

μ IR – An intermediate representation for transforming and optimizing the microarchitecture of application accelerators

Amirali Sharifian¹, Reza Hojabr¹, Navid Rahimi¹,
Sihao Liu², Apala Guha¹
Tony Nowatzki² and Arrvindh Shriraman¹



<https://github.com/sfu-arch/uir>

Simon Fraser University¹, UCLA²

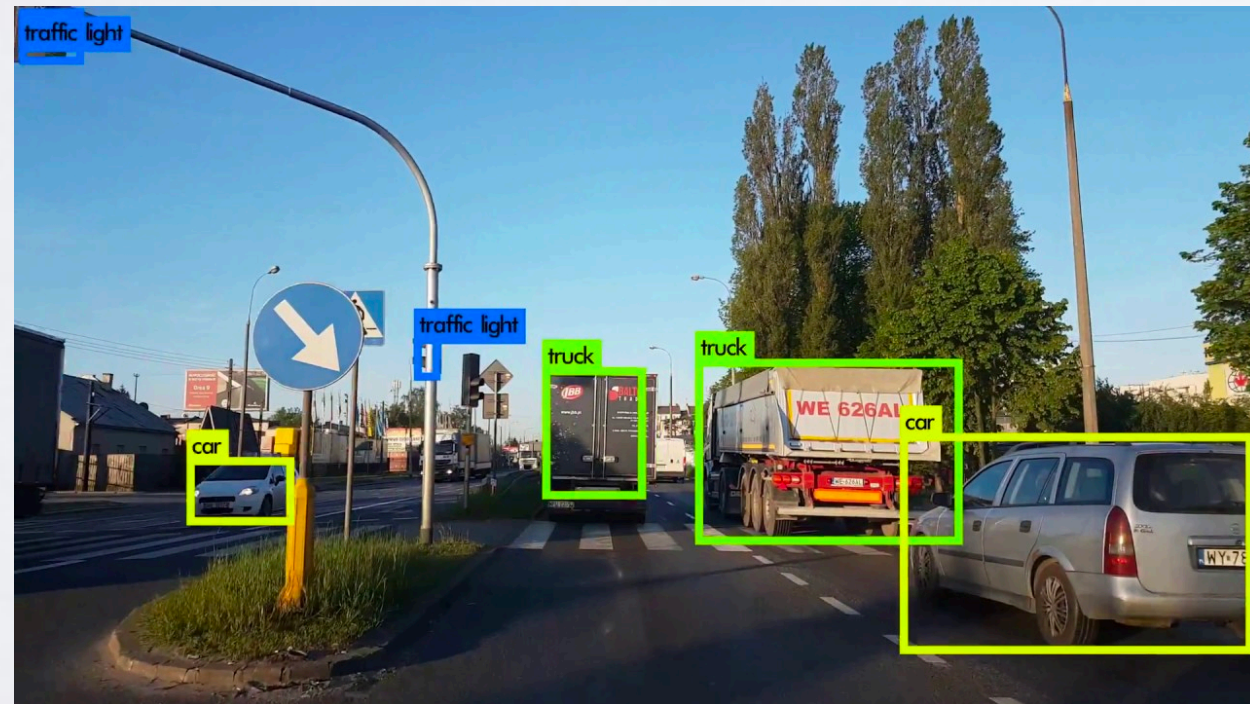
SFU

UCLA
Computer Science



The Accelerator Flowchart

Study the Application

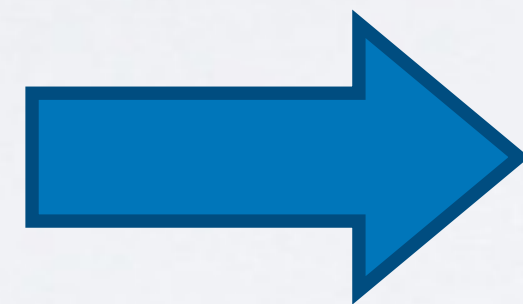


Object Detection Model



The Accelerator Flowchart

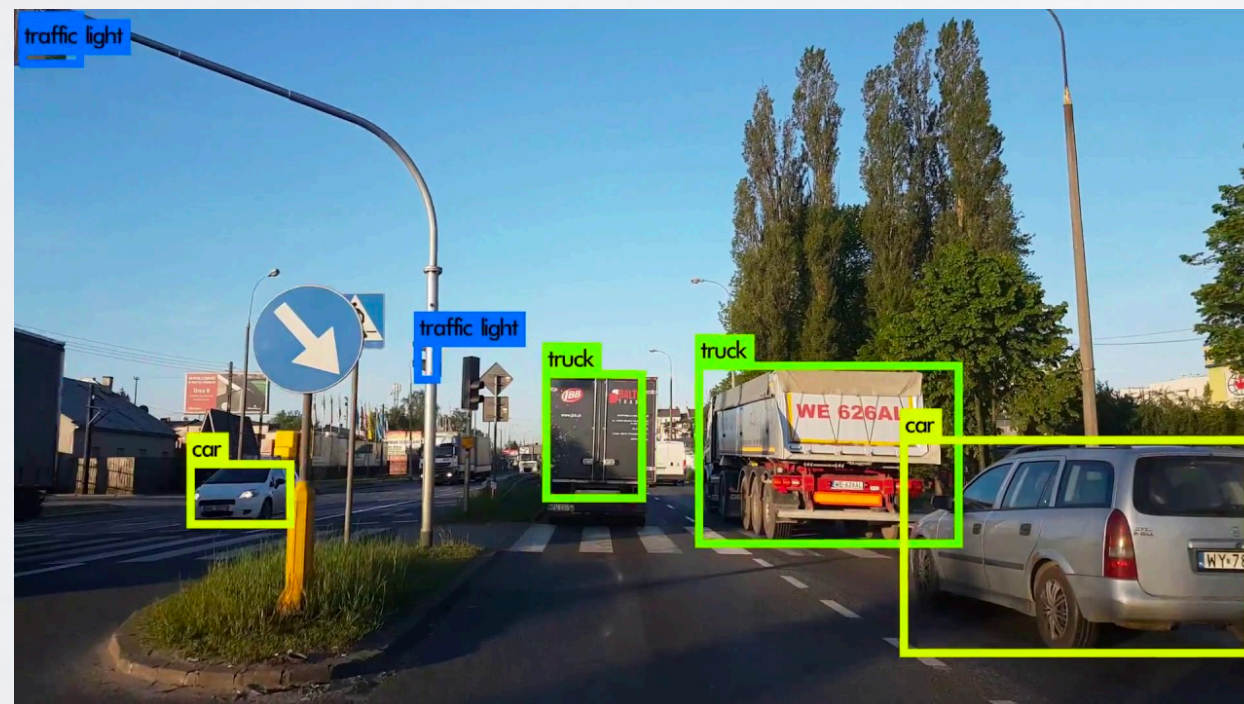
Study the Application



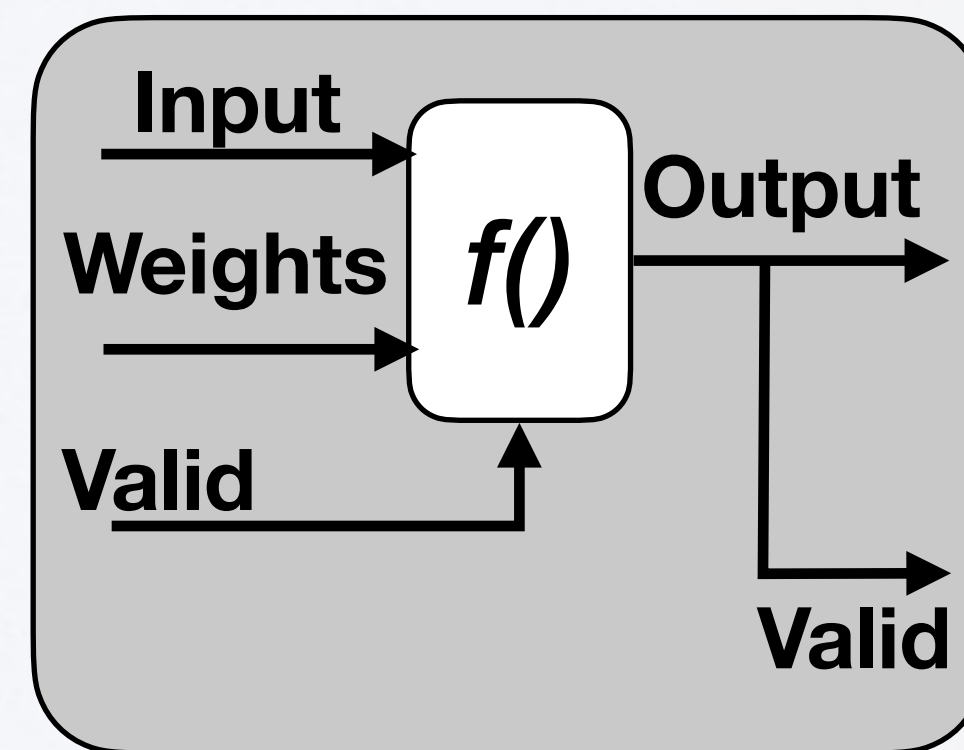
Design Hardware



RTL Design



Object Detection Model



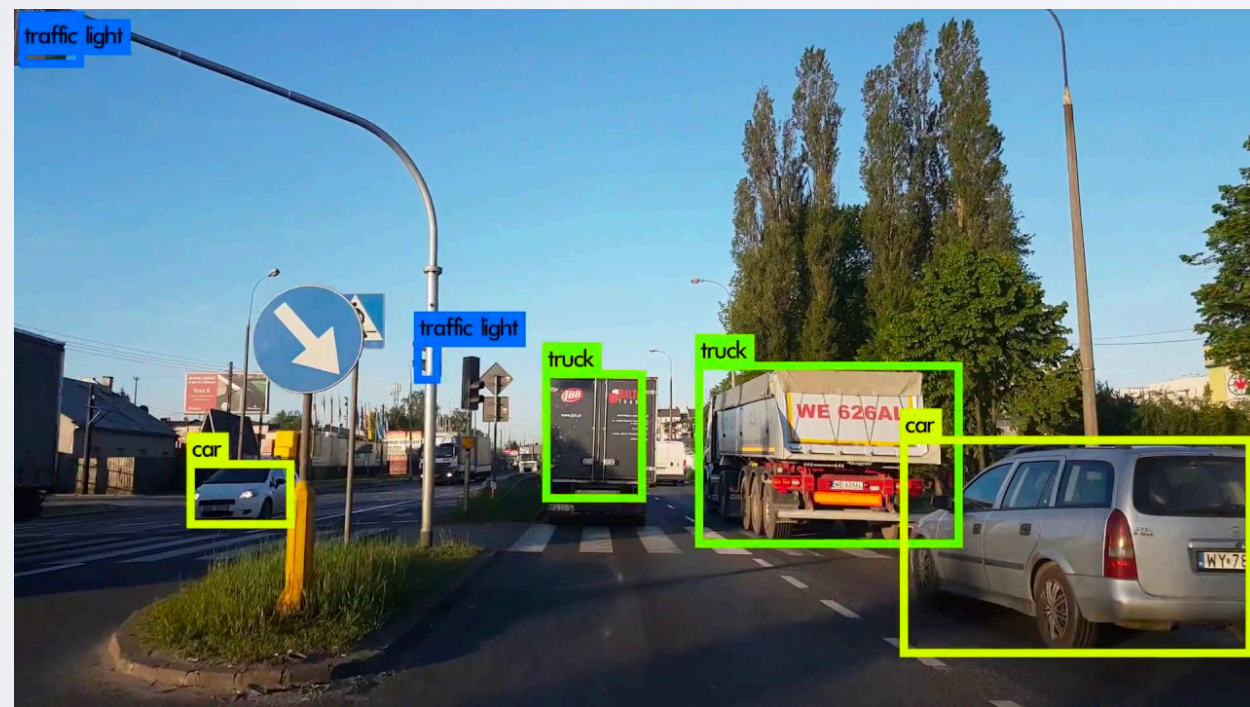
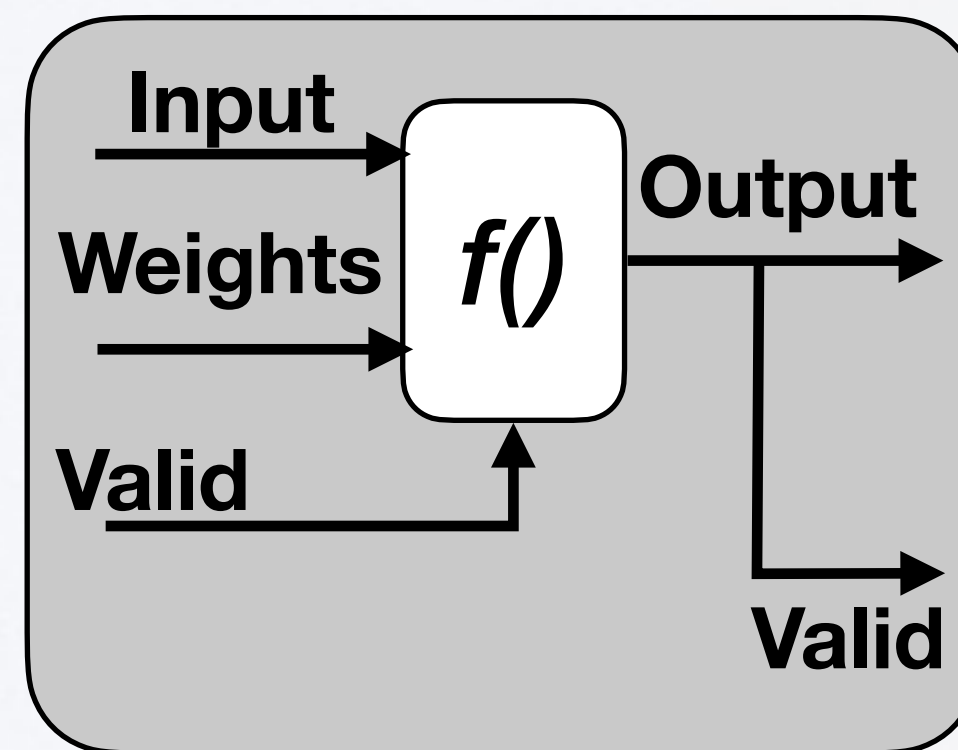
Clock
Reset



