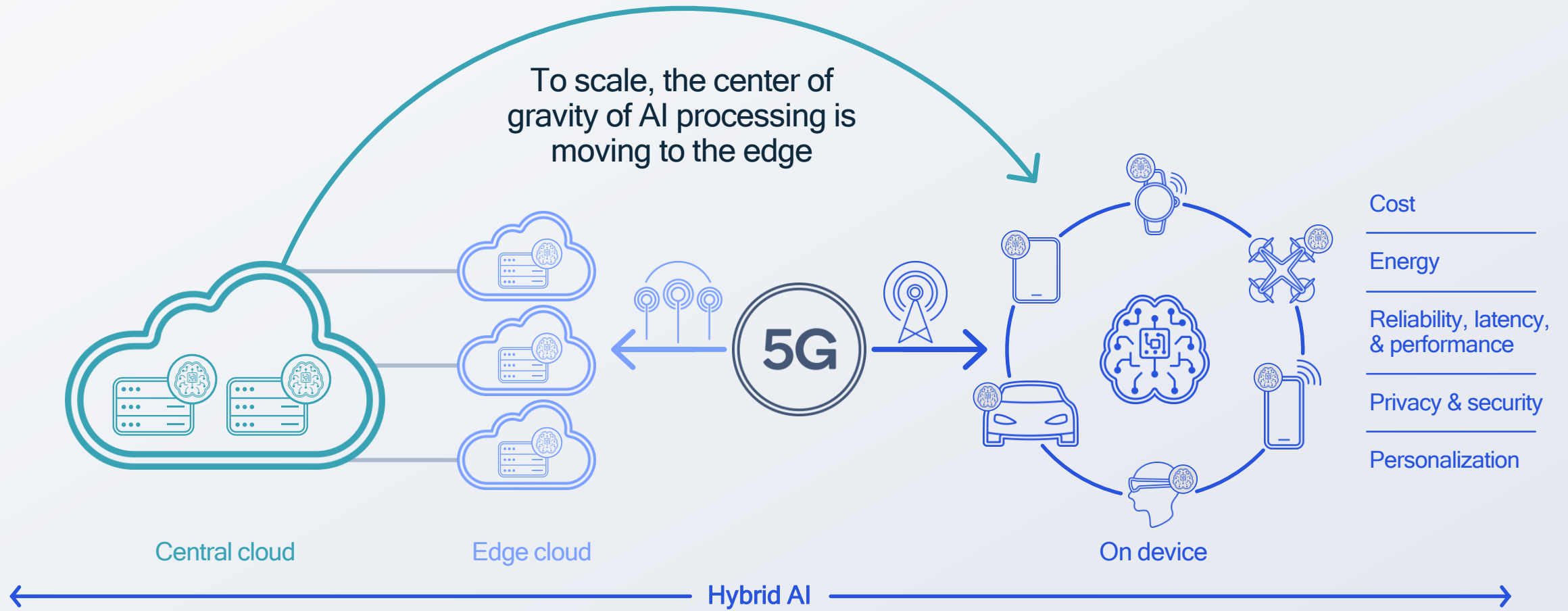# ML Inference at the Edge

## Felix Baum

Senior Director, Product Management
Qualcomm Technologies, Inc.

@qualcomm

To scale, the center of gravity of AI processing is moving to the edge

5G

Cost

Energy

Reliability, latency, & performance

Privacy & security

Personalization

Central cloud

Edge cloud

On device

Hybrid AI
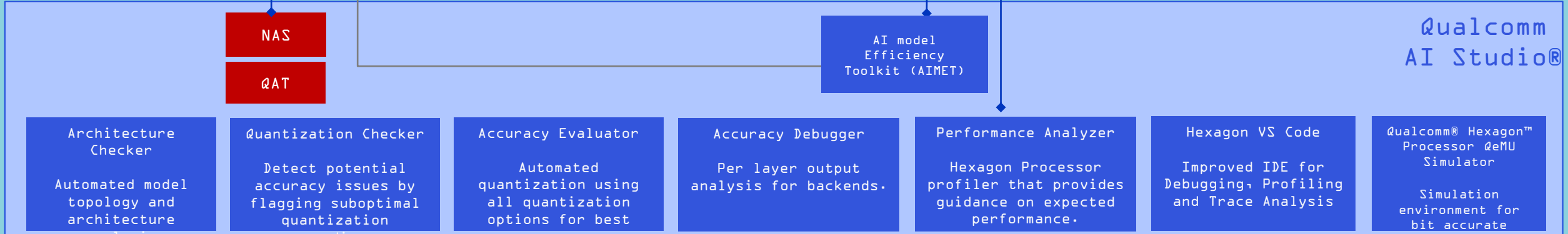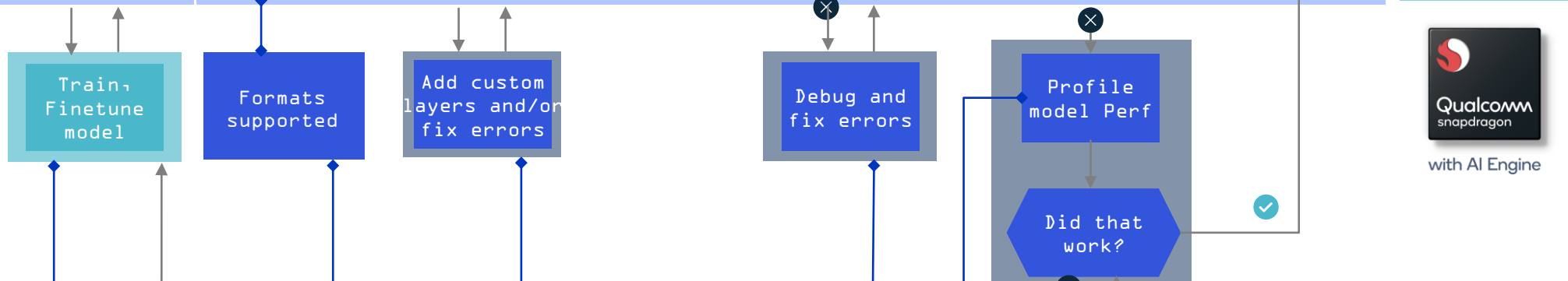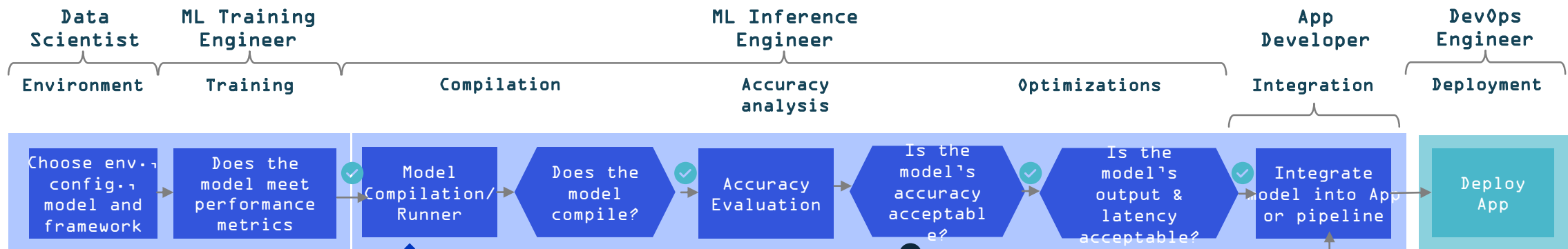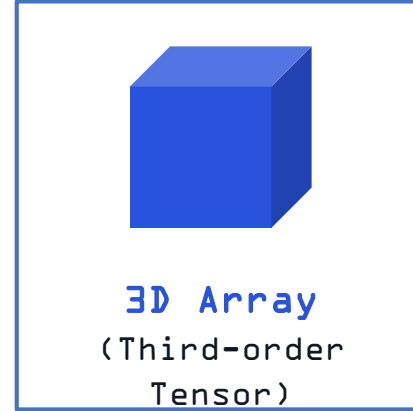
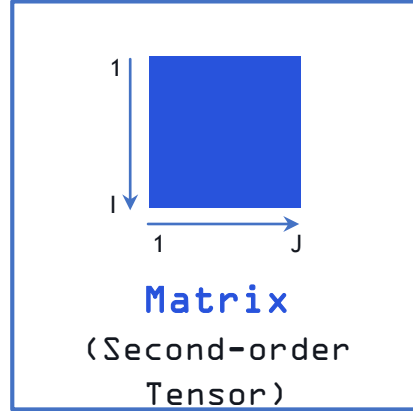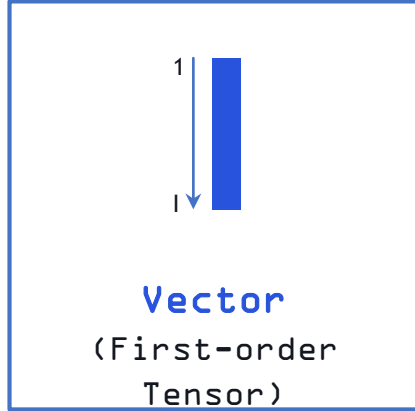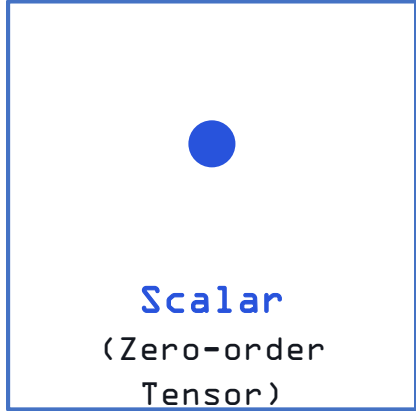We are leading the realization of the hybrid AI

Convergence of:

Wireless connectivity

Efficient computing

Distributed AI

Unlocking the data that will fuel our digital future and generative AI

2

**Roles (top):**
Data Scientist | ML Training Engineer | ML Inference Engineer | App Developer | DevOps Engineer

**Phases:**
Environment | Training | Compilation | Accuracy analysis | Optimizations | Integration | Deployment

**Flow boxes:**
- Choose env., config., model and framework
- Does the model meet performance metrics
- Model Compilation/Runner
- Does the model compile?
- Accuracy Evaluation
- Is the model's accuracy acceptable?
- Is the model's output & latency acceptable?
- Integrate model into App or pipeline
- Deploy App

**Supporting boxes:**
- Train, Finetune model
- Formats supported
- Add custom layers and/or fix errors
- Debug and fix errors
- Profile model Perf
- Did that work?
- Use advanced optimization techniques

**Red boxes:**
- Qualcomm® Neural Processing SDK and Qualcomm® AI Engine direct — Converter and Quantizer
- Custom Ops LLVM (C/C++) TVM (Python)
- NAS
- QAT
- AI model Efficiency Toolkit (AIMET)

