## Compute

- few cores
- few accelerators

## No OS

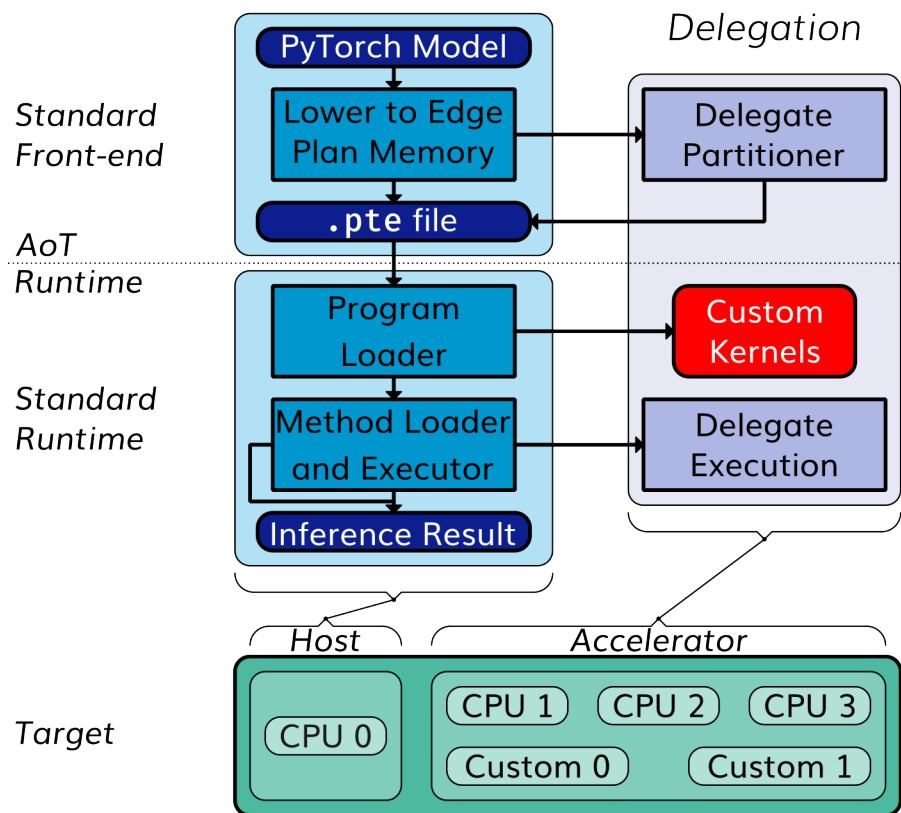
No standard I/O

No multi-core  
support

No run-time  
code loading

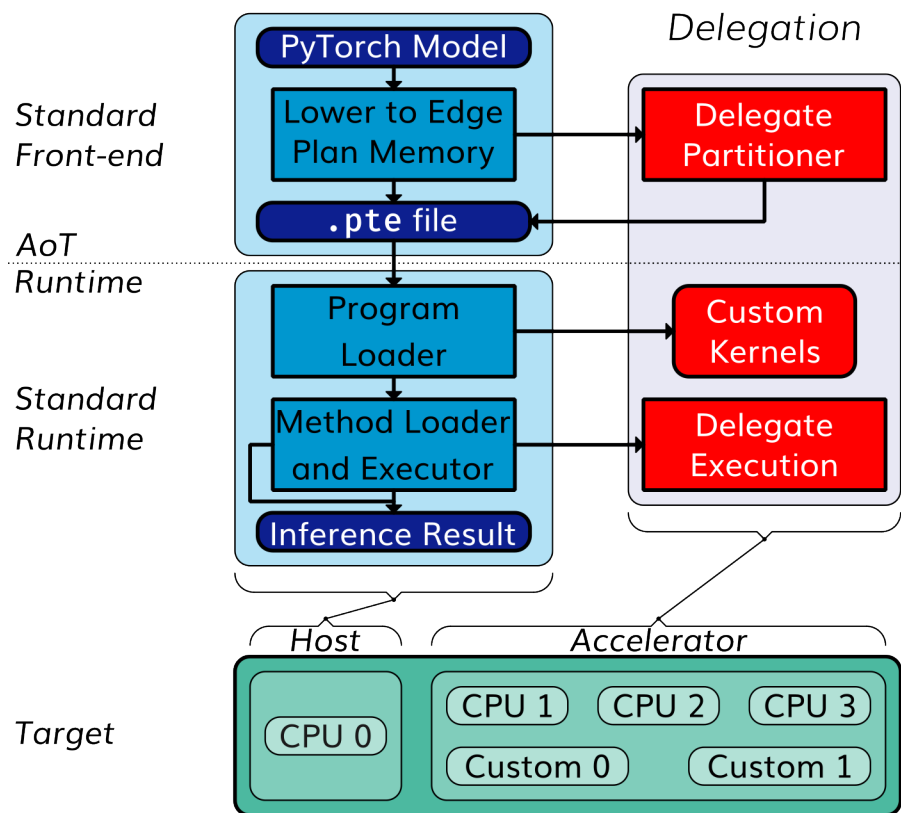
Success with minor changes

# Customize Performance



Existing functions

# Customize Performance

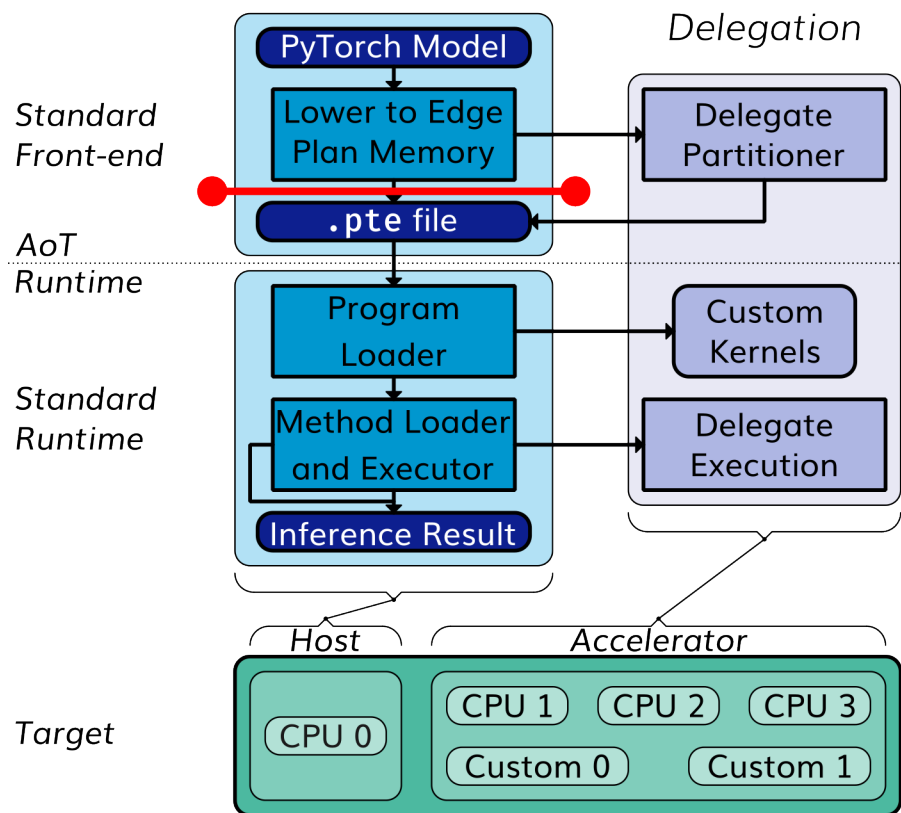


Existing functions