The Accelerator Flowchart



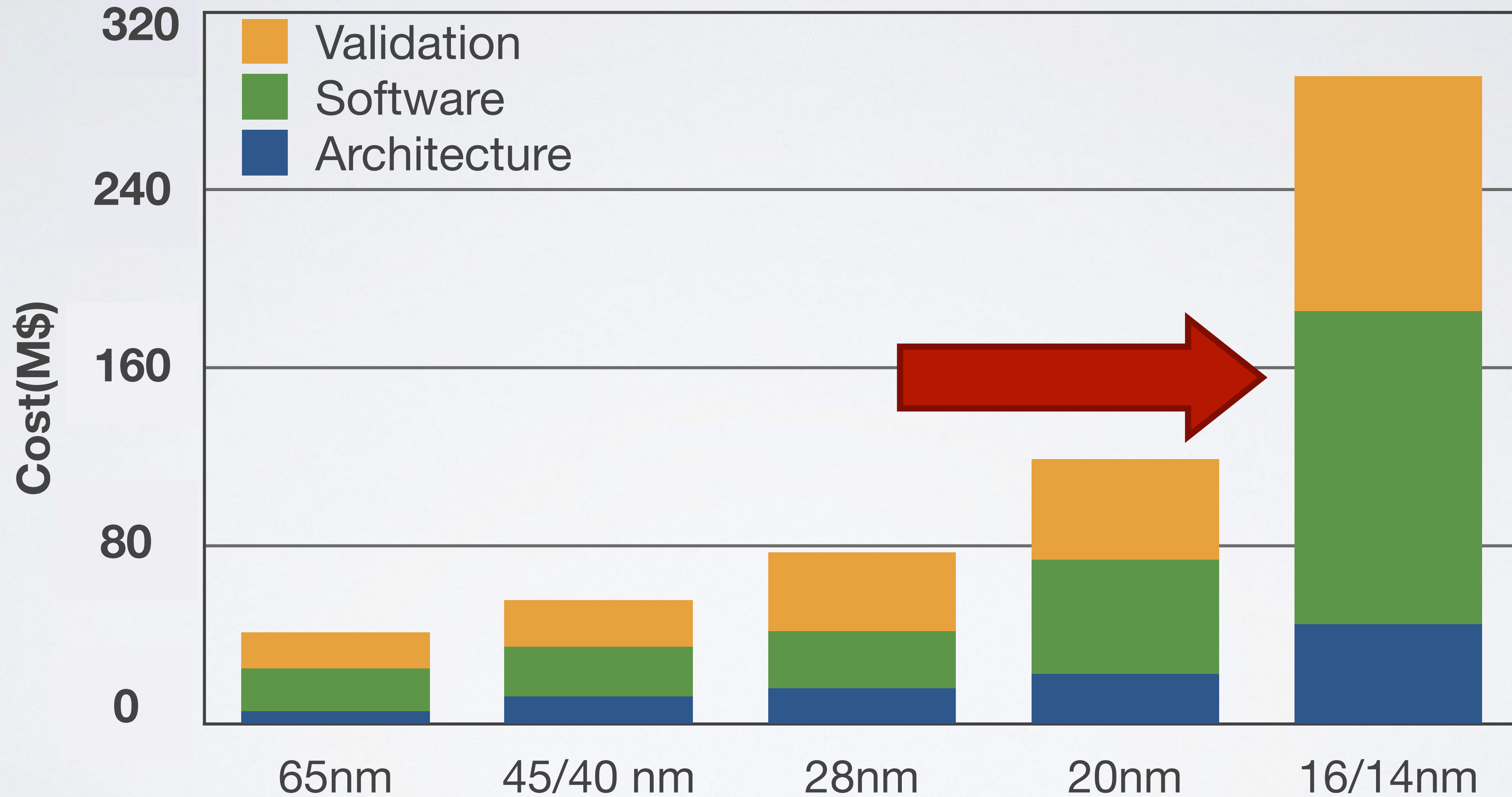
RTL Design



Object Detection Model

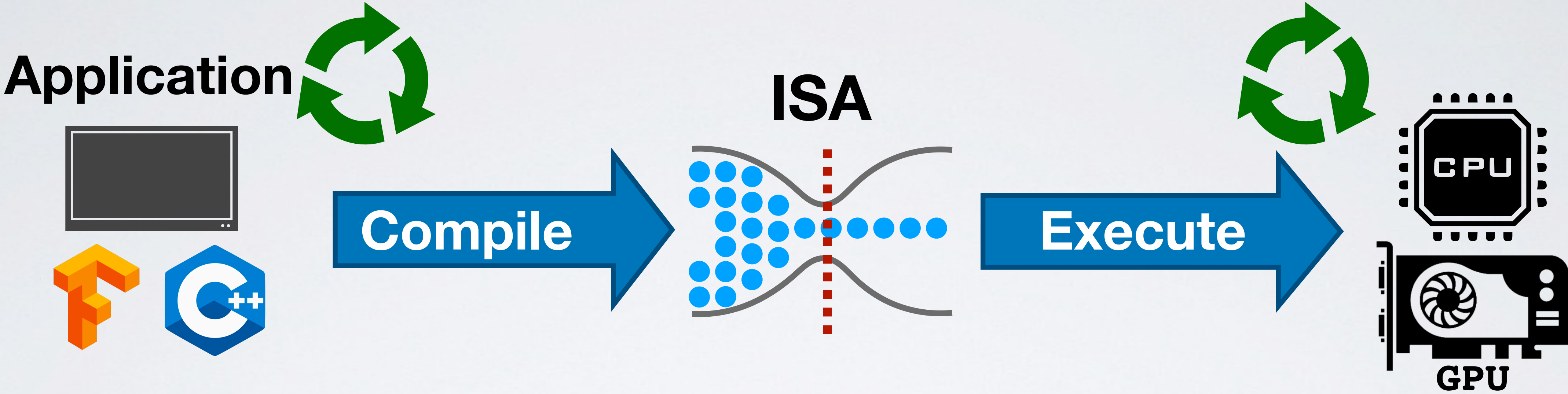
```
40  
41 # A map of names to methods that help  
42 MODEL_BUILD_UTIL_MAP = {  
43     'get_configs_from_pipeline_file':  
44         config_util.get_configs_from_  
45     'create_pipeline_proto_from_conf  
46         config_util.create_pipeline_p  
47     'save_pipeline_state_with_conf'
```

Problems With Design Flowchart



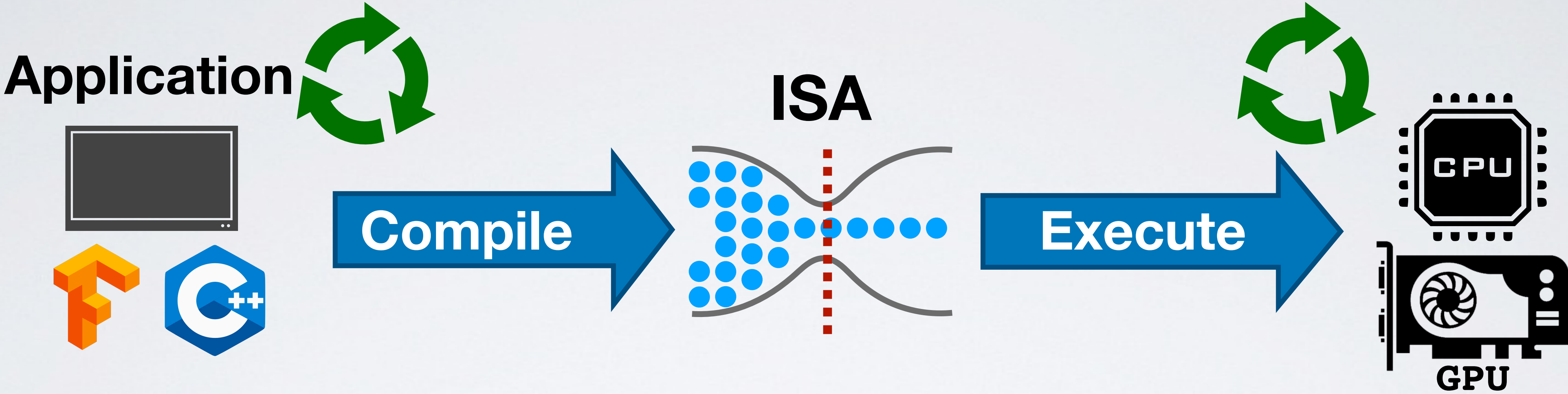
Feature Dimension (Transistor Count) — Source: IBS

ISA-based Flowchart



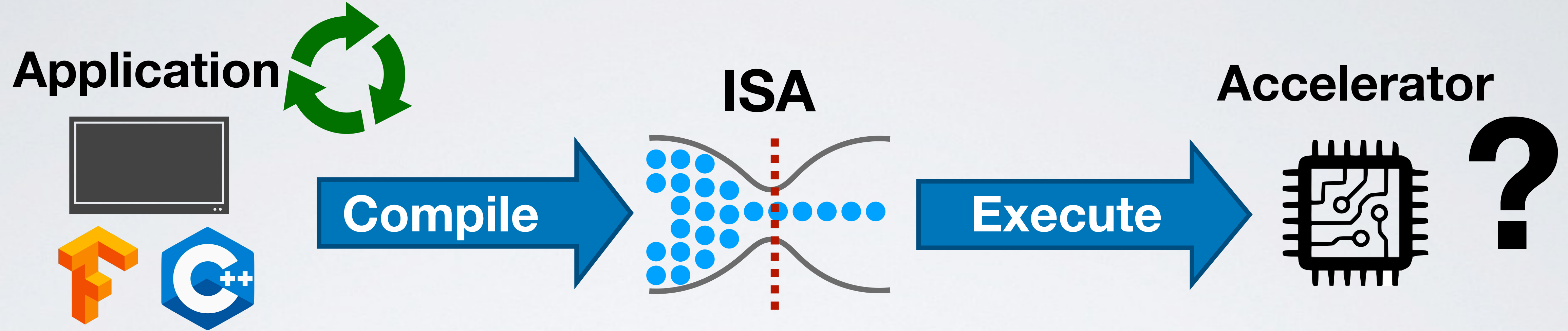
✓ Isolate Application from Architecture

ISA-based Flowchart



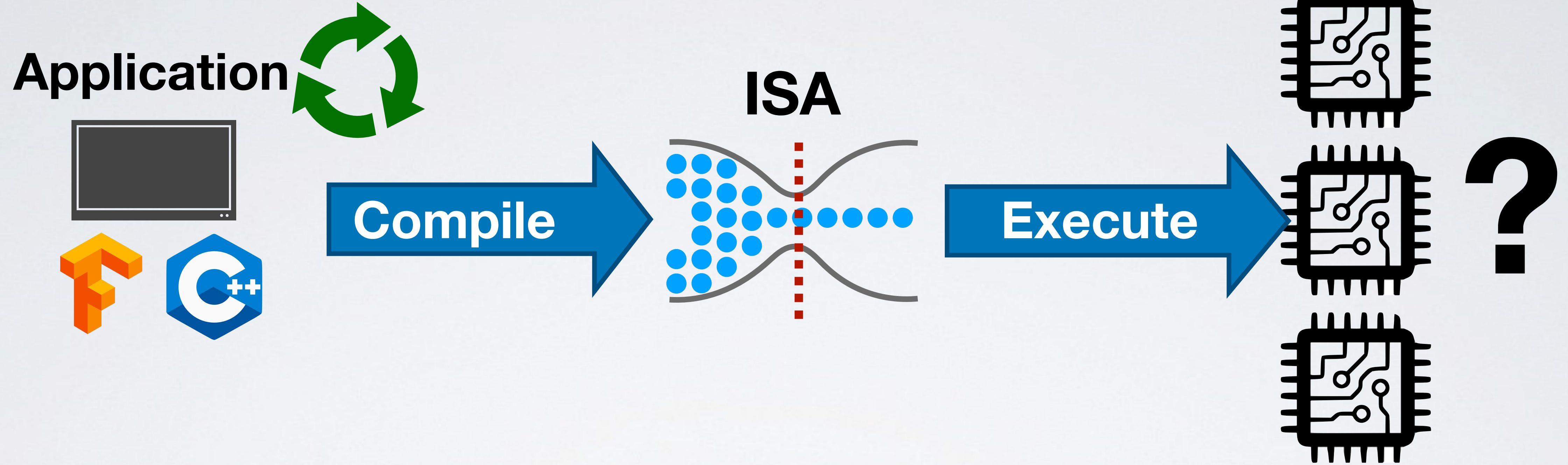
✓ Isolate Application from Architecture

ISA-based Flowchart



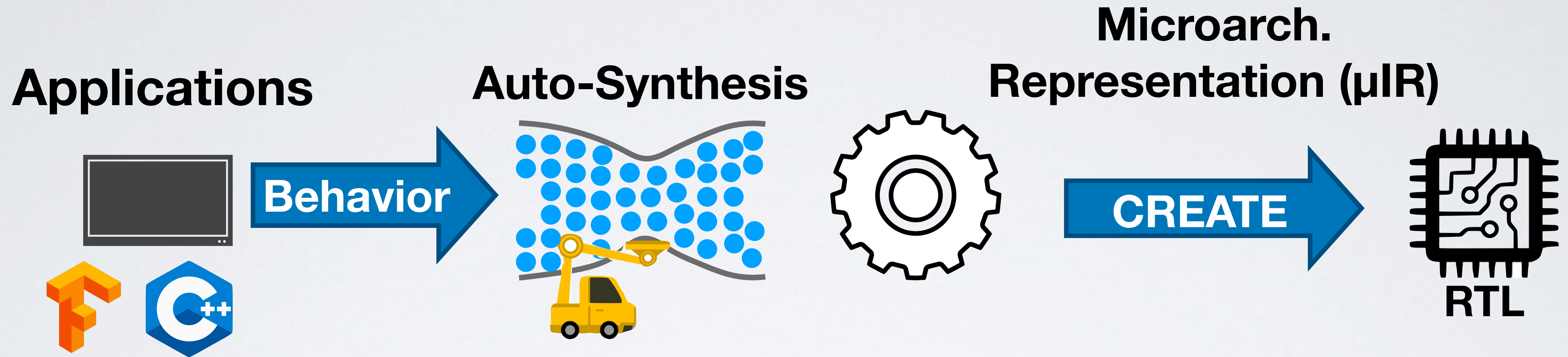
- ✓ Isolate Application from Architecture
- ▶ Not expressive enough to create hardware