**Qualcomm AI Stack / Qualcomm AI Studio®**

- Architecture Checker — Automated model topology and architecture analysis
- Quantization Checker — Detect potential accuracy issues by flagging suboptimal quantization encodings
- Accuracy Evaluator — Automated quantization using all quantization options for best accuracy
- Accuracy Debugger — Per layer output analysis for backends.
- Performance Analyzer — Hexagon Processor profiler that provides guidance on expected performance.
- Hexagon VS Code — Improved IDE for Debugging, Profiling and Trace Analysis
- Qualcomm® Hexagon™ Processor QeMU Simulator — Simulation environment for bit accurate validation

**Qualcomm snapdragon** with AI Engine

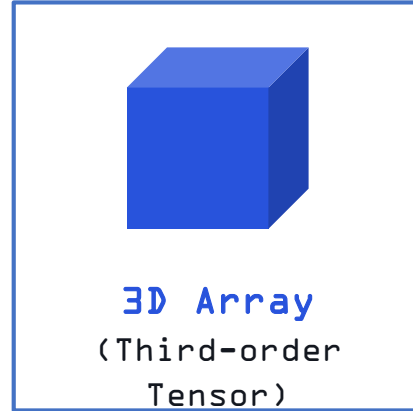AIMET is a product of Qualcomm Innovation Center, Inc

# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra

**Scalar**
(Zero-order
Tensor)

**Vector**
(First-order
Tensor)

**Matrix**
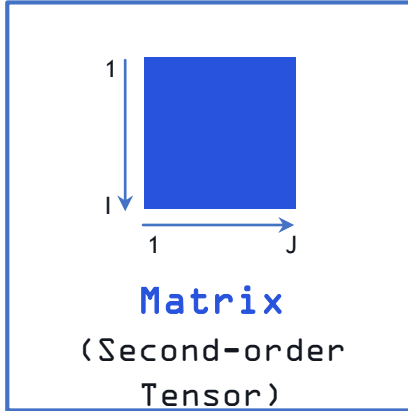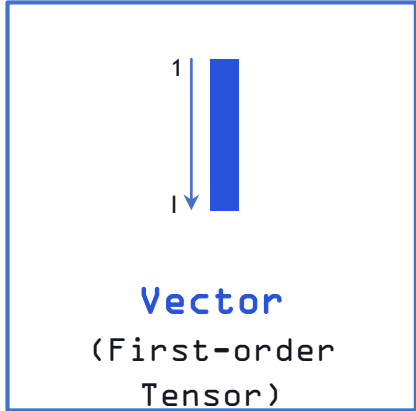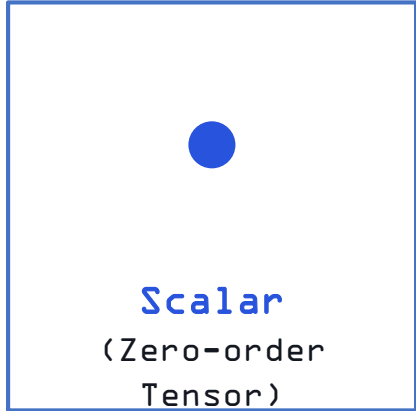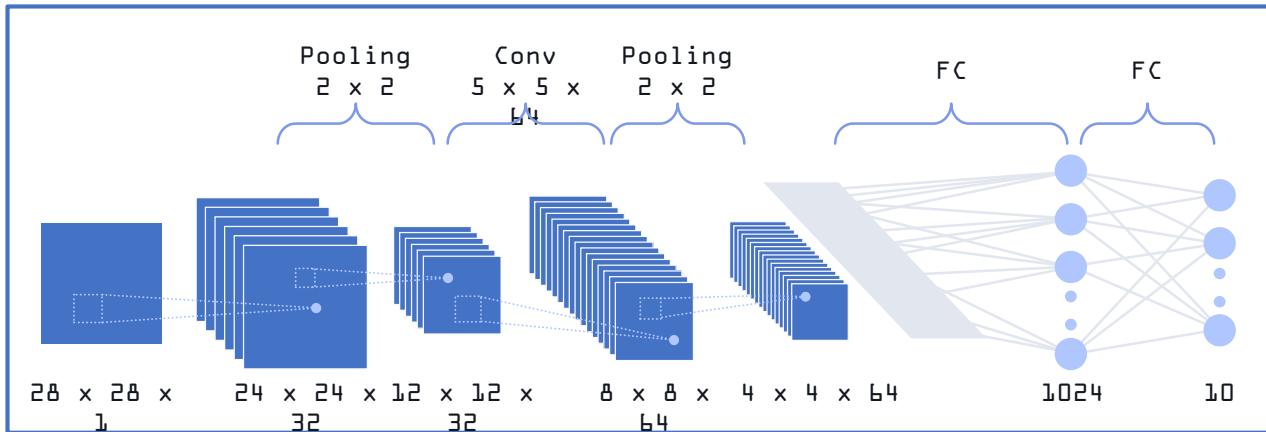(Second-order
Tensor)

**3D Array**
(Third-order
Tensor)

# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra



**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

**Matrix**
(Second-order Tensor)

**3D Array**
(Third-order Tensor)

Setup



Pooling 2 x 2   Conv 5 x 5 x 64   Pooling 2 x 2   FC   FC

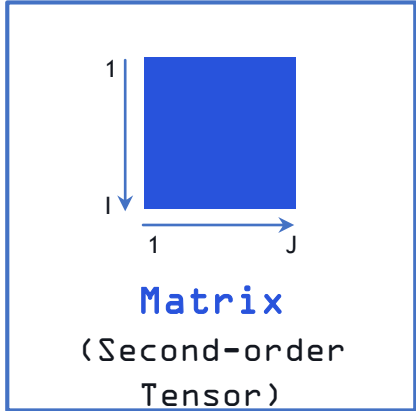28 x 28 x 1   24 x 24 x 32   12 x 12 x 32   8 x 8 x 64   4 x 4 x 64   1024   10

# Optimizing Hardware for AI
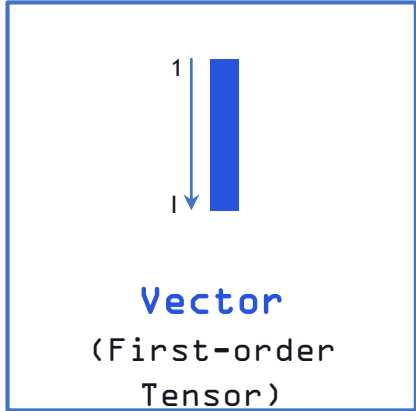
Neural Networks: A mundane pile of linear algebra

**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

1

I

**Matrix**
(Second-order Tensor)

1

I

1          J

**3D Array**
(Third-order Tensor)

Hexagon Processor

Setup

Pooling
2 x 2

Conv
5 x 5 x
64

Pooling
2 x 2

FC

FC

28 x 28 x
1

24 x 24 x 12 x 12 x
32         32

8 x 8 x    4 x 4 x 64
64

1024

10

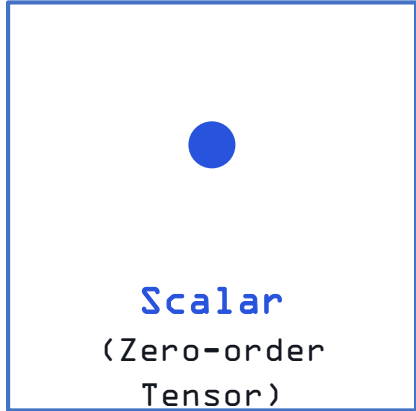# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra



**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

**Matrix**
(Second-order Tensor)

**3D Array**
(Third-order Tensor)

Setup

Pooling 2 x 2    Conv 5 x 5 x 64    Pooling 2 x 2    FC    FC

28 x 28 x 1    24 x 24 x 32    12 x 12 x 32    8 x 8 x 64    4 x 4 x 64    1024    10

Scalar Processor

**Hexagon Processor**

Scalar

Scalar Threads

# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra



**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

**Matrix**
(Second-order Tensor)

**3D Array**
(Third-order Tensor)

Hexagon Processor

Scalar

Scalar Threads

Vector

HVX

Setup

Pooling 2 x 2     Conv 5 x 5 x 64     Pooling 2 x 2     FC     FC

28 x 28 x 1     24 x 24 x 32     12 x 12 x 32     8 x 8 x 64     4 x 4 x 64     1024     10