New capabilities



# Customize Performance



Existing functions

New capabilities

Graph level optimization

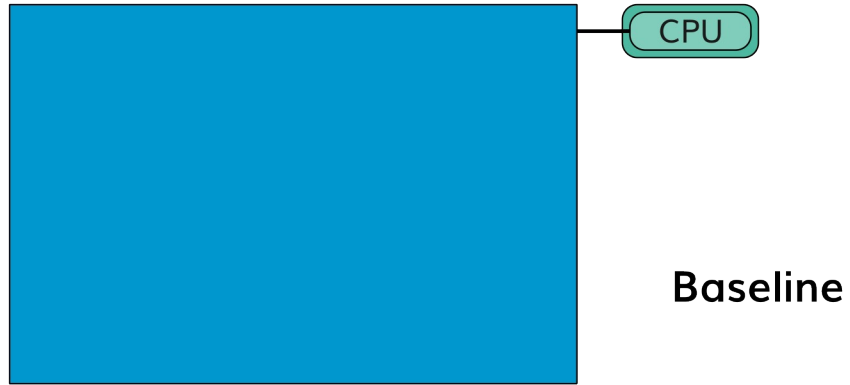


# Case Study: A RISC-V Processor with Custom NPU

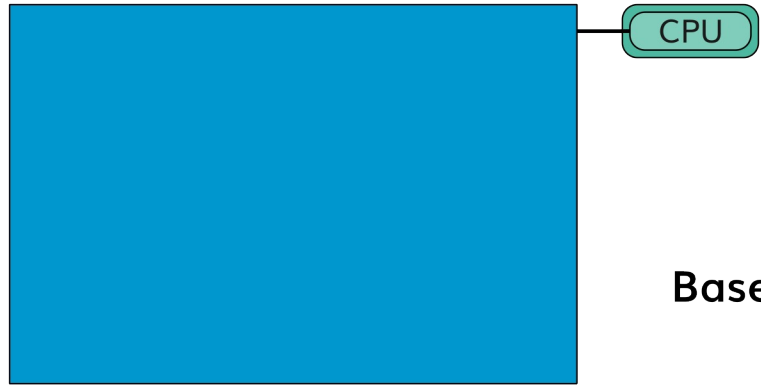
Copyright © 2026 Embecosm. Freely available under a  
Creative Commons Attribution-ShareAlike license.