ISA-based Flowchart



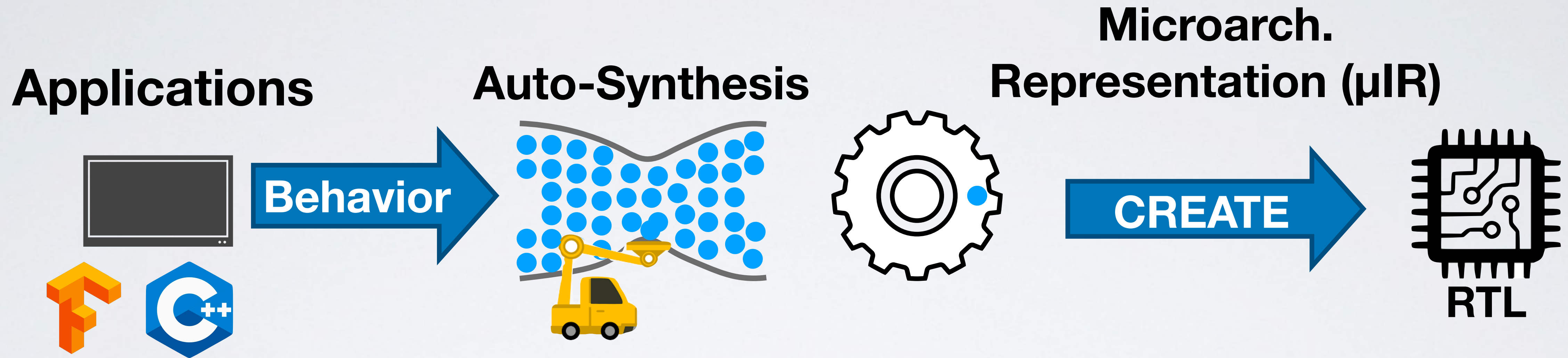
- ✓ Isolate Application from Architecture
- ▶ Not expressive enough to create hardware
- ▶ Not precise enough to explore hardware

μ IR — A New Accelerator Flowchart



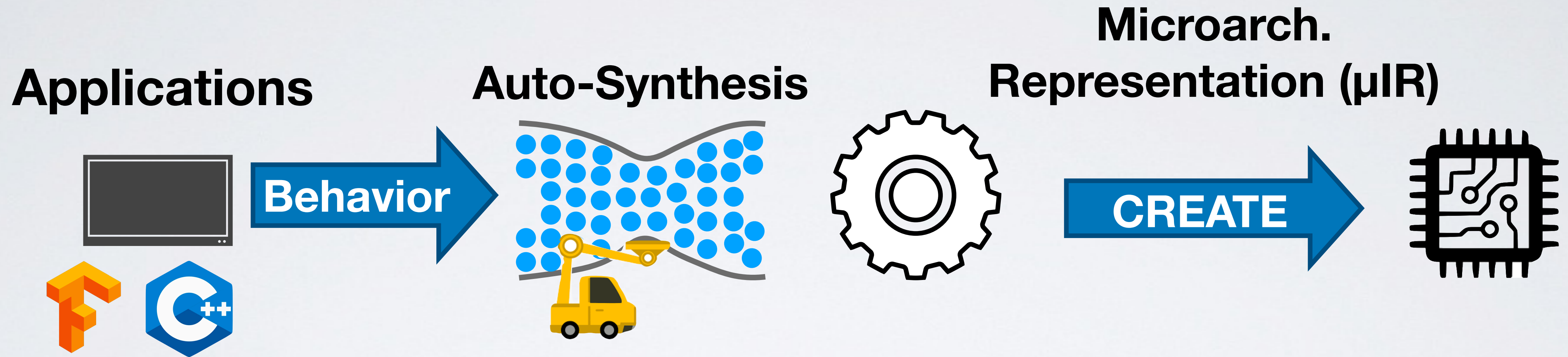
✓ End-to-End flow — Existing software for behavior/functionality

μ IR — A New Accelerator Flowchart



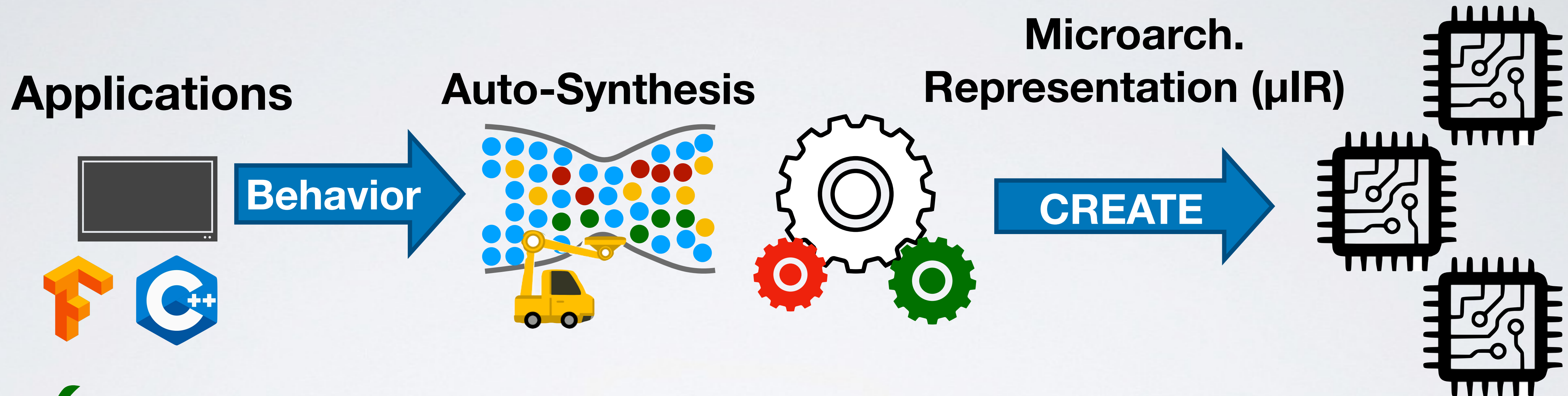
- ✓ End-to-End flow — Existing software for behavior/functionality
- ✓ Reduce effort — Compiler to extract behaviour

μ IR — A New Accelerator Flowchart



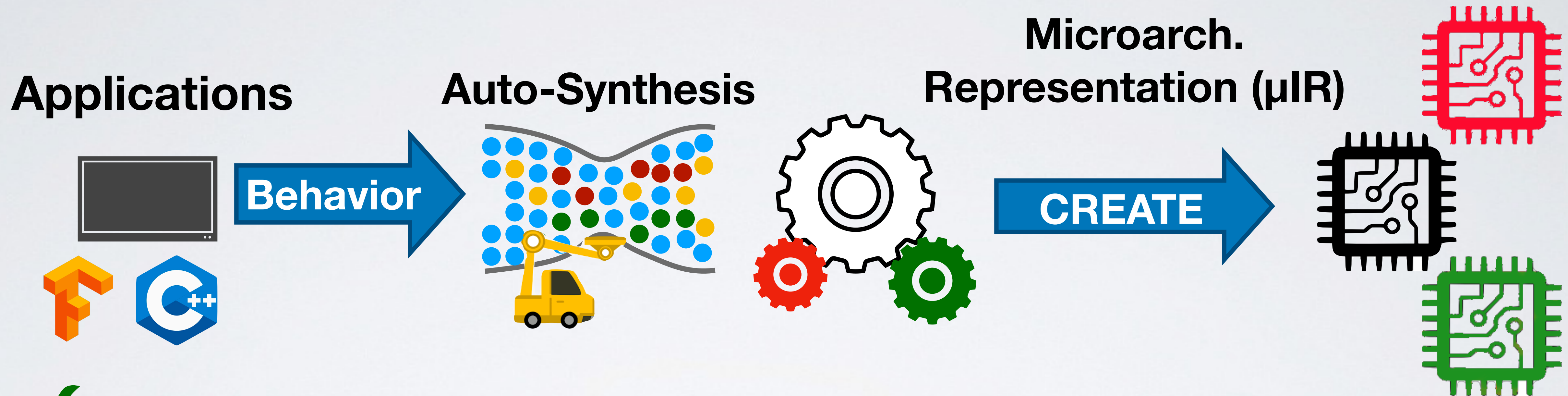
- ✓ End-to-End flow — Existing software for behavior/functionality
- ✓ Reduce effort — Compiler to extract behaviour
- ✓ Design exploration — New model for exploring architectures

μ IR — A New Accelerator Flowchart



- ✓ End-to-End flow — Existing software for behavior/functionality
- ✓ Reduce effort — Compiler to extract behaviour
- ✓ Design exploration — New model for exploring architectures
- ✓ Extensibility — Extensible to capture domain information