Scalar Processor     Vector Processor     Vector Processor
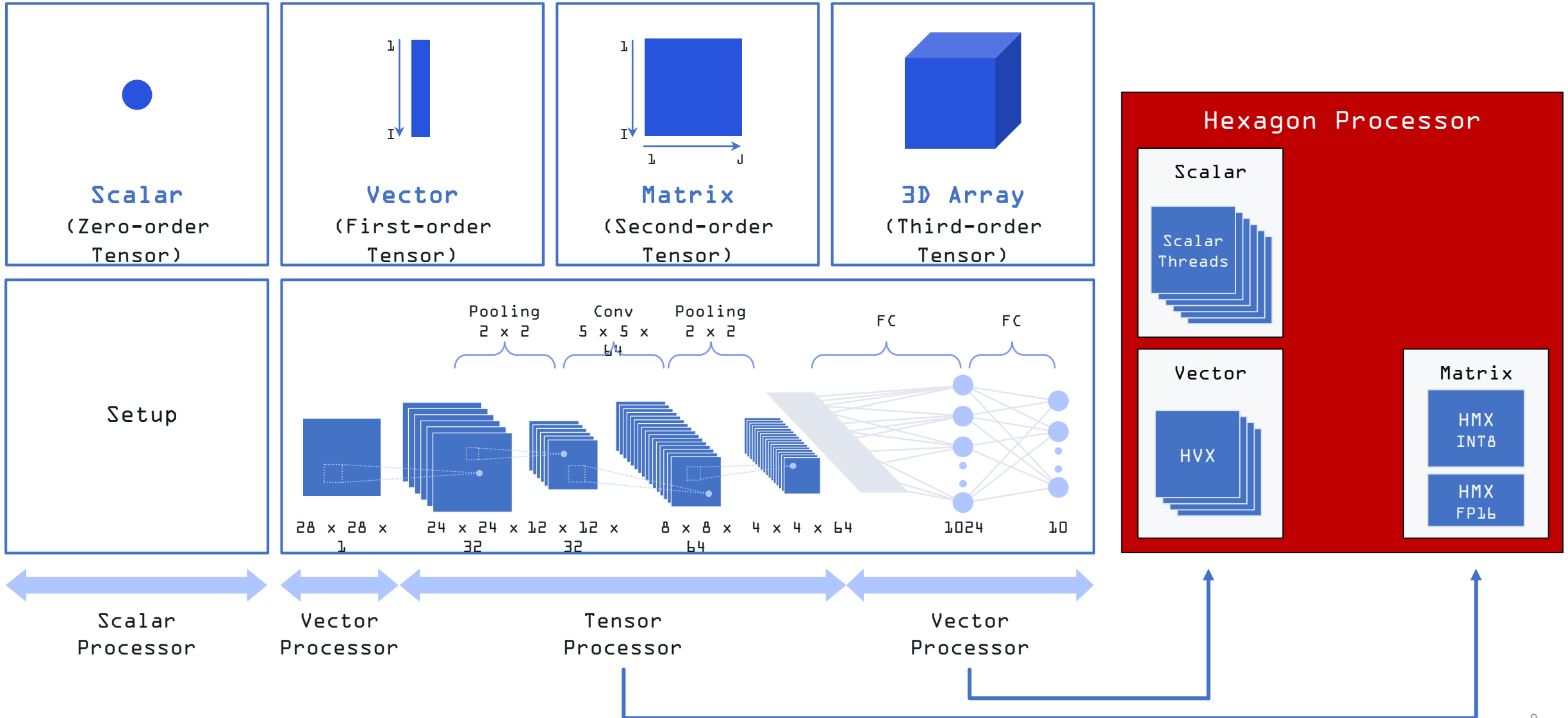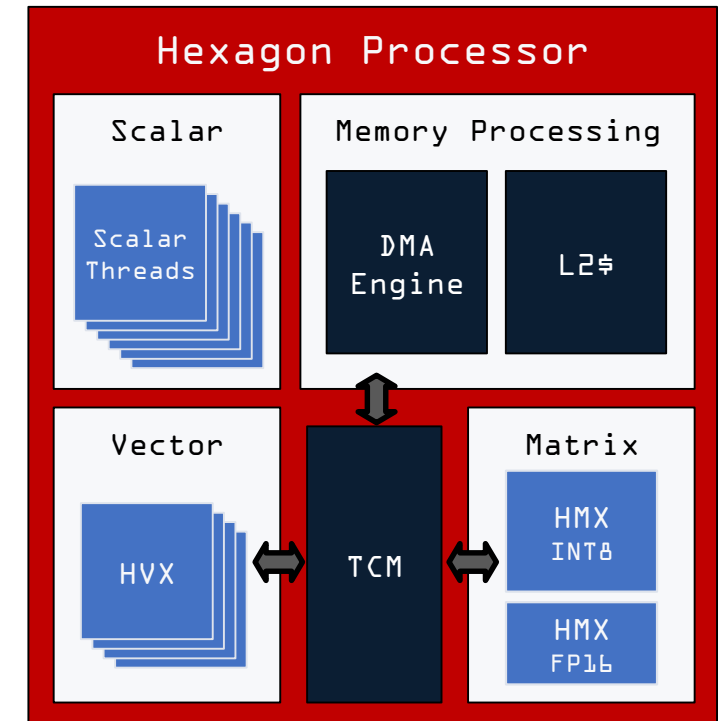
# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra

# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra

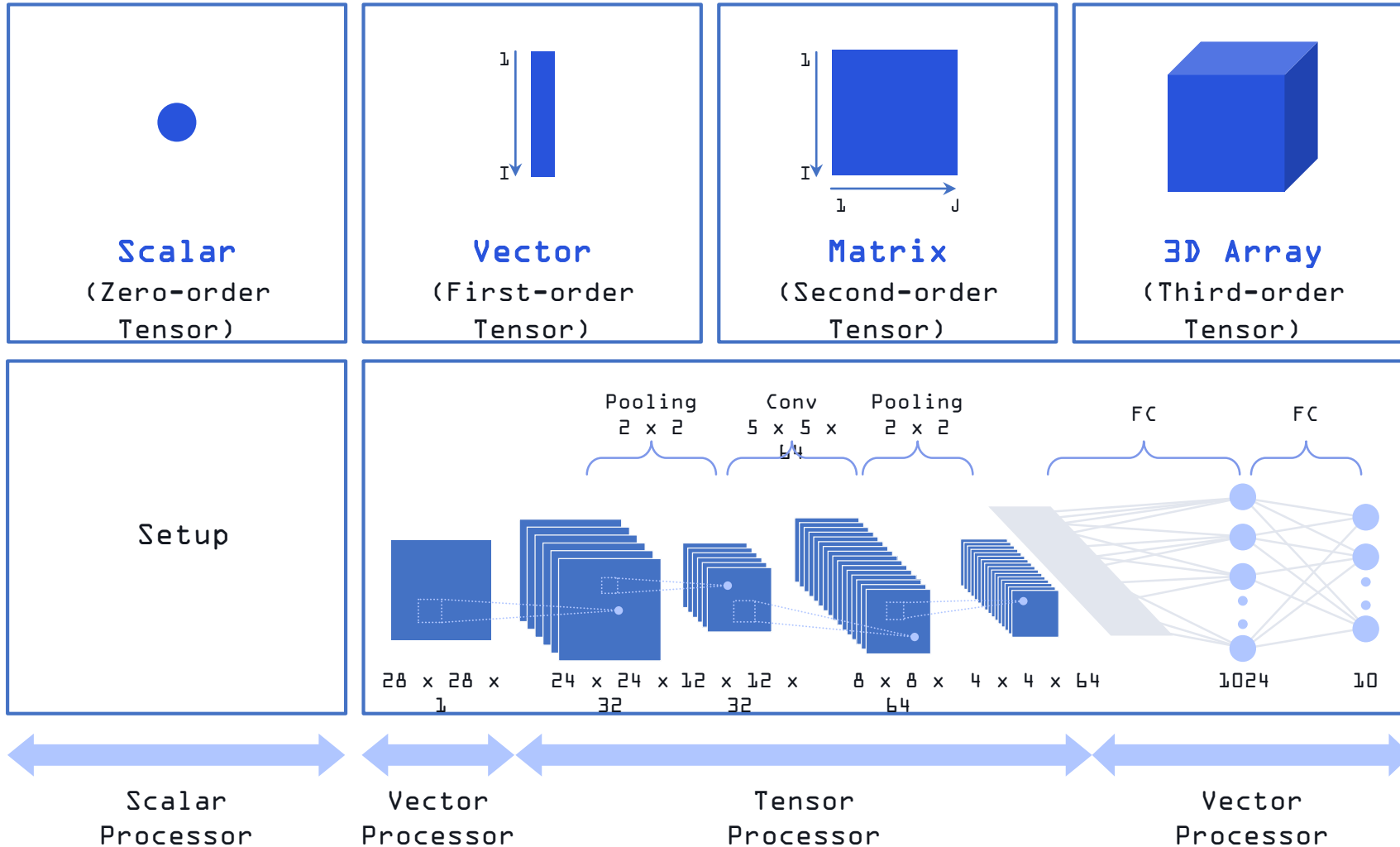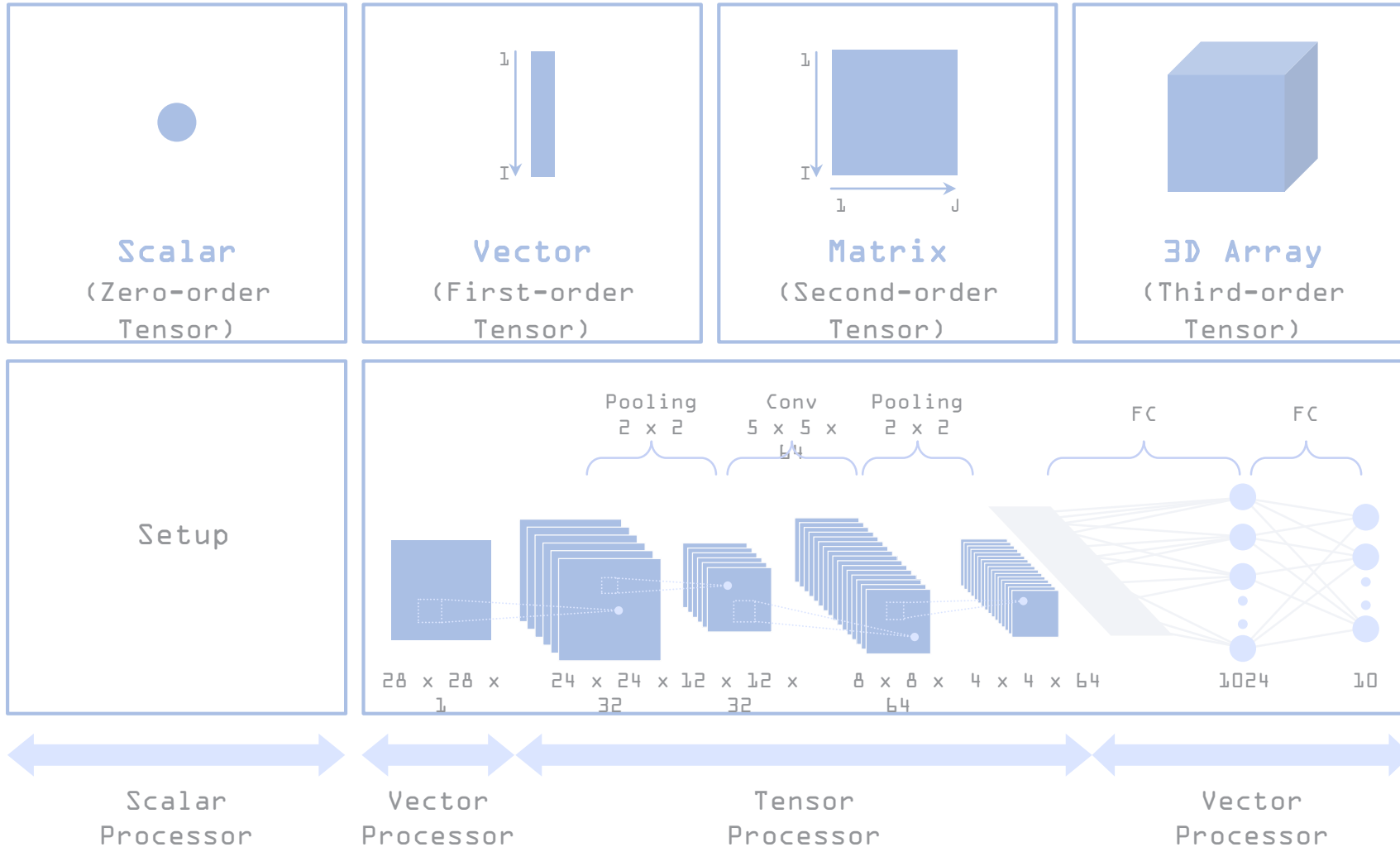# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra