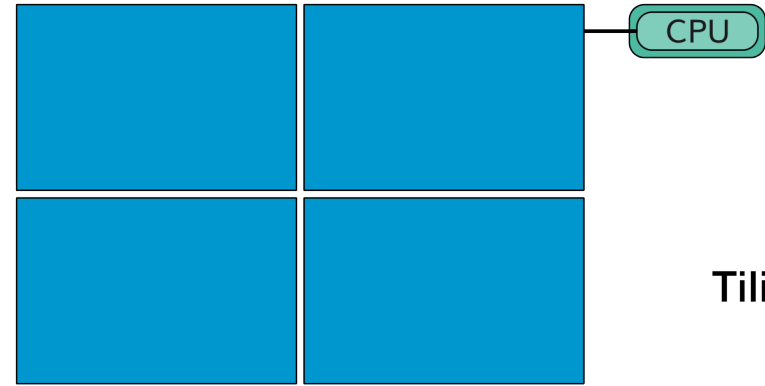
# Optimization Strategies



# Optimization Strategies

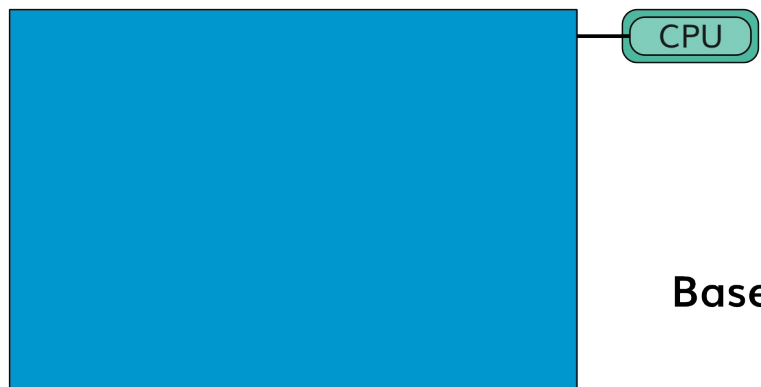


Baseline

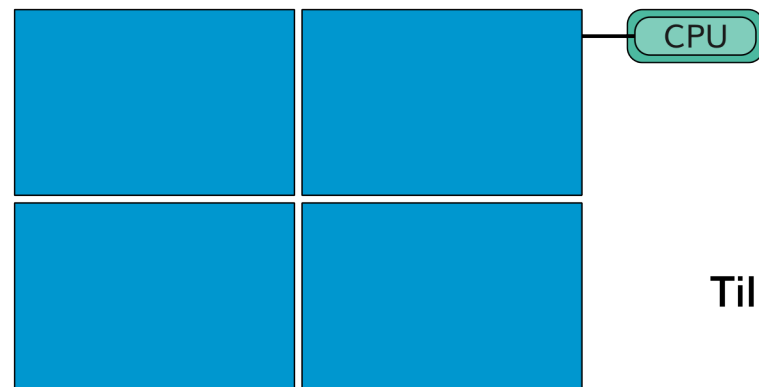


Tiling

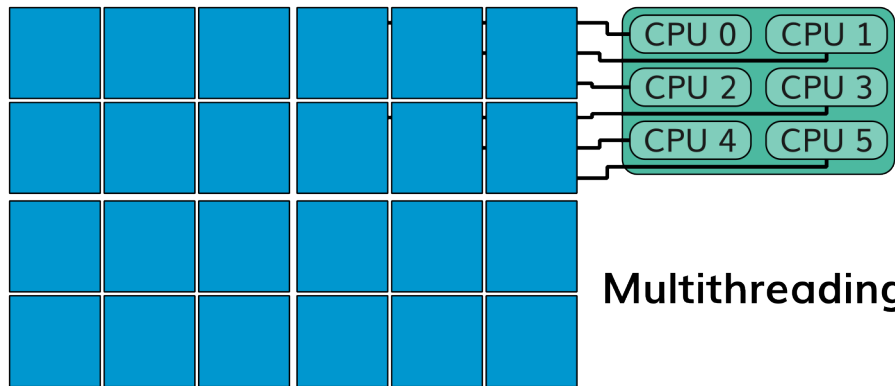
# Optimization Strategies



Baseline