μ IR — A New Accelerator Flowchart

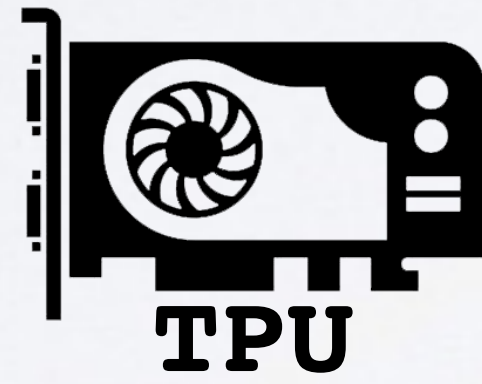
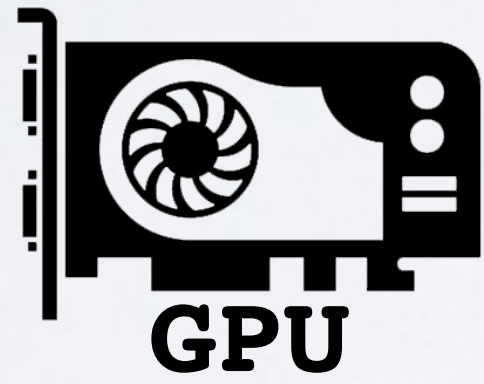
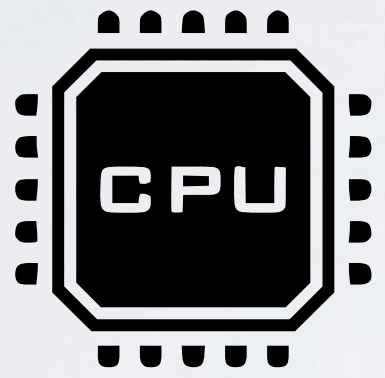


- ✓ End-to-End flow — Existing software for behavior/functionality
- ✓ Reduce effort — Compiler to extract behaviour
- ✓ Design exploration — New model for exploring architectures
- ✓ Extensibility — Extensible to capture domain information