# Optimizing Hardware for AI

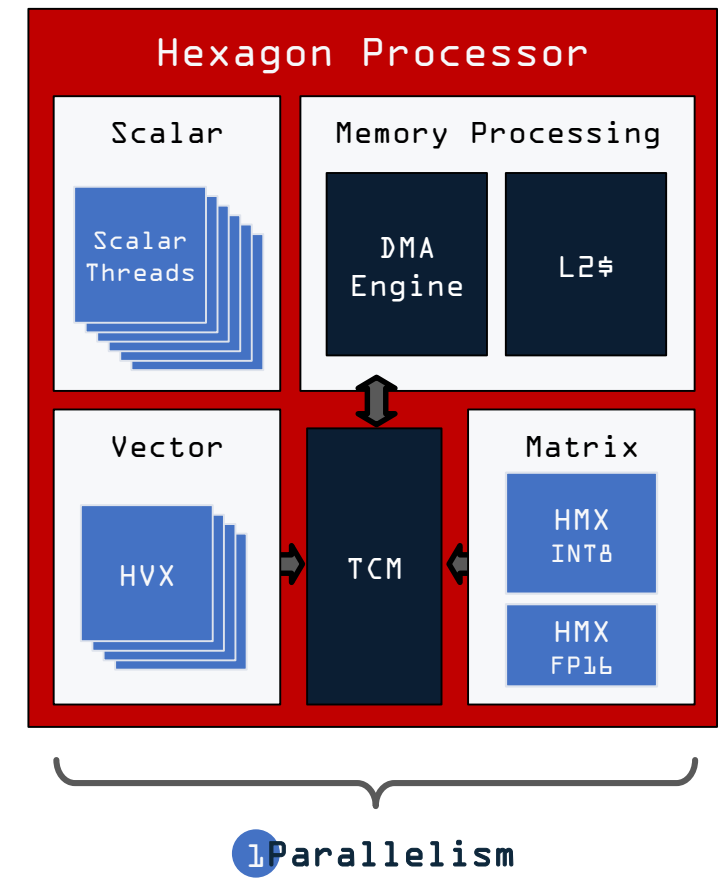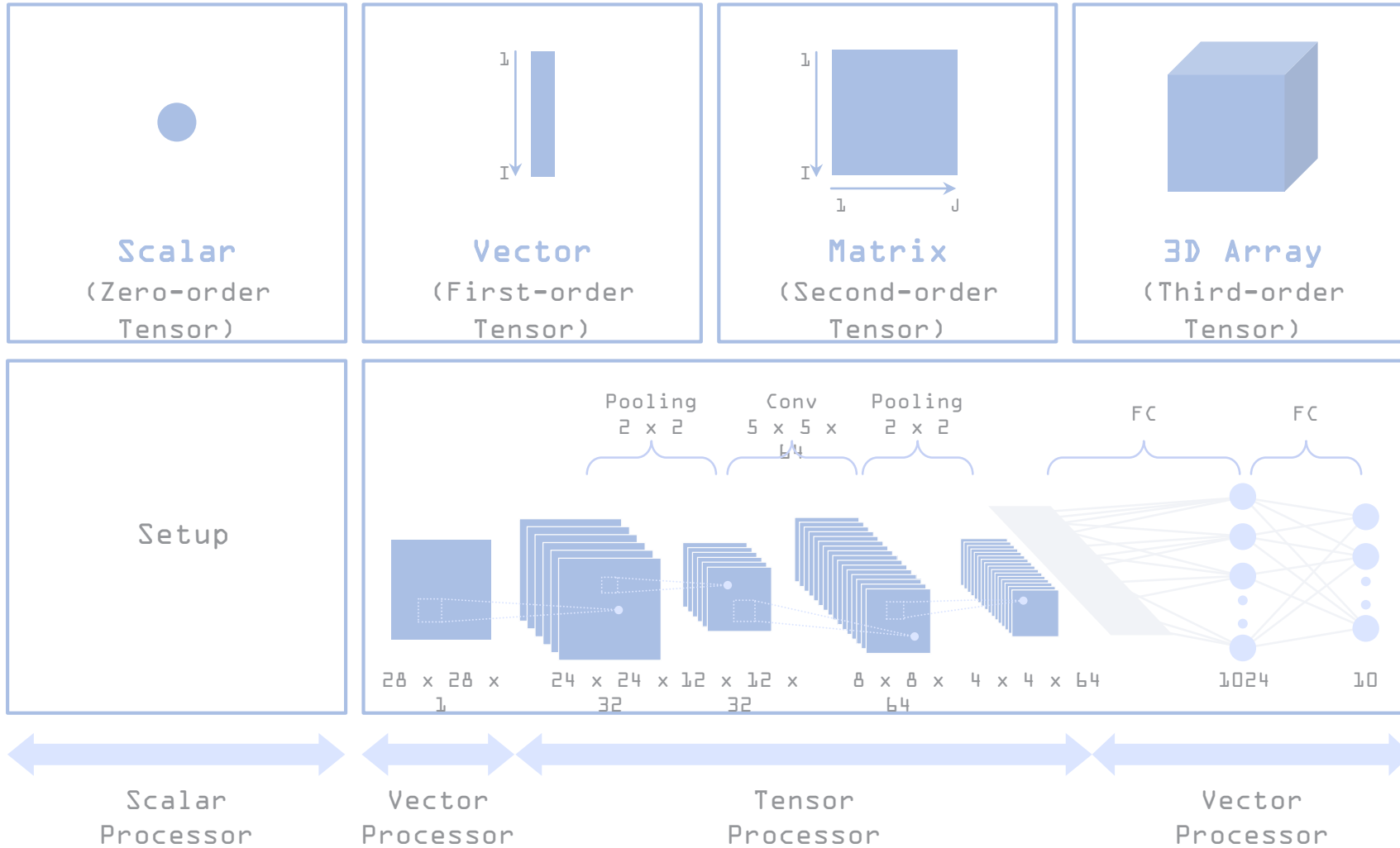Neural Networks: A mundane pile of linear algebra



**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

**Matrix**
(Second-order Tensor)

**3D Array**
(Third-order Tensor)

Setup

Pooling 2 x 2

Conv 5 x 5 x 64

Pooling 2 x 2

FC

FC

28 x 28 x 1

24 x 24 x 32

12 x 12 x 32

8 x 8 x 64

4 x 4 x 64

1024

10

Scalar Processor

Vector Processor

Tensor Processor

Vector Processor

**Optimization Goals**

1. Maximize parallelism
2. Minimize data movement

**Hexagon Processor**

Scalar

Scalar Threads

Memory Processing

DMA Engine

L2$

Vector

HVX

TCM

Matrix

HMX INT8

HMX FP16

DRAM

2 **Memory Access**

12

# Optimizing Hardware for AI

Neural Networks: A mundane pile of linear algebra



Scalar
(Zero-order Tensor)

Vector
(First-order Tensor)

Matrix
(Second-order Tensor)

3D Array
(Third-order Tensor)

Setup

Pooling 2 x 2    Conv 5 x 5 x 64    Pooling 2 x 2    FC    FC

28 x 28 x 1    24 x 24 x 32    12 x 12 x 32    8 x 8 x 64    4 x 4 x 64    1024    10

Scalar Processor    Vector Processor    Tensor Processor    Vector Processor

Optimization Goals

1. Maximize parallelism
2. Minimize data movement

Hexagon Processor

Scalar
Scalar Threads

Memory Processing
DMA Engine    L2$

Vector
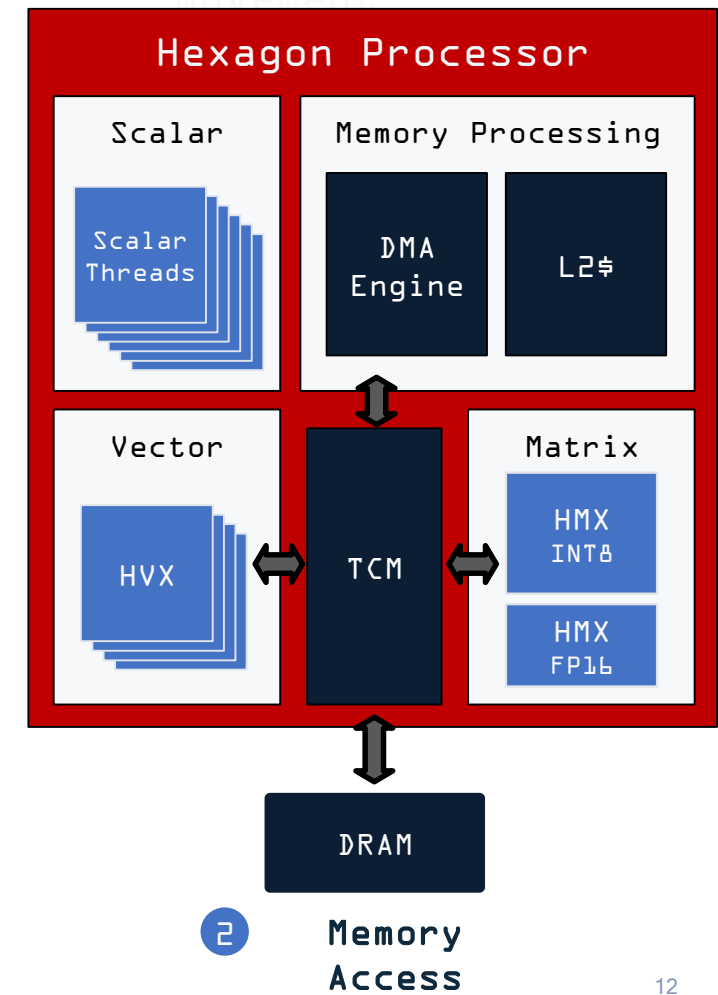HVX    TCM

Matrix
HMX INT8
HMX FP16

13

# Optimizing Hardware for AI: Transformers

Neural Networks: A mundane pile of linear algebra

## Optimization Goals
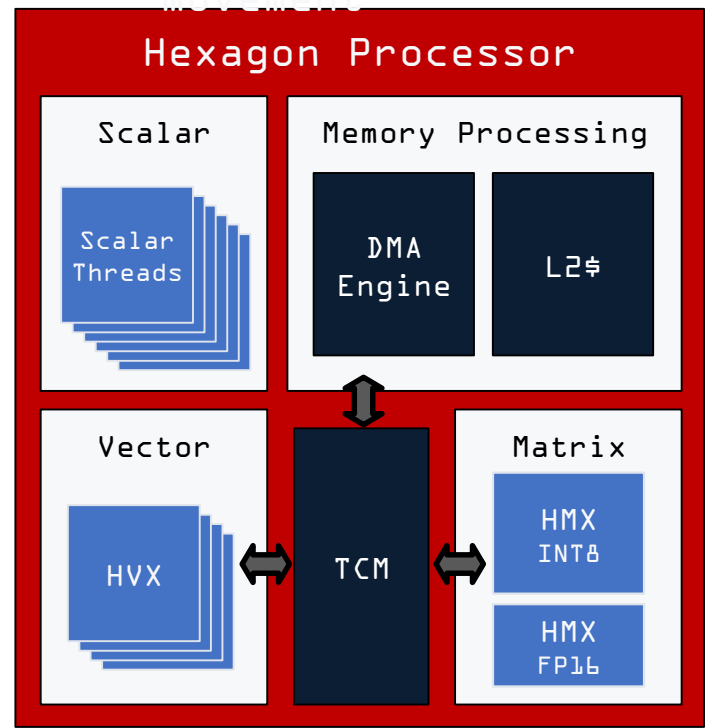
1. Maximize parallelism
2. Minimize data movement



Scalar
(Zero-order Tensor)