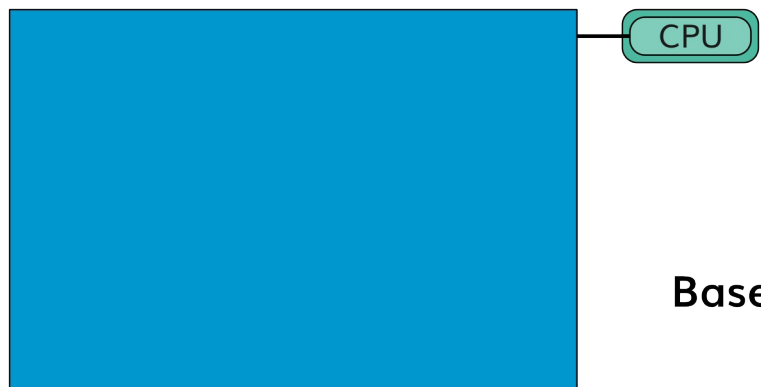


Tiling

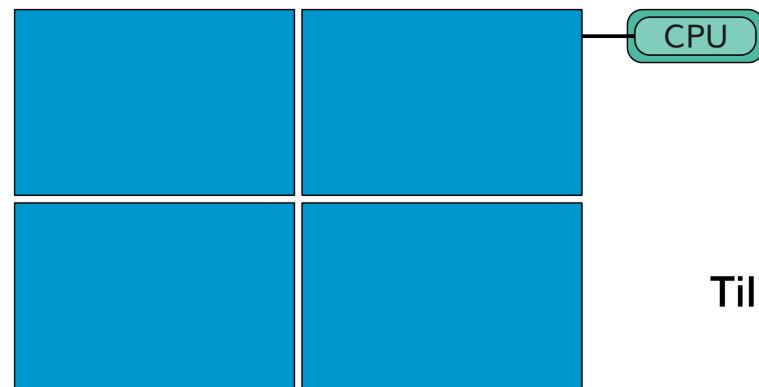


Multithreading

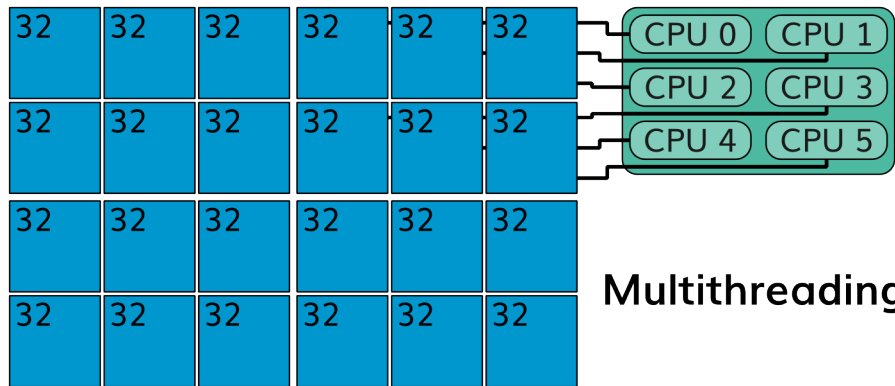
# Optimization Strategies



Baseline



Tiling

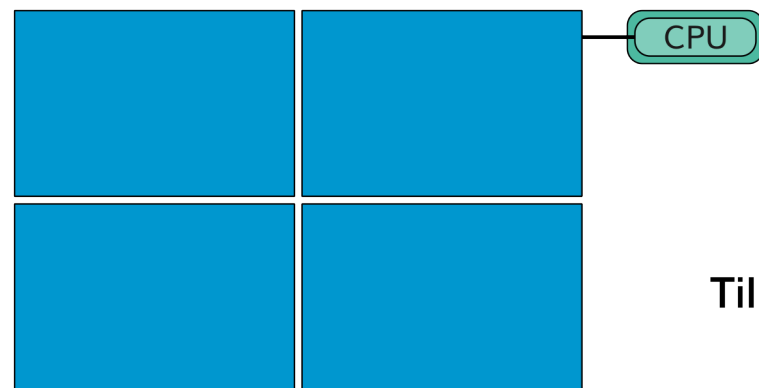


Multithreading