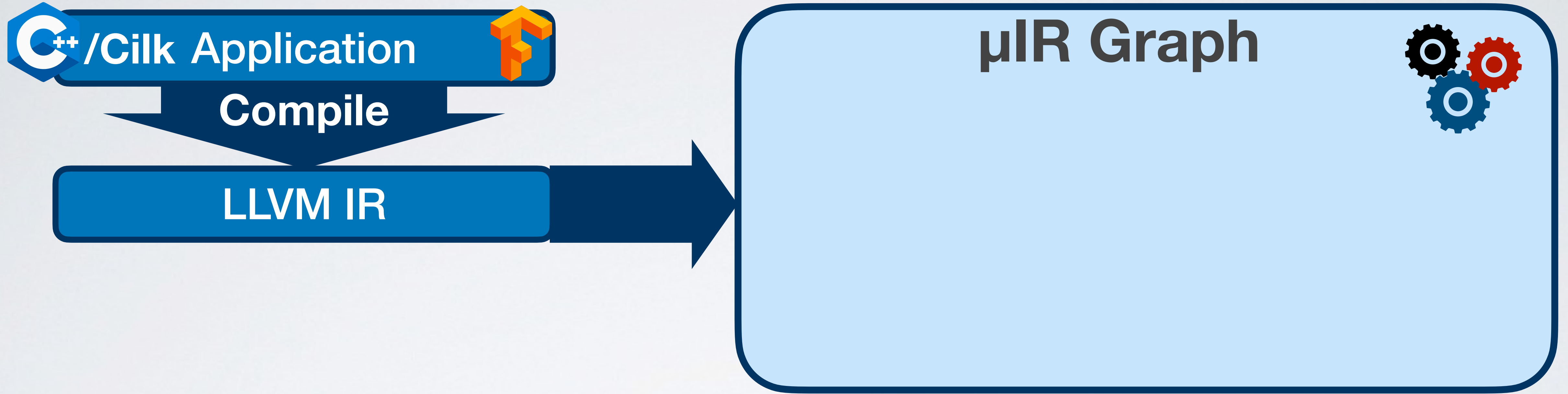
 /Cilk Application 

Compile

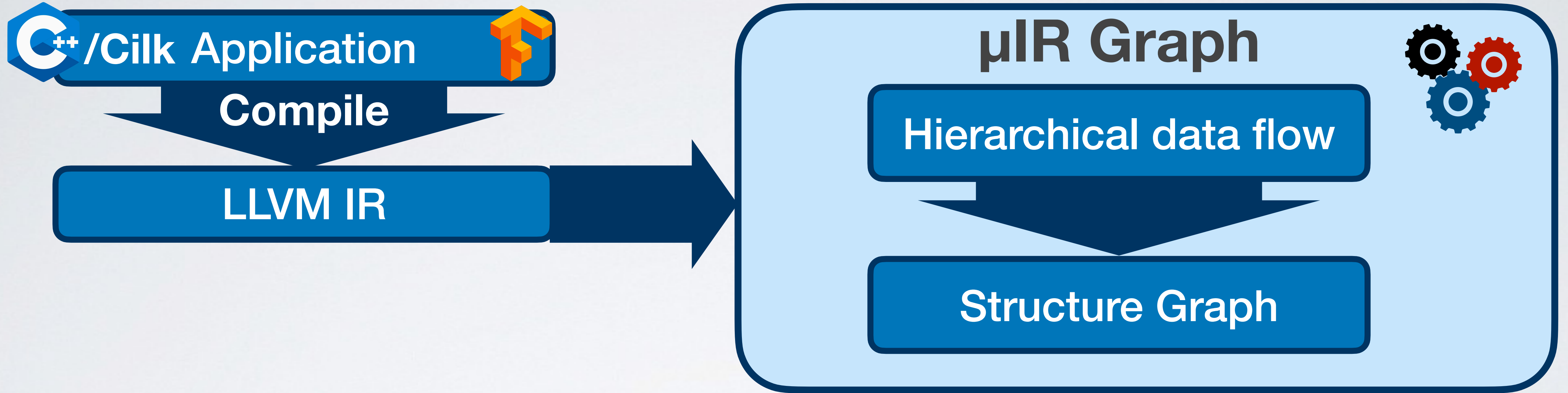
LLVM IR



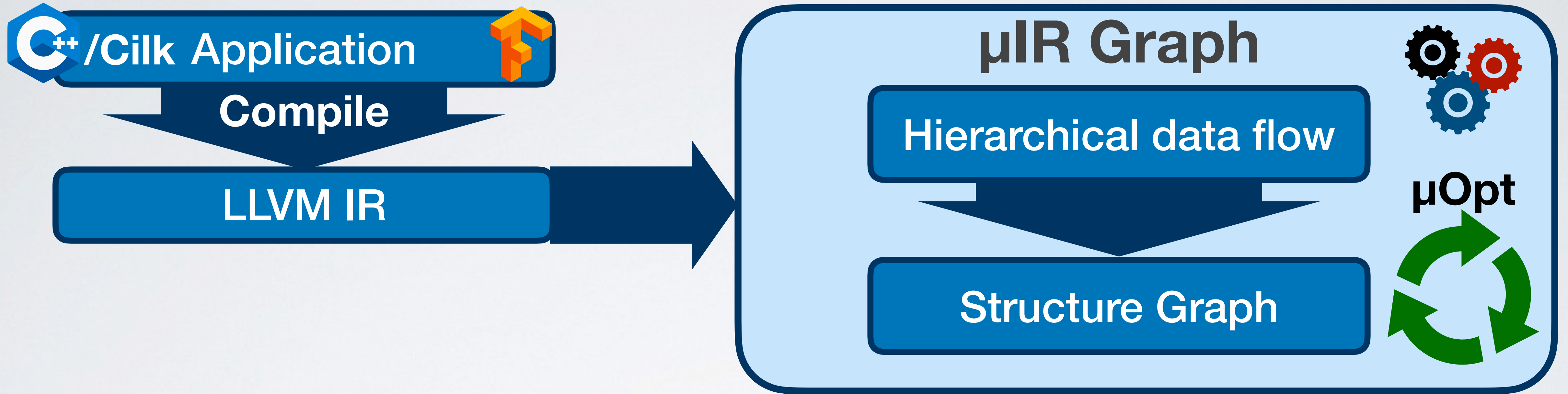
μIR A New Accelerator Flowchart



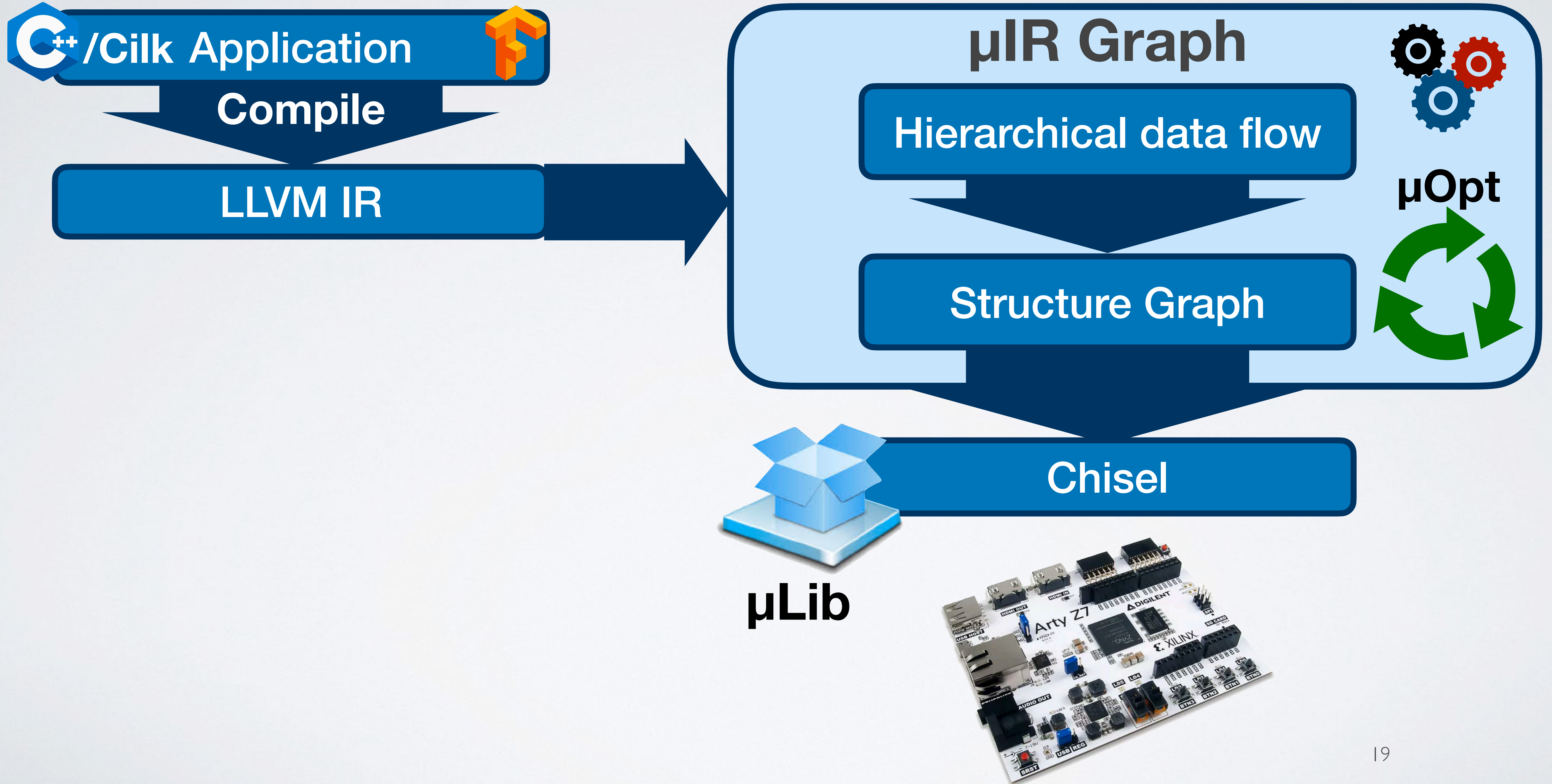
μIR A New Accelerator Flowchart



μIR A New Accelerator Flowchart

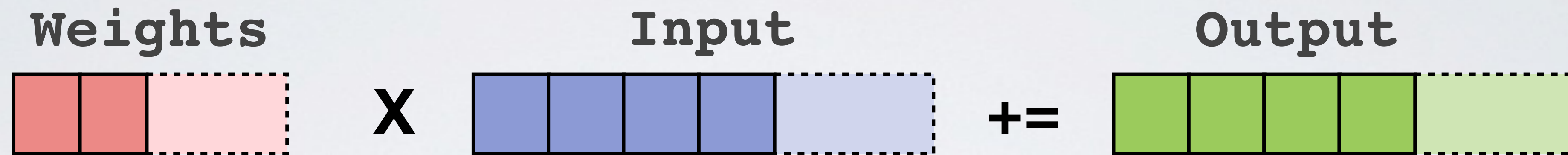


μIR A New Accelerator Flowchart



- Motivation
- **μ IR behaviour graph**
- μ IR structural graph
- Evaluation
- Summary

μ IR Behavioural Graph (1/3)



```