Vector
(First-order Tensor)

Matrix
(Second-order Tensor)

3D Array
(Third-order Tensor)

Setup

Scalar Processor

Transformer Architecture

HVX (~70%)

HMX (~30%)

Hexagon Processor

Scalar — Scalar Threads

Memory Processing — DMA Engine, L2$

Vector — HVX

TCM

Matrix — HMX INT8, HMX FP16

# Optimizing Hardware for AI: Super Resolution

Neural Networks: A mundane pile of linear algebra



**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

**Matrix**
(Second-order Tensor)

**3D Array**
(Third-order Tensor)

Setup

Scalar Processor

ABPN, SESR, XLSR, Q-SRNet

$f$ — number of channels

$m$ — number of intermediate conv layers

------► collapsible residual connection

--[p]--► "partial" collapsible residual connection

--[ri]--► "repeat-interleaving" collapsible residual connection

Super Resolution Architecture



## Optimization Goals

1. Maximize parallelism
2. Minimize data movement

HVX (~30%)

HMX (~90%)

### Hexagon Processor

**Scalar**

Scalar Threads

**Memory Processing**

DMA Engine

L2$

**Vector**

HVX

TCM

**Matrix**

HMX INT8

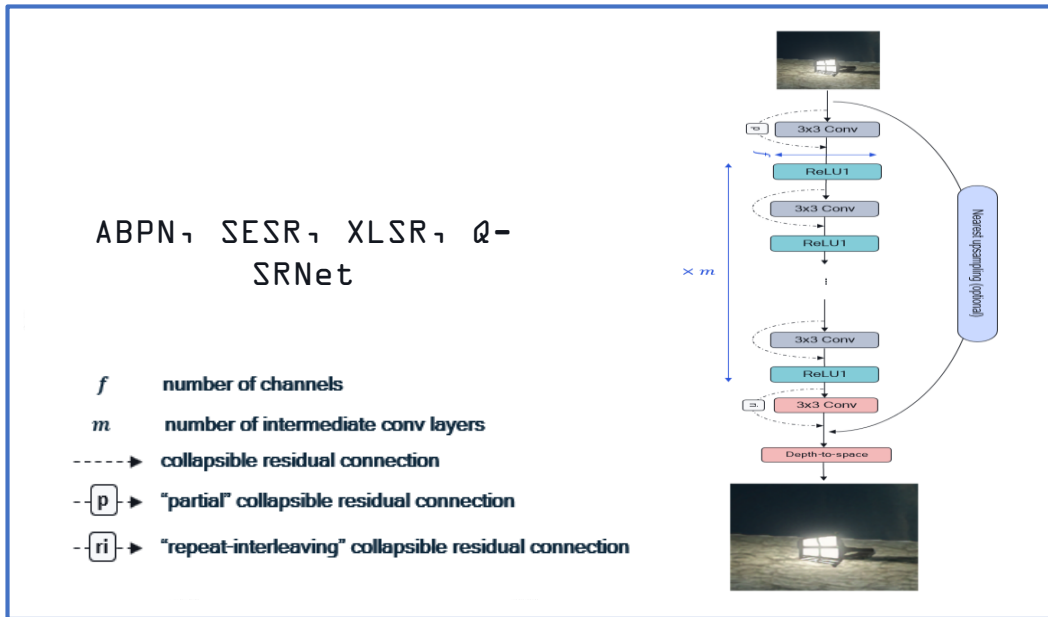HMX FP16

15

# Hexagon Processor: Execution of ML use cases

**Optimization Goals**

1. Maximize parallelism
2. Minimize data movement

**Scalar**
(Zero-order Tensor)

**Vector**
(First-order Tensor)

**Matrix**
(Second-order Tensor)

**3D Array**
(Third-order Tensor)

Setup

Pooling 2 x 2
Conv 5 x 5 x 64
Pooling 2 x 2
FC
FC

28 x 28 x 1
24 x 24 x 32
12 x 12 x 32
8 x 8 x 64
4 x 4 x 64
1024
10

Scalar Processor

Vector Processor

Tensor Processor

Vector Processor

## Hexagon Processor

**Scalar**

Scalar Threads

**Memory Processing**

DMA Engine

L2$

**Vector**

HVX

TCM

**Matrix**

HMX INT8

HMX FP16

# Hexagon Processor: Execution of end-to-end use cases

# Hexagon Processor: Concurrency Model

# Hexagon Processor: Concurrency Model

# AI Model Compilation: Steps

Conversion / Quantization



```
┌─────────────────────────────────┐
│  Framework optimizations        │
│            ↓                    │
│         Converter               │
│            ↓                    │
│         Quantizer               │
│            ↓                    │
│     Common   Graph              │
│      optimizations              │
└─────────────────────────────────┘
```

# AI Model Compilation: Steps

**Conversion / Quantization**

| |
|---|
| Framework optimizations |
| ↓ |
| Converter |
| ↓ |
| Quantizer |
| ↓ |
| Common Graph optimizations |

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

# AI Model Compilation: Steps

**Conversion / Quantization**

| Framework optimizations |
|:---:|
| ↓ |
| Converter |
| ↓ |
| Quantizer |
| ↓ |
| Common Graph optimizations |

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

# AI Model Compilation: Steps

**Conversion / Quantization**

| Framework optimizations |
|---|
| ↓ |
| Converter |
| ↓ |
| Quantizer |
| ↓ |
| Common Graph optimizations |

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

If required, graph can be quantized according to various config. parameters

# AI Model Compilation: Steps

**Conversion / Quantization**

| Framework optimizations | → | Framework level (Pytorch, TF, etc) optimizations, op folding, etc. |
|---|---|---|
| Converter | → | Framework graph is translated to the IR Graph |
| Quantizer | → | If required, graph can be quantized according to various config. parameters |
| Common Graph optimizations | → | Framework agnostic graph optimizations are applied such as batchnorm folding |

# AI Model Compilation: Steps

**Conversion / Quantization**
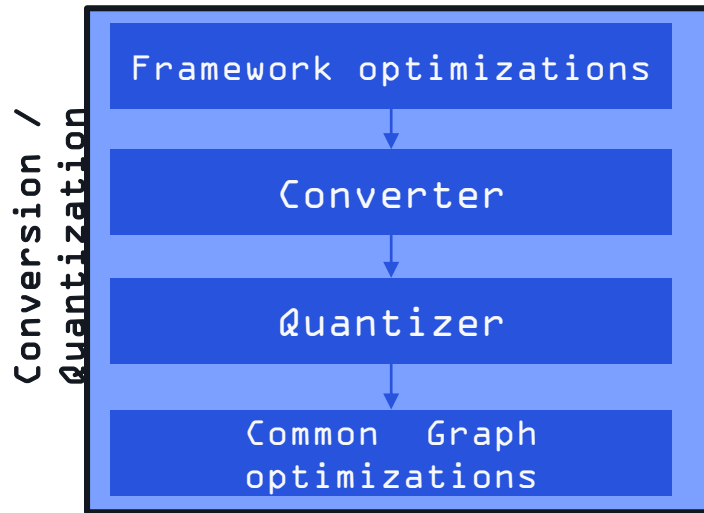
- Framework optimizations
- Converter
- Quantizer
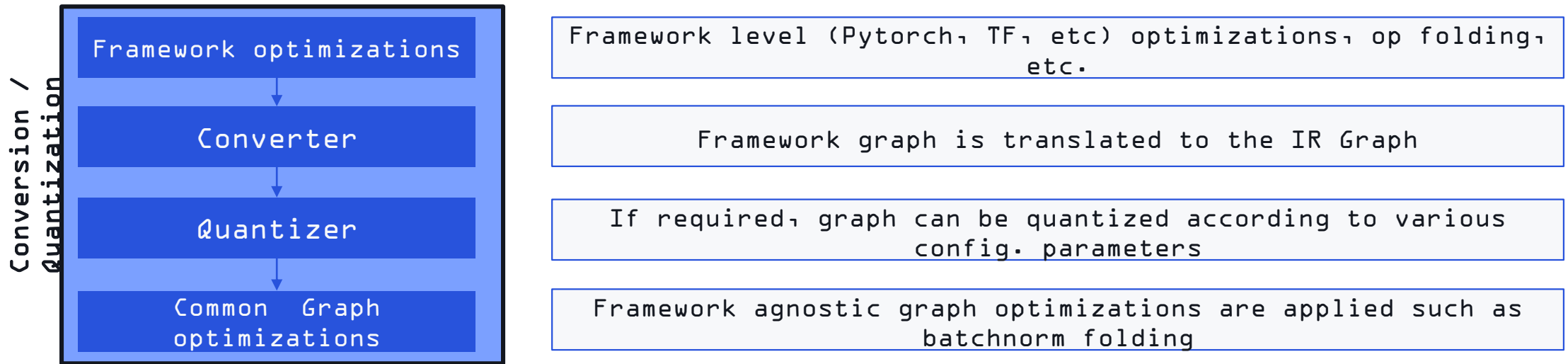- Common Graph optimizations
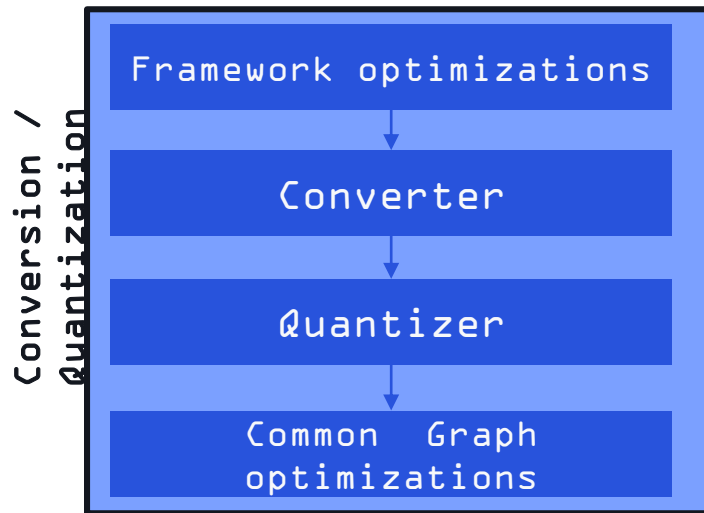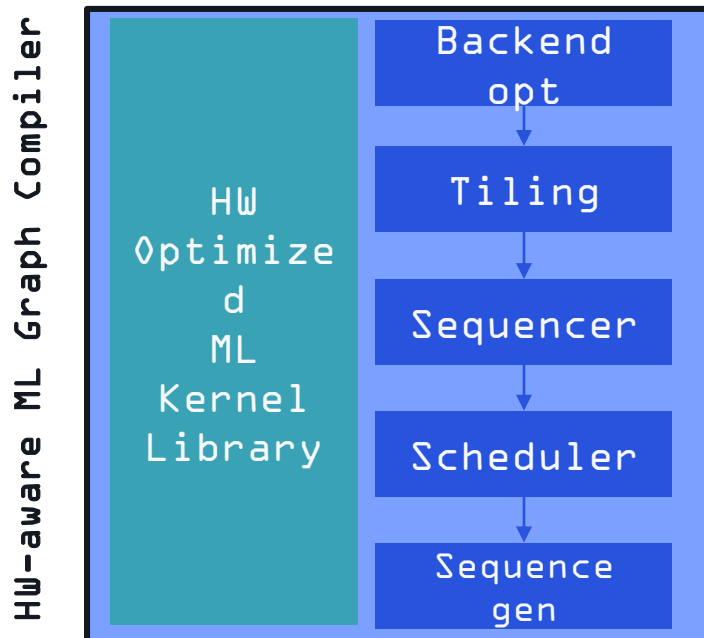
Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