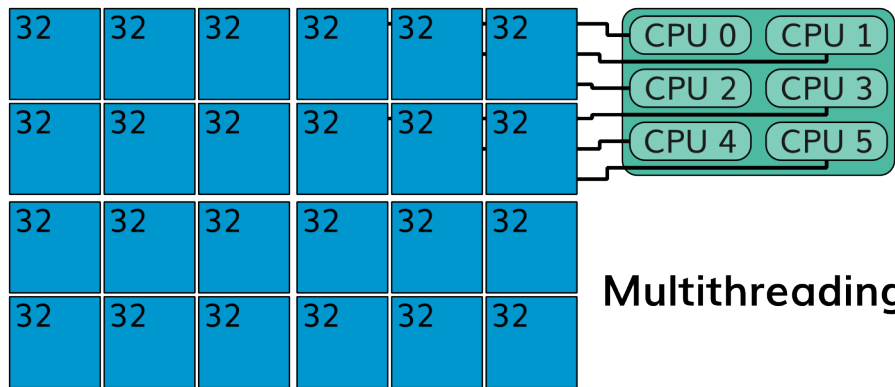
# Optimization Strategies



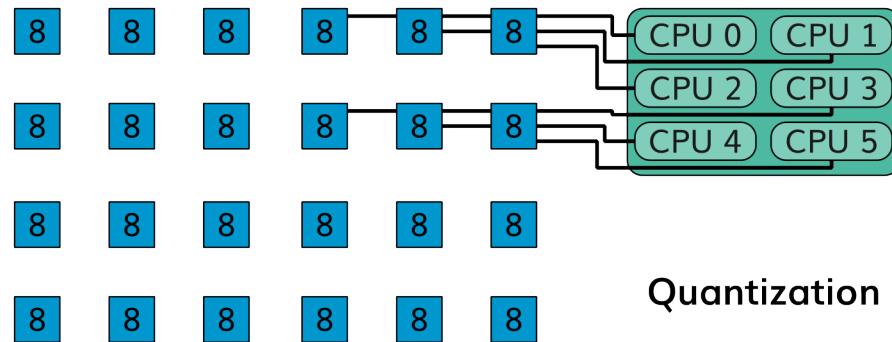
Baseline



Tiling

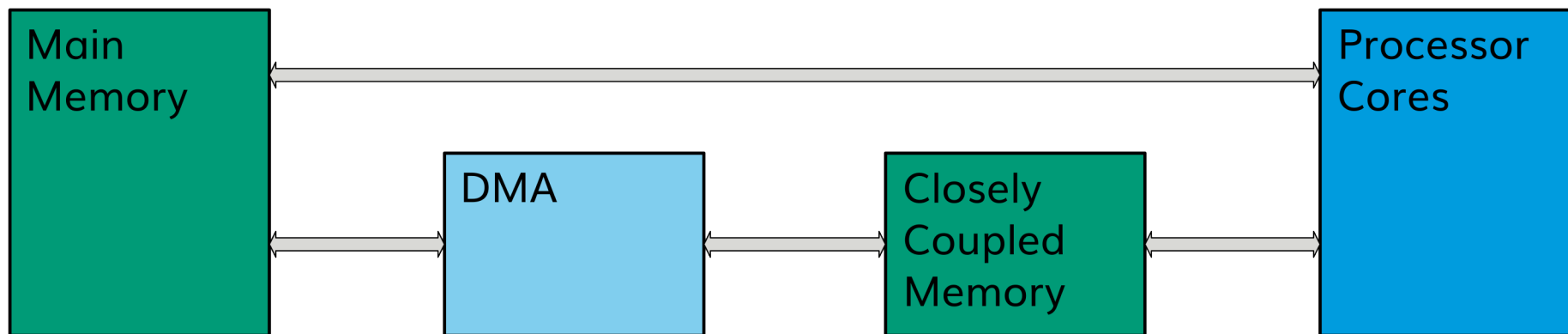


Multithreading



Quantization

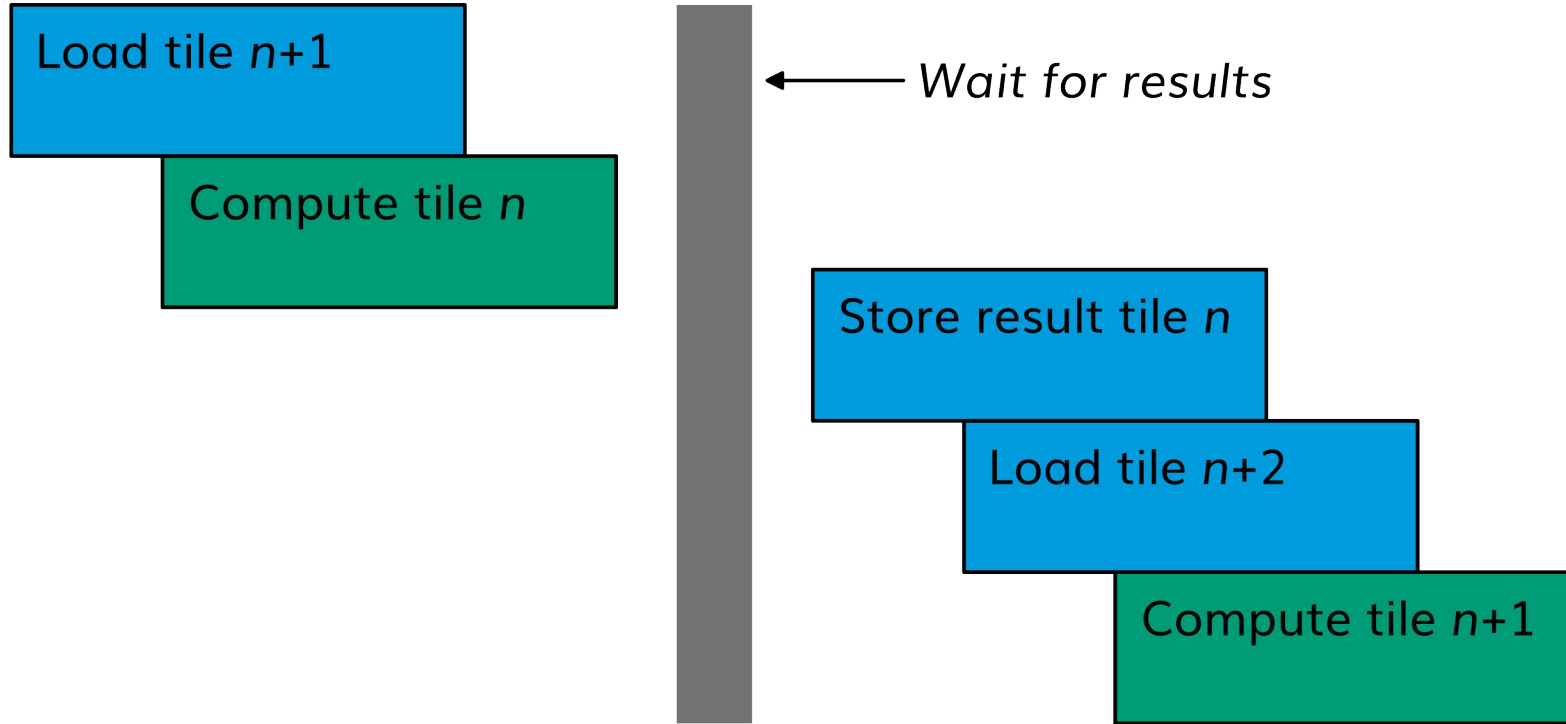