parallel_for(i=0; I<(M-W); i++)  
  parallel_for(j=0; j<W; j++)  
    output[i] += input[i+j] * weight[j];
```

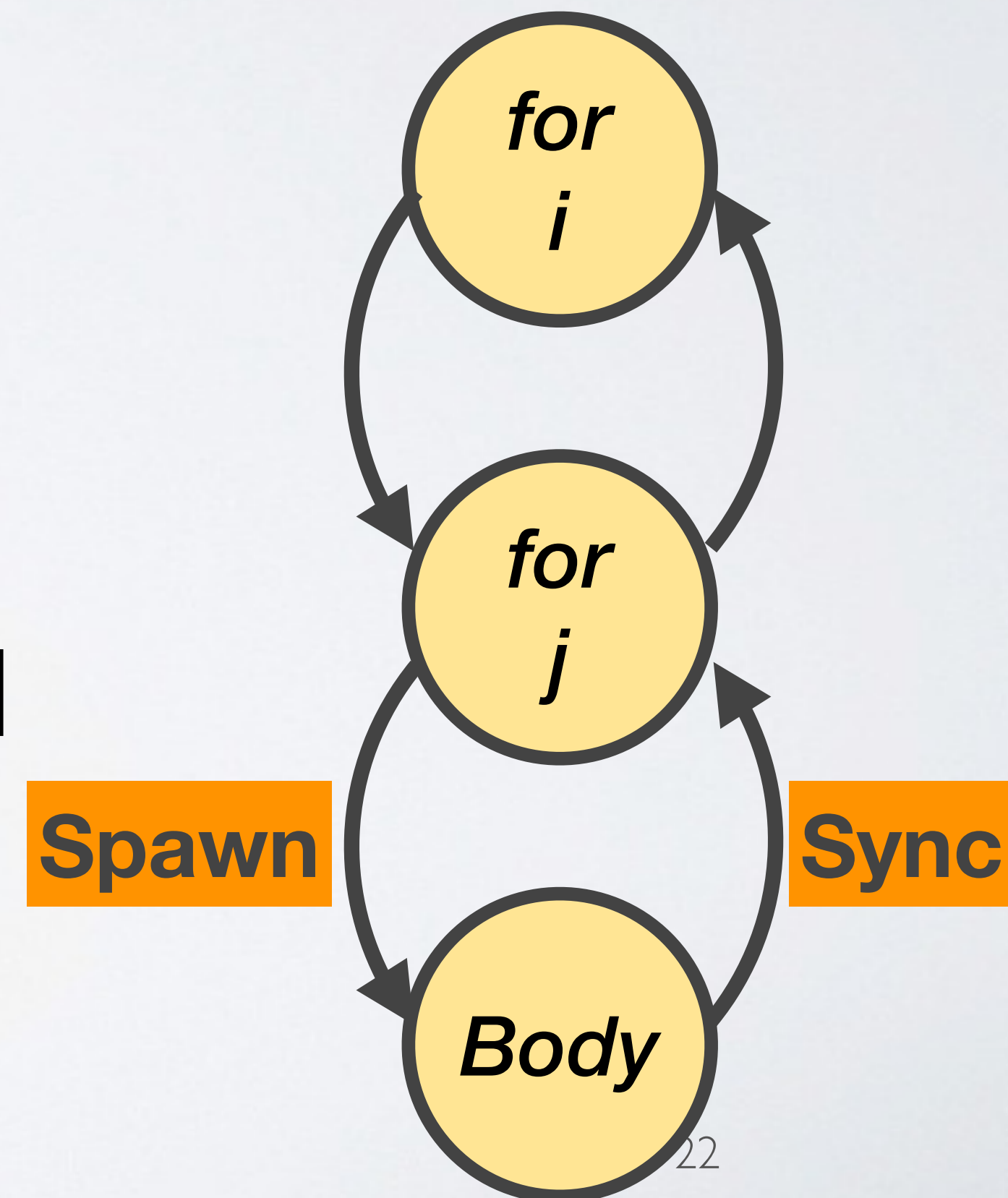
$$Conv = \sum_{i=0}^{N-W} \sum_{j=0}^W W_{ij} * IN_{ij}$$

μ IR Behavioural Graph (2/3)

```
parallel_for(i=0; I<(M-W); i++)  
  parallel_for(j=0; j<W; j++)  
    output[i] += input[i+j] * weight[j];
```

- **Hierarchical and heterogeneous task graph**

- ▶ Decompose the input program to task blocks.
- ▶ Graph of task blocks
 - Implement arbitrary heterogeneous parallel and serial