If required, graph can be quantized according to various config. parameters

Framework agnostic graph optimizations are applied such as batchnorm folding

**HW-aware ML Graph Compiler**

- HW Optimized ML Kernel Library
- Backend opt
- Tiling
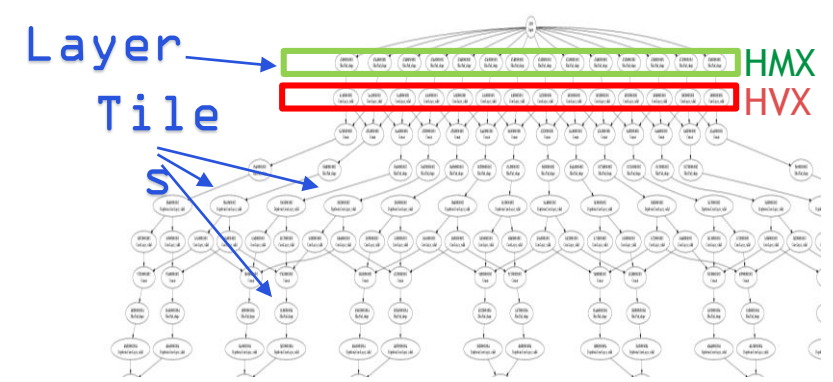- Sequencer
- Scheduler
- Sequence gen

Naive sequencers executed Nets "Layer-by-Layer", sequentially. Sometimes 1-3 layers can be aggregated (e.g., conv followed by RELU). "Layer by Layer" leaves performance and memory bandwidth on the table:
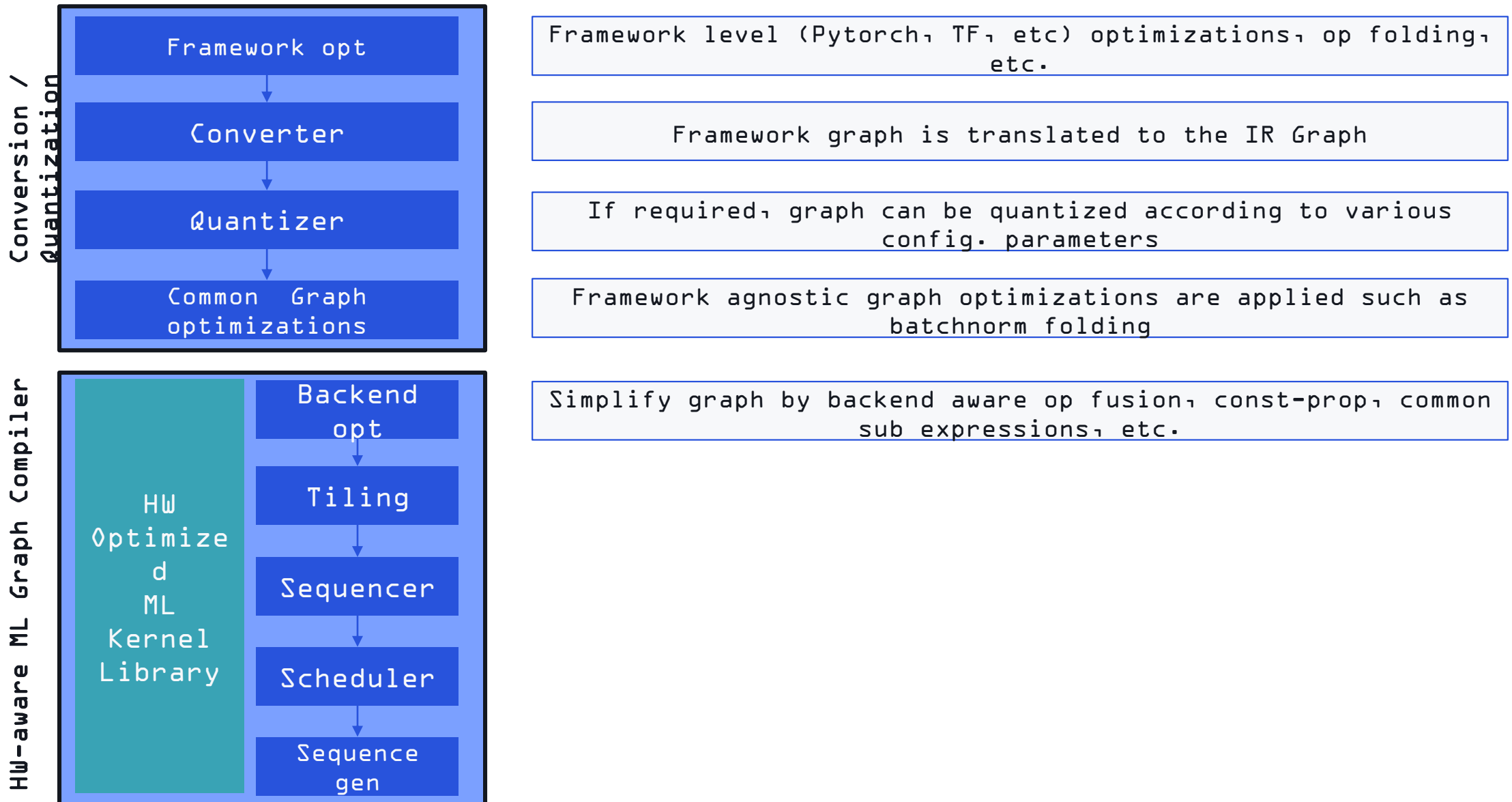
- If we exploit concurrencies and simultaneously operate on data from multiple layers, execution finishes faster
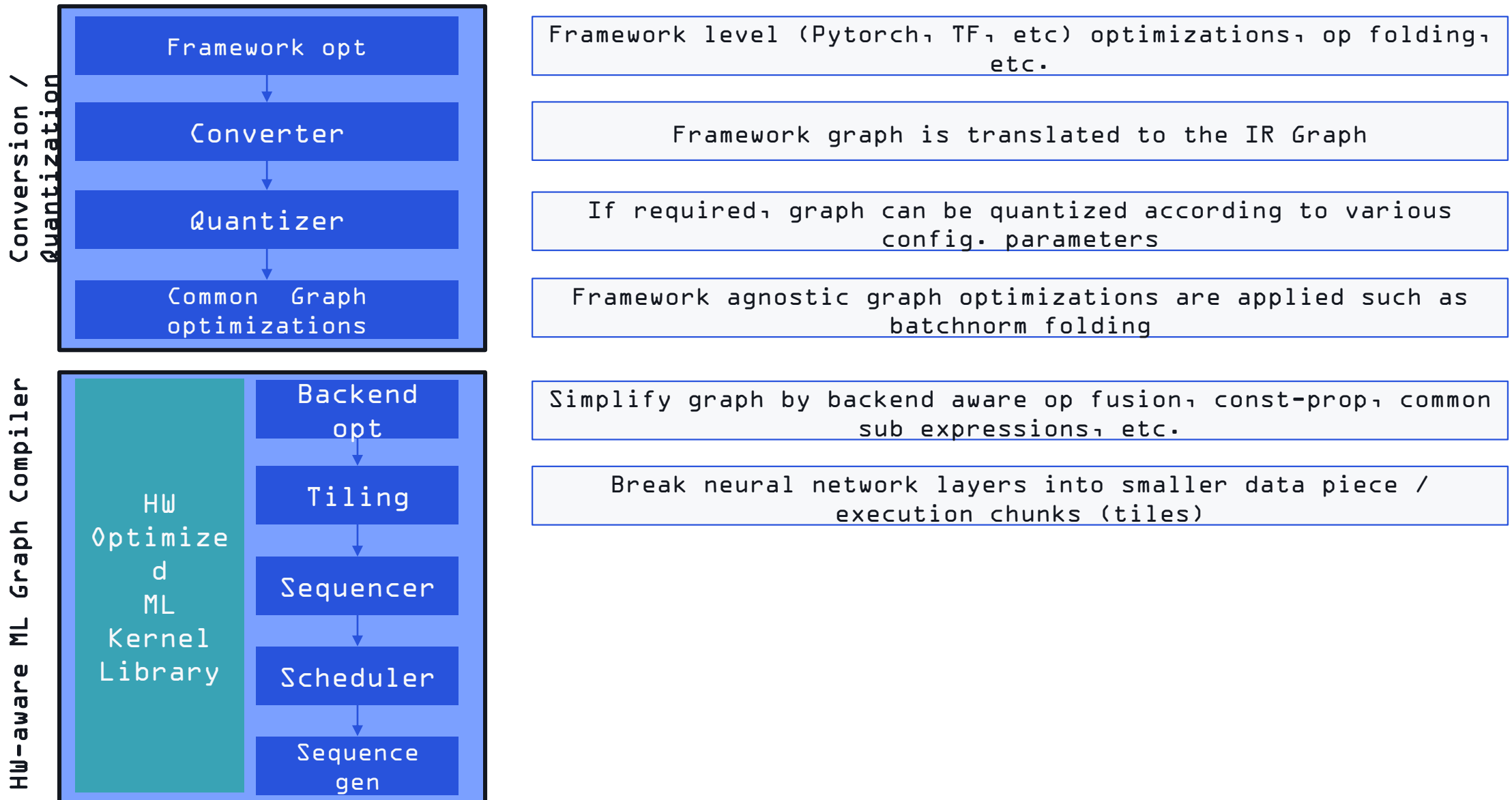- A layer's output, once consumed by next layer, is discardable. This saves DDR bandwidth, but TCM must be large enough, or data unit small enough, to store intermediate output

Layer

Tile

s

HMX

HVX

# AI Model Compilation: Steps

**Conversion / Quantization**

| Framework opt | Framework level (Pytorch, TF, etc) optimizations, op folding, etc. |
|---|---|
| Converter | Framework graph is translated to the IR Graph |
| Quantizer | If required, graph can be quantized according to various config. parameters |
| Common Graph optimizations | Framework agnostic graph optimizations are applied such as batchnorm folding |