# Optimization in Depth: Memory



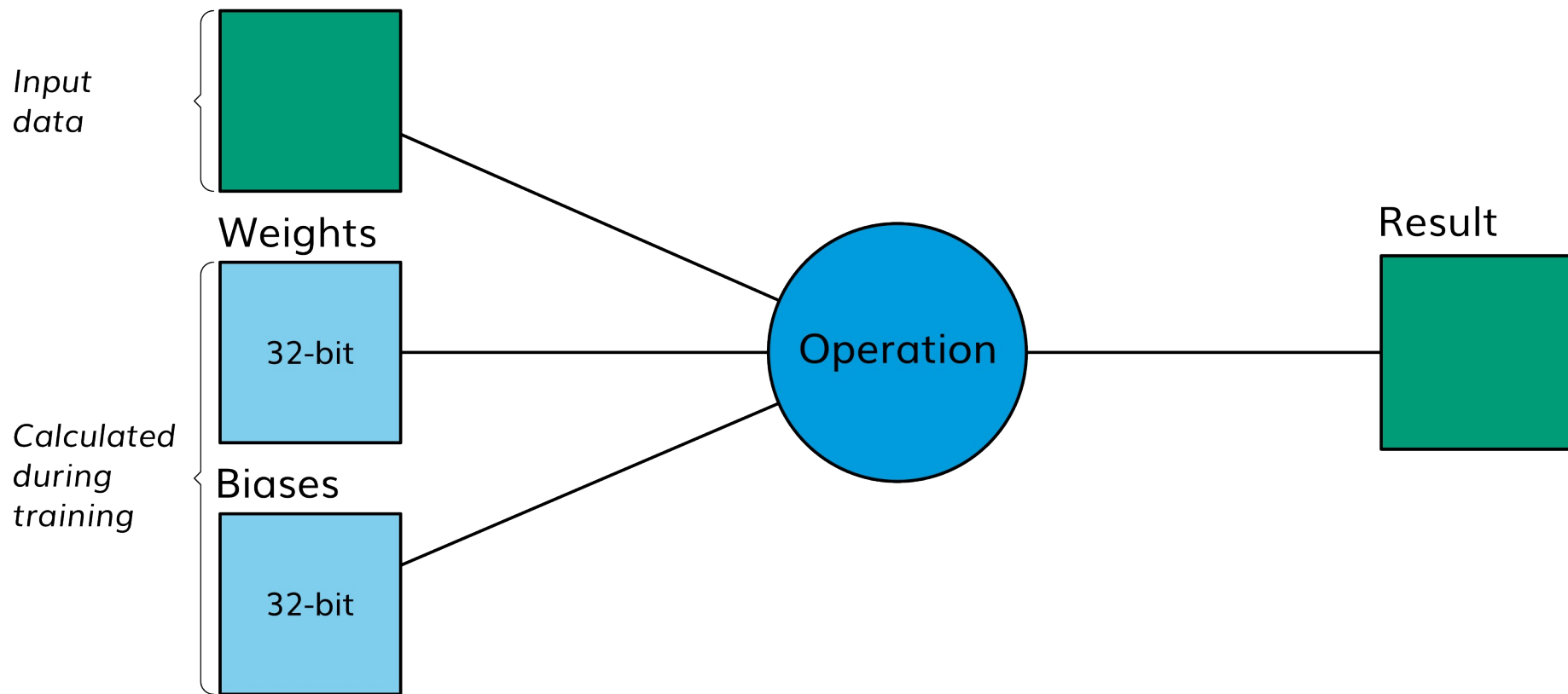
- Cache frequently used data
  - e.g. initial layers
- Break data into "tiles"
  - e.g. sub-areas of images
- Algorithms are predictable
  - regular loops
  - pre-fetch data
  - e.g. image convolution