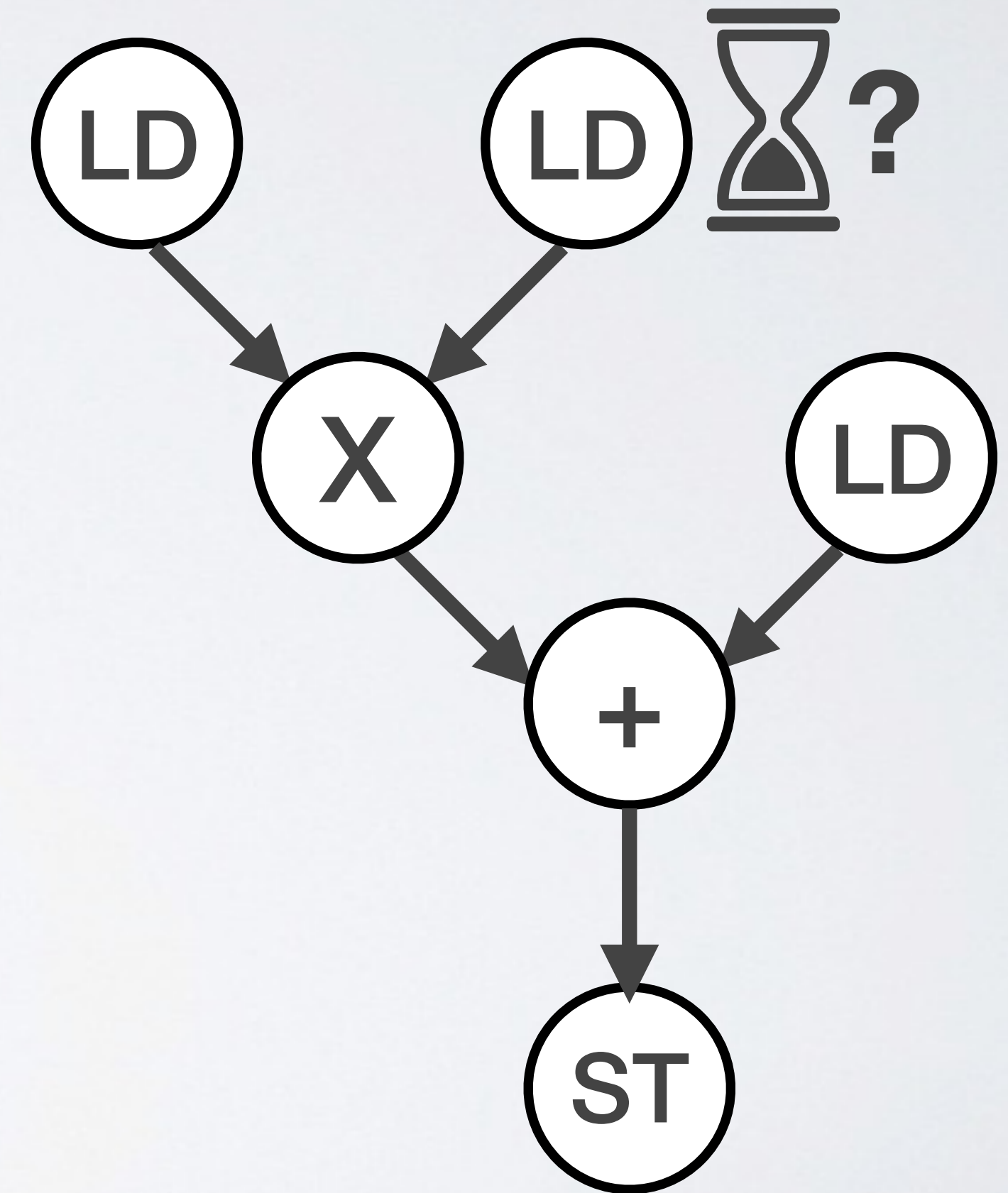


μIR Behavioural Graph (3/3)

```
output[i] += input[i+j] * weight[j];
```

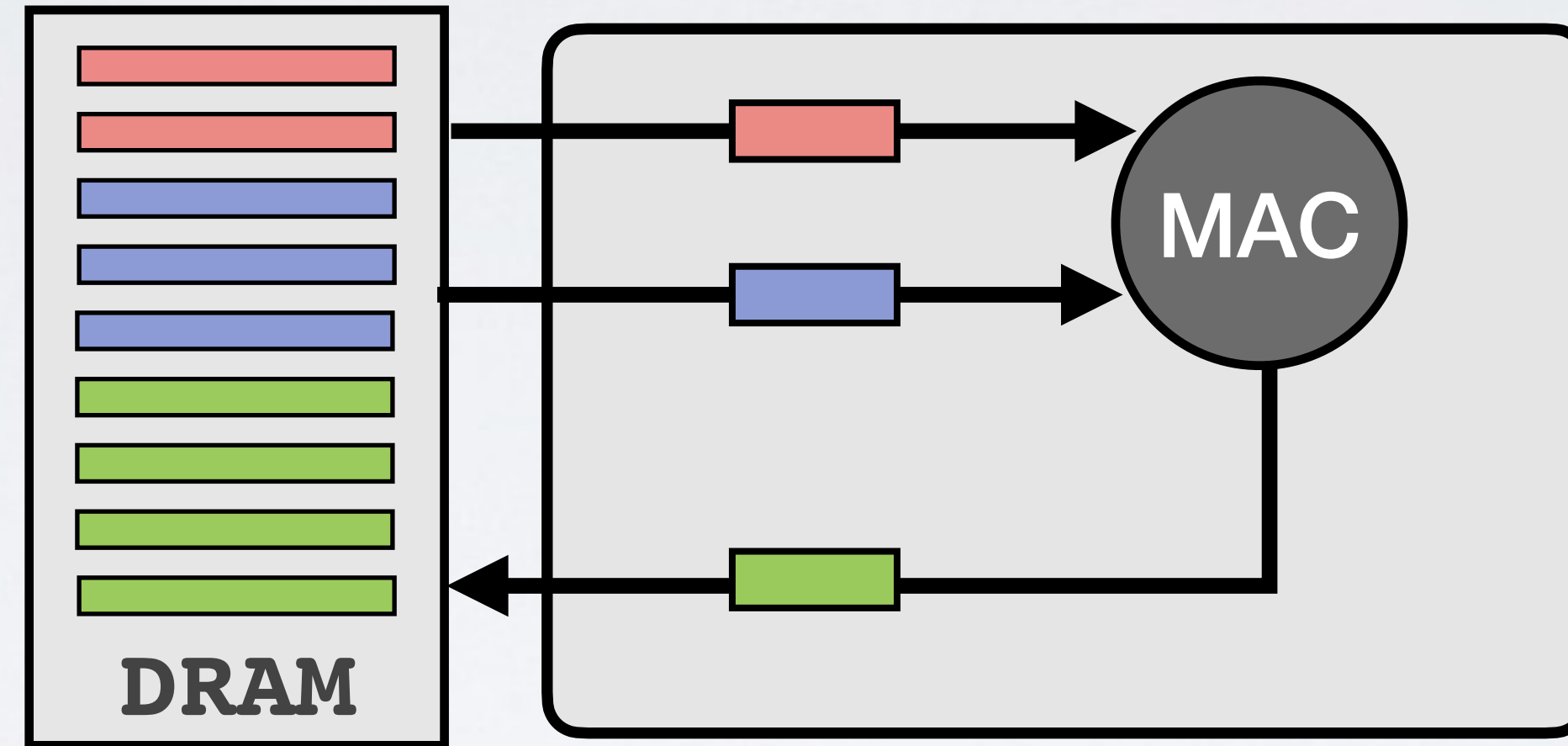
• μIR Dataflow graph for each task block

- ▶ Pipelined (single and multi-cycle)
- ▶ Typed (e.g., FP precision)
- ▶ Non-deterministic cycle ops (e.g., Mem ops)
- ▶ Permits local transformations because tasks are asynchronous

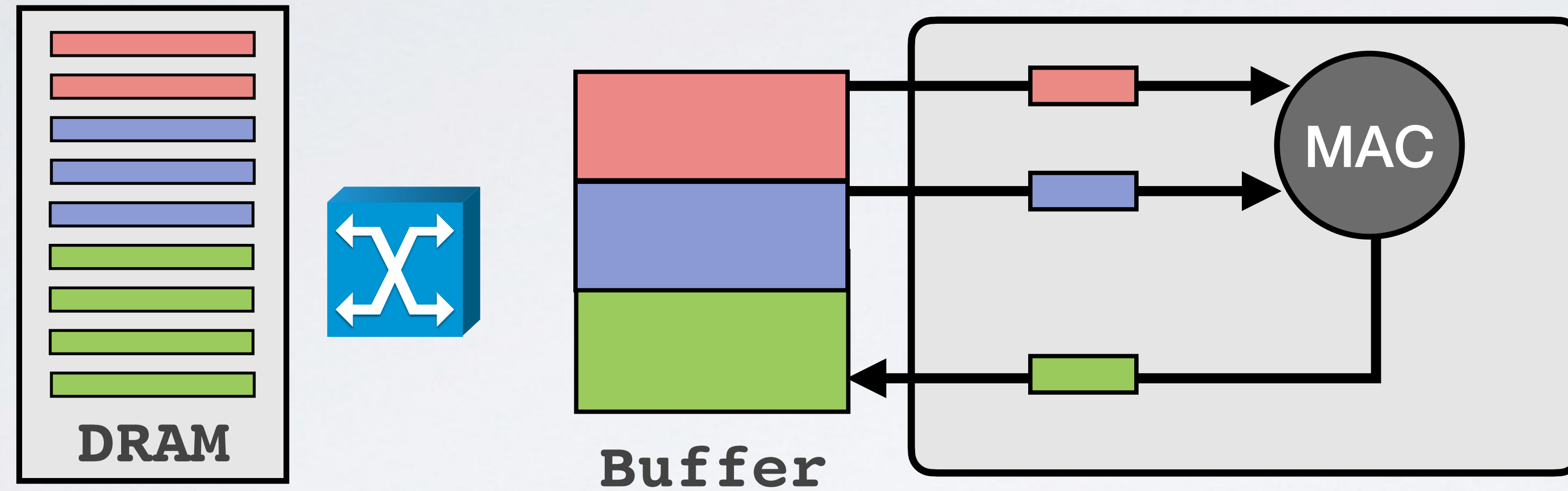


- Motivation
- μ IR behaviour graph
- **μ IR structural graph**
- Evaluation
- Summary