**HW-aware ML Graph Compiler**

HW Optimized ML Kernel Library

| Backend opt | Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc. |
|---|---|
| Tiling | |
| Sequencer | |
| Scheduler | |
| Sequence gen | |

# AI Model Compilation: Steps

**Conversion / Quantization**

| | |
|---|---|
| Framework opt | Framework level (Pytorch, TF, etc) optimizations, op folding, etc. |
| Converter | Framework graph is translated to the IR Graph |
| Quantizer | If required, graph can be quantized according to various config. parameters |
| Common Graph optimizations | Framework agnostic graph optimizations are applied such as batchnorm folding |

**HW-aware ML Graph Compiler**

HW Optimized ML Kernel Library

| | |
|---|---|
| Backend opt | Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc. |
| Tiling | Break neural network layers into smaller data piece / execution chunks (tiles) |
| Sequencer | |
| Scheduler | |
| Sequence gen | |

# AI Model Compilation: Steps

**Conversion / Quantization**

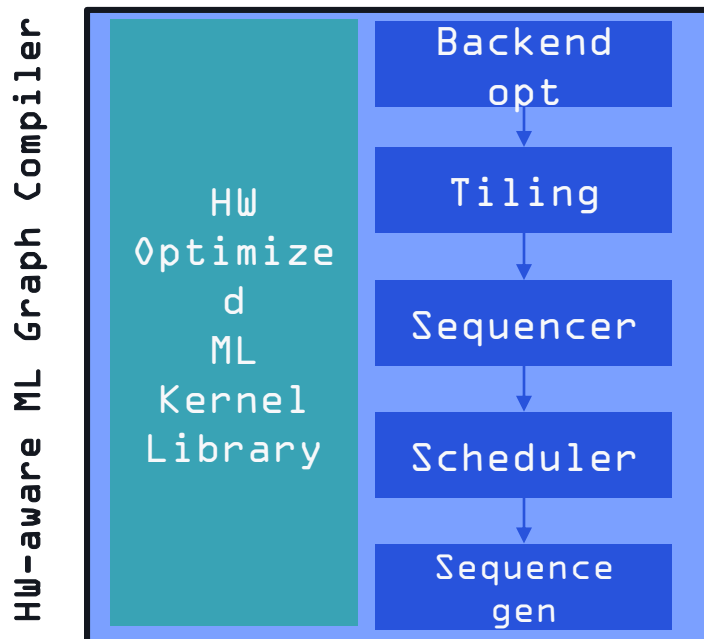| Framework opt |
|---|
| Converter |
| Quantizer |
| Common Graph optimizations |

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

If required, graph can be quantized according to various config. parameters

Framework agnostic graph optimizations are applied such as batchnorm folding

**HW-aware ML Graph Compiler**

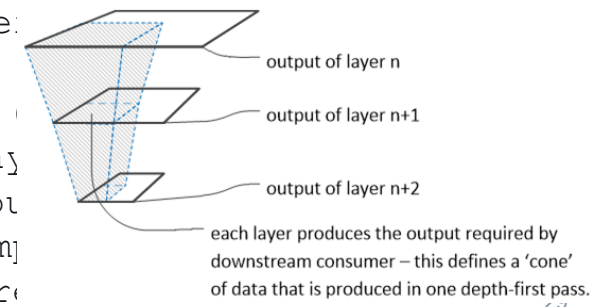| HW Optimized ML Kernel Library | Backend opt |
|---|---|
| | Tiling |
| | Sequencer |
| | Scheduler |
| | Sequence gen |

Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc.

Break neural network layers into smaller data piece / execution chunks (tiles)

**Scheduler** and **Sequencer** tools dictate order of tile execution to get best performance (completion time) & reduce DDR BW and power. To handle variety of network architectures, different types of sequencers are created. Each Sequencer is composed of cooperating algos & heuristics where some can be non-linear – small changes in networks give different results

# AI Model Compilation: Steps

**Conversion / Quantization**

- Framework opt
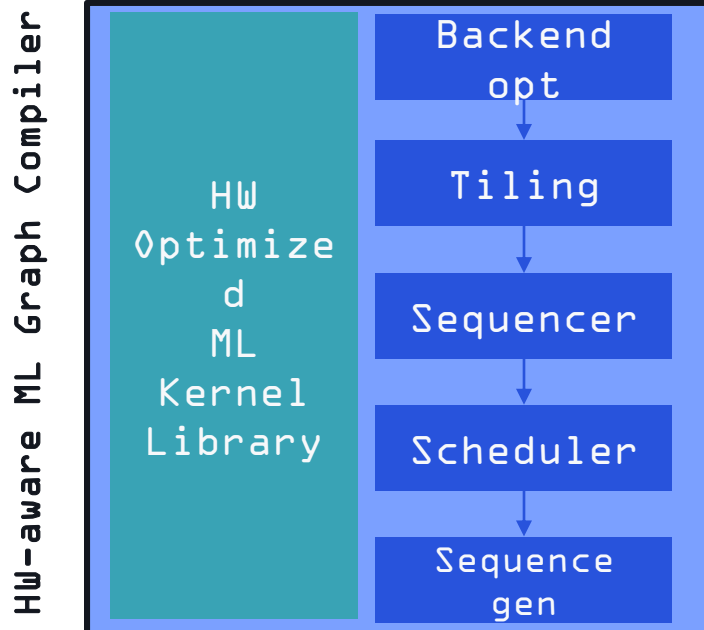- Converter
- Quantizer
- Common Graph optimizations

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

If required, graph can be quantized according to various config. parameters

Framework agnostic graph optimizations are applied such as batchnorm folding

**HW-aware ML Graph Compiler**

- HW Optimized ML Kernel Library
- Backend opt
- Tiling
- Sequencer
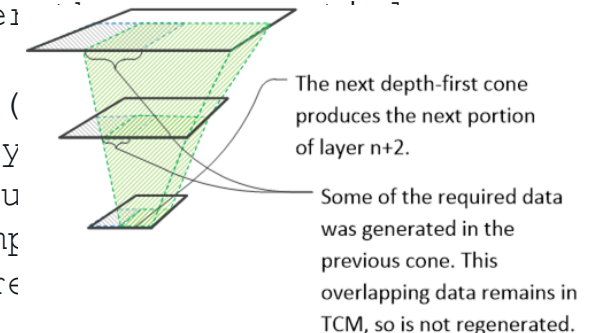- Scheduler
- Sequence gen

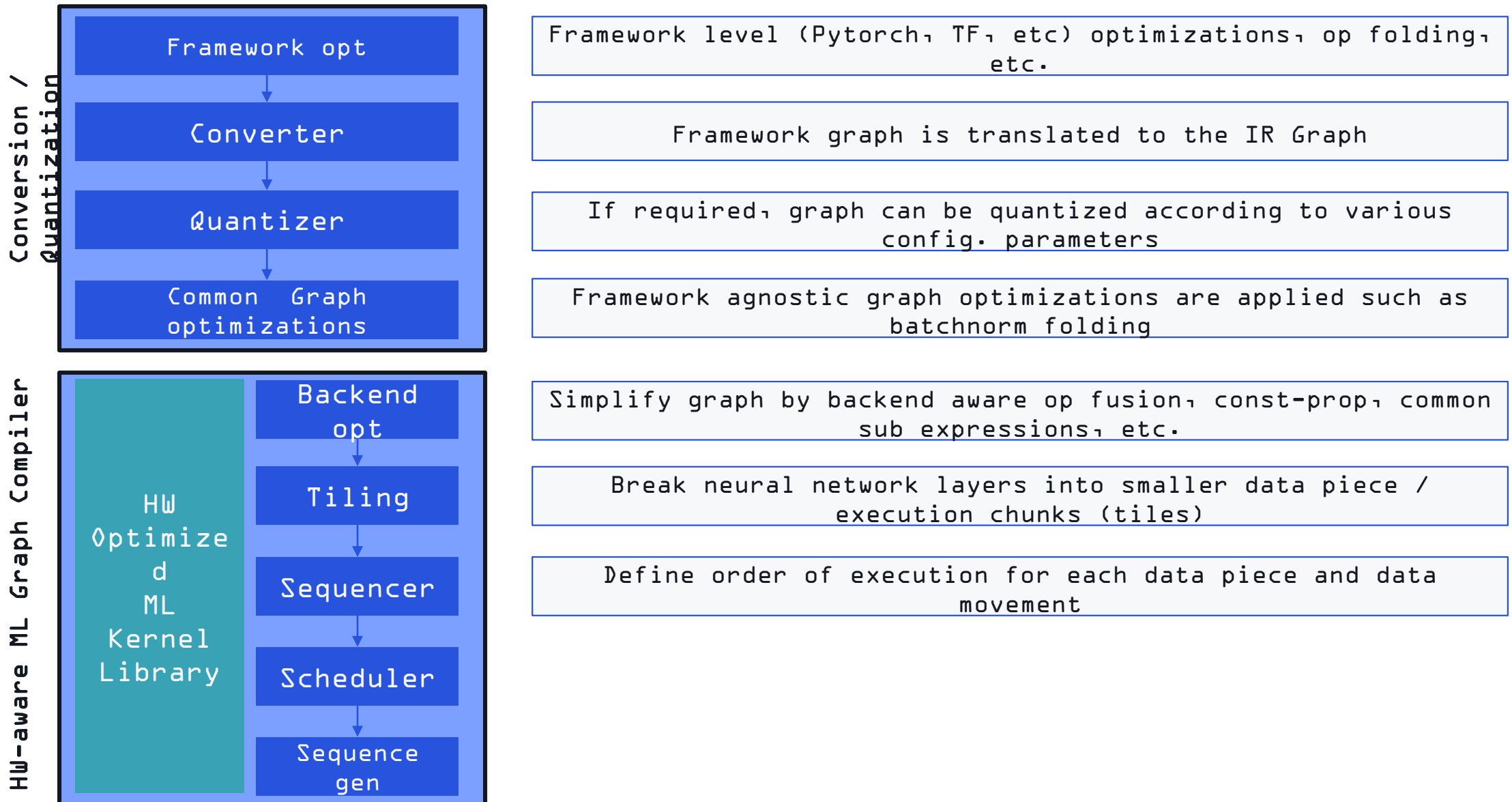Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc.