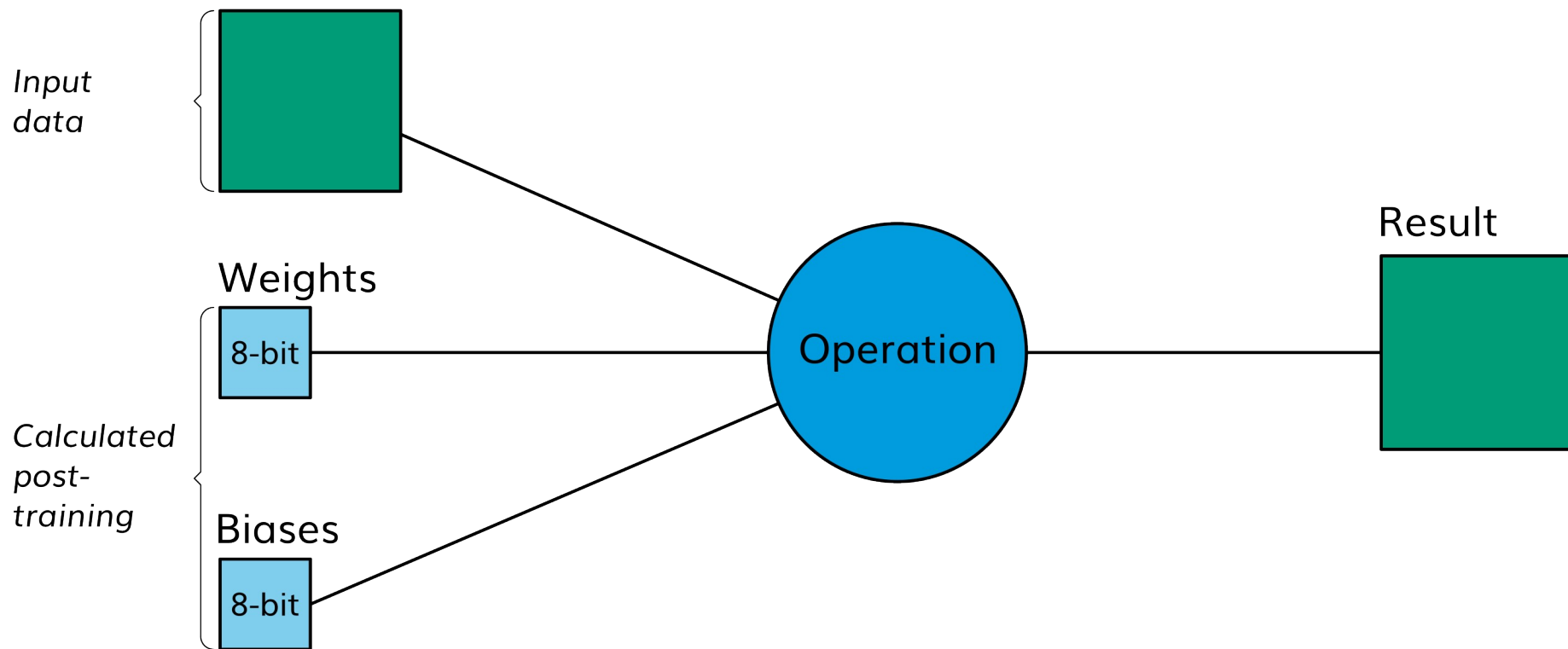
# Optimization in Depth: DMA



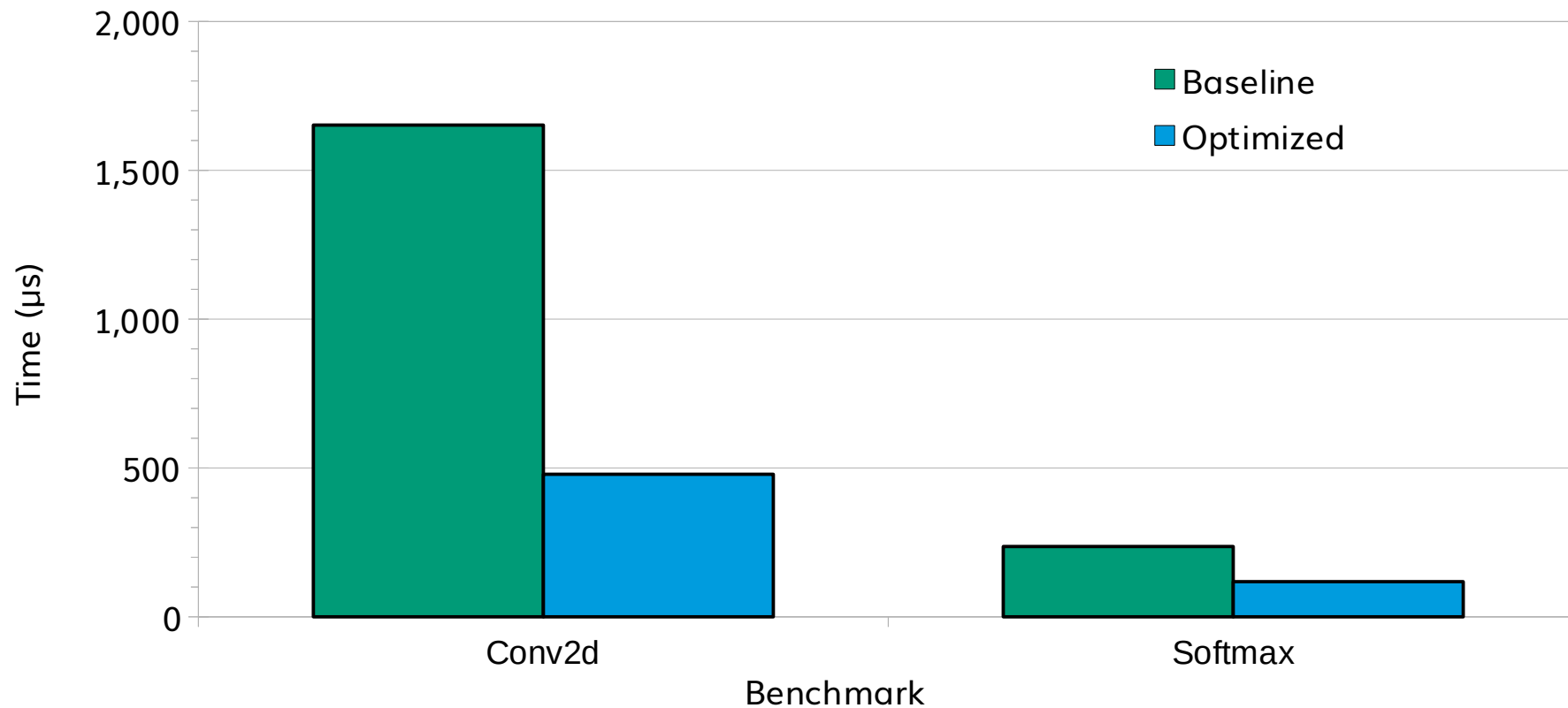
# Optimization in Depth: Quantization



# Optimization in Depth: Quantization



# The Performance Benefit



# Acknowledgement

The Embecosm team would like to thank our colleagues at Mosaic SoC in Switzerland who have supported the work presented here

mosaic



# Thank You

shane.slattery@embecoscsm.com

pietra.ferreira@embecoscsm.com

william.jones@embecoscsm.com

jeremy.bennett@embecoscsm.com



Shane Slattery

Pietra Ferreira

William Jones

Jeremy Bennett