Break neural network layers into smaller data piece / execution chunks (tiles)

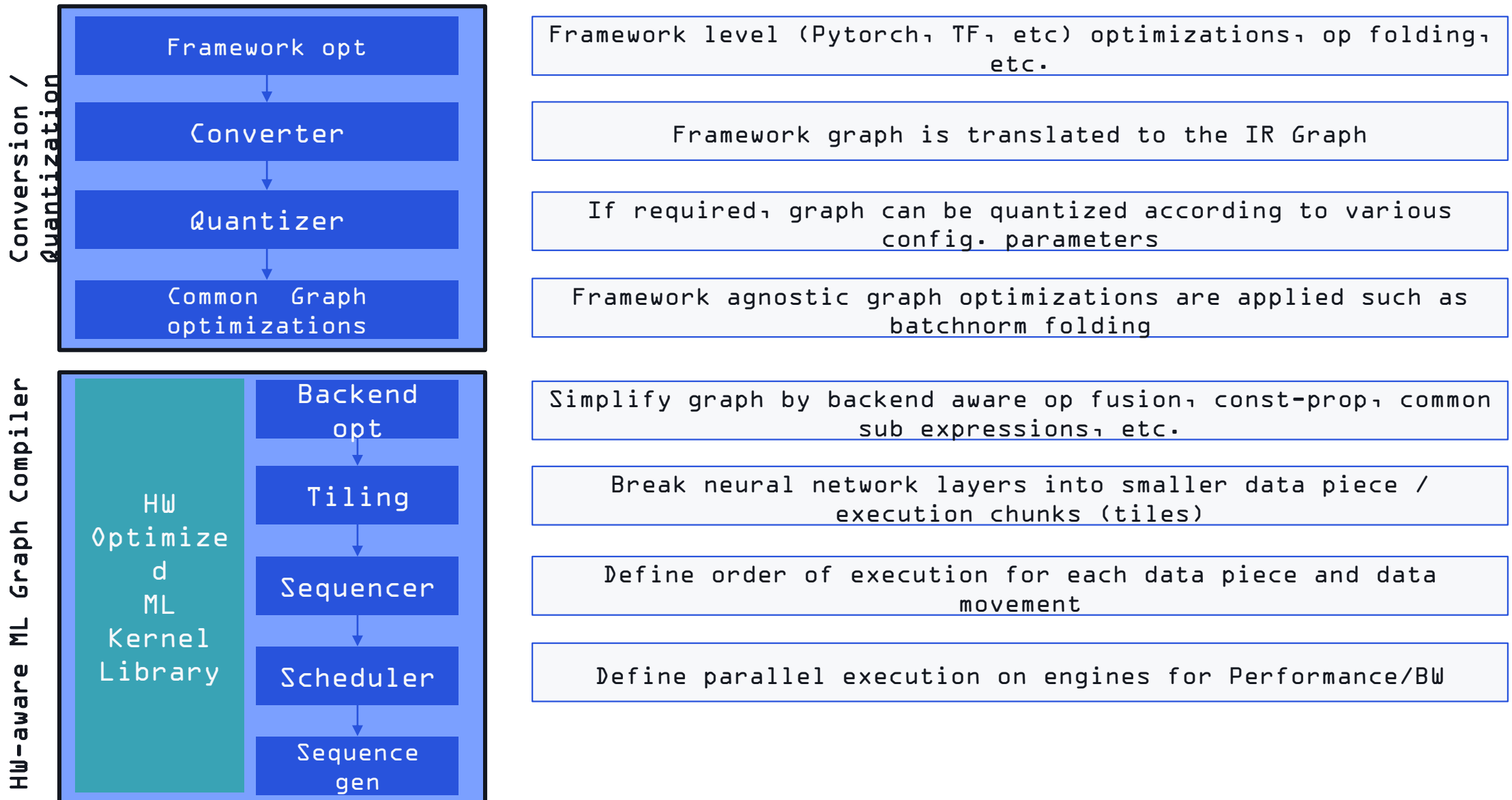To minimize DDR Bandwidth pressure, utilize locality between successive layers to reduce DDR BW. Conside~ layers:
- Output of layer n is the input of layer ~
- Output of layer (n+1) is the input of lay~
- Output of Layer (n+2) is divided into fou~

Each portion results in a separate into comp~
Intermediate results within a cone are stor~
do not consume DDR bandwidth.



output of layer n

output of layer n+1

output of layer n+2

each layer produces the output required by downstream consumer – this defines a 'cone' of data that is produced in one depth-first pass.

# AI Model Compilation: Steps

**Conversion / Quantization**

| Framework opt |
| Converter |
| Quantizer |
| Common Graph optimizations |

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

If required, graph can be quantized according to various config. parameters

Framework agnostic graph optimizations are applied such as batchnorm folding

**HW-aware ML Graph Compiler**

HW Optimized ML Kernel Library

| Backend opt |
| Tiling |
| Sequencer |
| Scheduler |
| Sequence gen |

Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc.

Break neural network layers into smaller data piece / execution chunks (tiles)

To minimize DDR Bandwidth pressure, utilize locality between successive layers to reduce DDR BW. Consider ~~the~~ ~~~~~~ ~~~~~~ layers:
- Output of layer n is the input of layer ~~(~~
- Output of layer (n+1) is the input of lay~~~~
- Output of Layer (n+2) is divided into fou~~~~
Each portion results in a separate into comp~~~~
Intermediate results within a cone are store~~~~
do not consume DDR bandwidth.



The next depth-first cone produces the next portion of layer n+2.

Some of the required data was generated in the previous cone. This overlapping data remains in TCM, so is not regenerated.

# AI Model Compilation: Steps

**Framework opt**

↓

**Converter**

↓

**Quantizer**

↓

**Common Graph optimizations**

Framework level (Pytorch, TF, etc) optimizations, op folding, etc.

Framework graph is translated to the IR Graph

If required, graph can be quantized according to various config. parameters

Framework agnostic graph optimizations are applied such as batchnorm folding

**HW-aware ML Graph Compiler**

**HW Optimized ML Kernel Library**

**Backend opt**

↓

**Tiling**

↓

**Sequencer**

↓

**Scheduler**

↓

**Sequence gen**

Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc.

Break neural network layers into smaller data piece / execution chunks (tiles)

Define order of execution for each data piece and data movement

# AI Model Compilation: Steps

**Conversion / Quantization**

| Framework opt | Framework level (Pytorch, TF, etc) optimizations, op folding, etc. |
| Converter | Framework graph is translated to the IR Graph |
| Quantizer | If required, graph can be quantized according to various config. parameters |
| Common Graph optimizations | Framework agnostic graph optimizations are applied such as batchnorm folding |

**HW-aware ML Graph Compiler**

HW Optimized ML Kernel Library

| Backend opt | Simplify graph by backend aware op fusion, const-prop, common sub expressions, etc. |
| Tiling | Break neural network layers into smaller data piece / execution chunks (tiles) |
| Sequencer | Define order of execution for each data piece and data movement |
| Scheduler | Define parallel execution on engines for Performance/BW |
| Sequence gen | |

# AI Model Compilation: Steps

**Conversion / Quantization**

| | |
|---|---|
| Framework opt | Framework level (Pytorch, TF, etc) optimizations, op folding, etc. |
| Converter | Framework graph is translated to the IR Graph |
| Quantizer | If required, graph can be quantized according to various config. parameters |
| Common Graph optimizations | Framework agnostic graph optimizations are applied such as batchnorm folding |

**HW-aware ML Graph Compiler**

HW Optimized ML Kernel Library

| | |
|---|---|
| Backend opt | Simplify graph by op fusion, const-prop, common sub expressions, etc. |
| Tiling | Break neural network into smaller data pieces |
| Sequencer | Define order of execution for each data piece and data movement |
| Scheduler | Define parallel execution on engines for Performance/BW |
| Sequence gen | Generate the optimized list of functions to run on hardware |

# AI Model Compilation: Sequencer determines optimal order

What order do I execute each operation?

All orders must follow a topological sort.



Very simple network for illustration with only 10 operations

Red lines show 3 potential valid topological sorts

1102 Valid topological sorts for this simple network of 10 operations!

Compiler algos trade-off DDR BW & Performance (latency) for each network.

# AI Model Compilation: Optimal Execution Order

Threads, Run Orders, Timelines



Run Order:

A B C

Foreground Process

Background Process 1

Background Process 2

0    1    2    3    4    5    6    7    8

Execution Cycles

# AI Model Compilation: Optimal Execution Order

Threads, Run Orders, Timelines



Run Order:

A B C

Foreground Process

Background Process 1

Background Process 2

A    B

C

Execution Cycles

0   1   2   3   4   5   6   7   8

# AI Model Compilation: Optimal Execution Order

Tiling

# AI Model Compilation: Optimal Execution Order

Tiling



Run Order:

A B1 B2 C1 C2 C3 C4

Foreground Process

A     B

Background Process 1

C

Background Process 2

Execution Cycles

0   1   2   3   4   5   6   7   8

# AI Model Compilation: Optimal Execution Order

Tiling



Run Order:

A B1 B2 C1 C2 C3 C4

Foreground Process: A, B1, B2

Background Process 1: C1, C2, C3, C4

Background Process 2

Execution Cycles

0  1  2  3  4  5  6  7  8

# AI Model Compilation: Optimal Execution Order

Scheduling

C1 and C2 only
depend on B1, so
they can be
reordered with B2.

Run Order:

A B1 C1 C2 B2 C3 C4

| | Foreground Process | Background Process 1 | Background Process 2 |
|---|---|---|---|

Execution Cycles

# AI Model Compilation: Optimal Execution Order

Optimal ordering

C1 and C2 only depend on B1, so they can be reordered with B2.

Run Order:

A B1 C1 C2 B2 C3 C4



| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Foreground Process** | A | B1 | B2 | | | | | |
| **Background Process 1** | | | C1 | C3 | | | | |
| **Background Process 2** | | | C2 | C4 | | | | |

0    1    2    3    4    5    6    7    8

Execution Cycles

# AI Model Performance: inf/sec

Snapdragon 8 Gen 2 VS Competitor A / Competitor B

## Super resolution (RDN)

| | |
|---|---|
| Snapdragon 8 Gen2 | ████████████████████████ |
| Competitor A | ██████████████ |
| Competitor B | ████████ |

## Face recognition (FaceNet)

| | |
|---|---|
| Snapdragon 8 Gen2 | ████████████████████████ |
| Competitor A | ████████████ |
| Competitor B | ██████████ |

## Bokeh (DeeplabV3+)

| | |
|---|---|
| Snapdragon 8 Gen2 | ████████████████████ |
| Competitor A | ███████ |
| Competitor B | ██ |

## Natural language processing (MobileBERT)

| | |
|---|---|
| Snapdragon 8 Gen2 | ████████████████████████ |
| Competitor A | ██ |
| Competitor B | ██████████████ |

# AI Model Performance: inf/sec per watt

**Snapdragon 8 Gen 2**

VS

Competitor A

Competitor B

## Super resolution (RDN)

Snapdragon 8 Gen2

Competitor A

Competitor B

## Face recognition (FaceNet)

Snapdragon 8 Gen2

Competitor A

Competitor B

## Bokeh (DeeplabV3+)

Snapdragon 8 Gen2

Competitor A

Competitor B

## Natural language processing (MobileBERT)

Snapdragon 8 Gen2

Competitor A

Competitor B

# Thank you