μIR Structural Graph : Stage 0 — Unoptimized

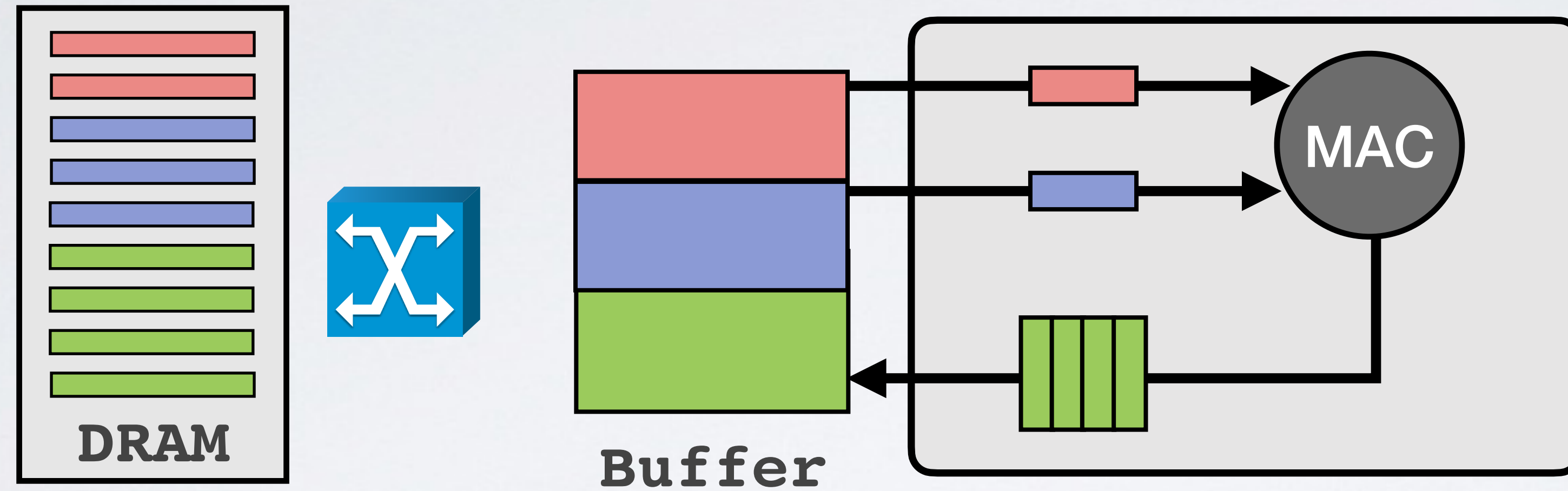


μ IR Structural Graph : Stage 1 — Locality



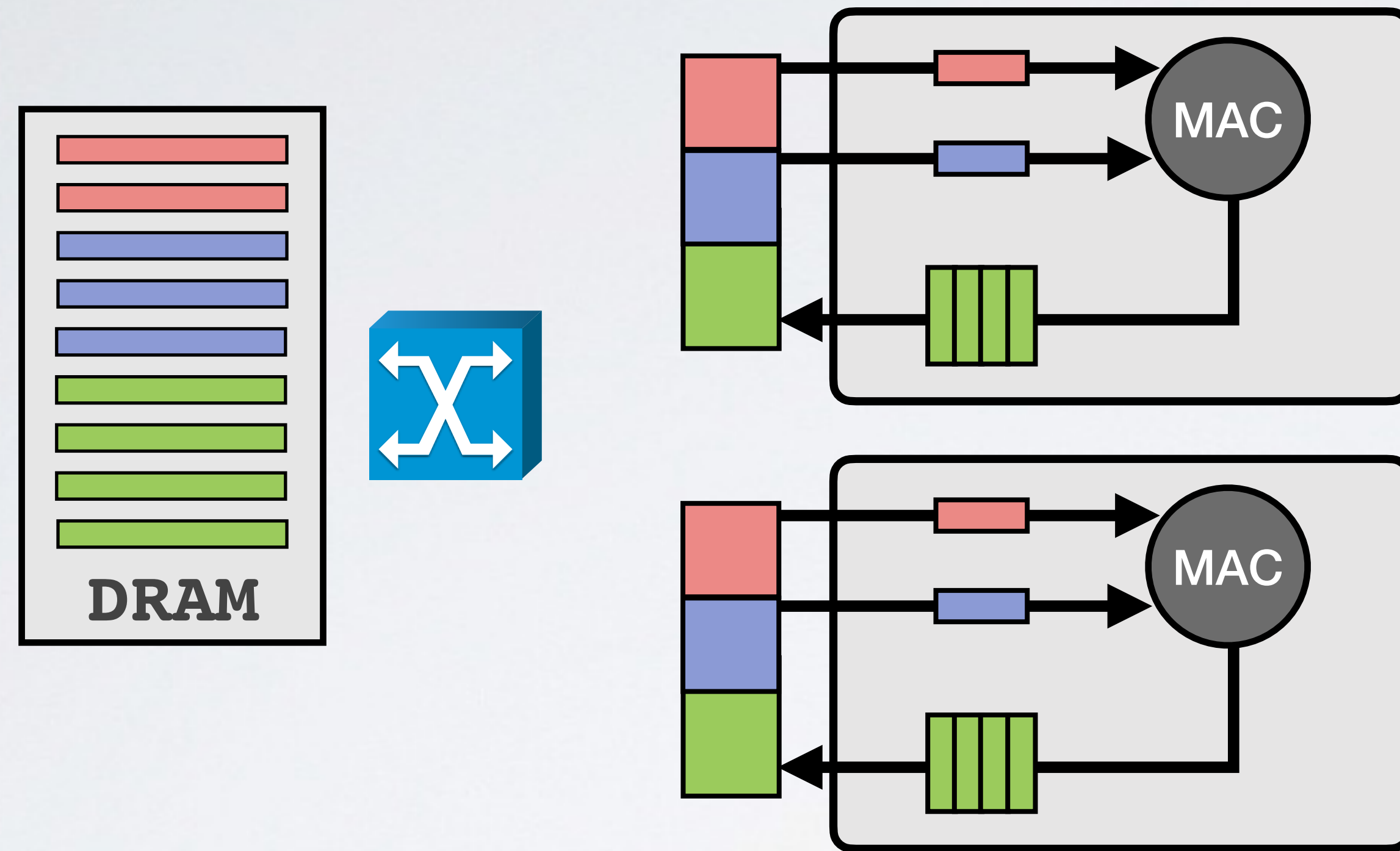
- Define custom memory hierarchy — Buffer, FIFO

μ IR Structural Graph : Stage 1 — Locality



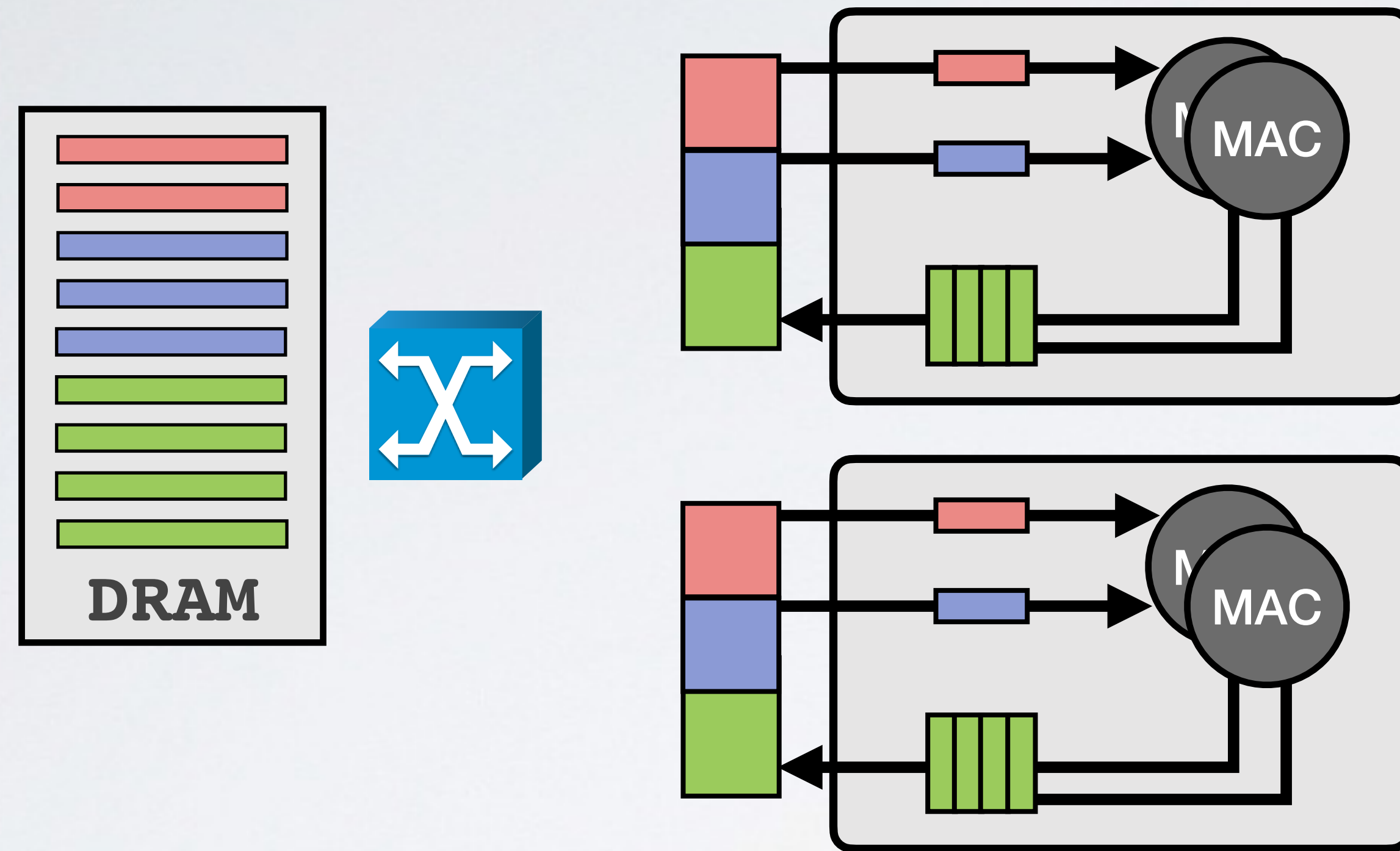
- Define custom memory hierarchy — Buffer, FIFO

μ IR Structural Graph : Stage 2 — Tiling



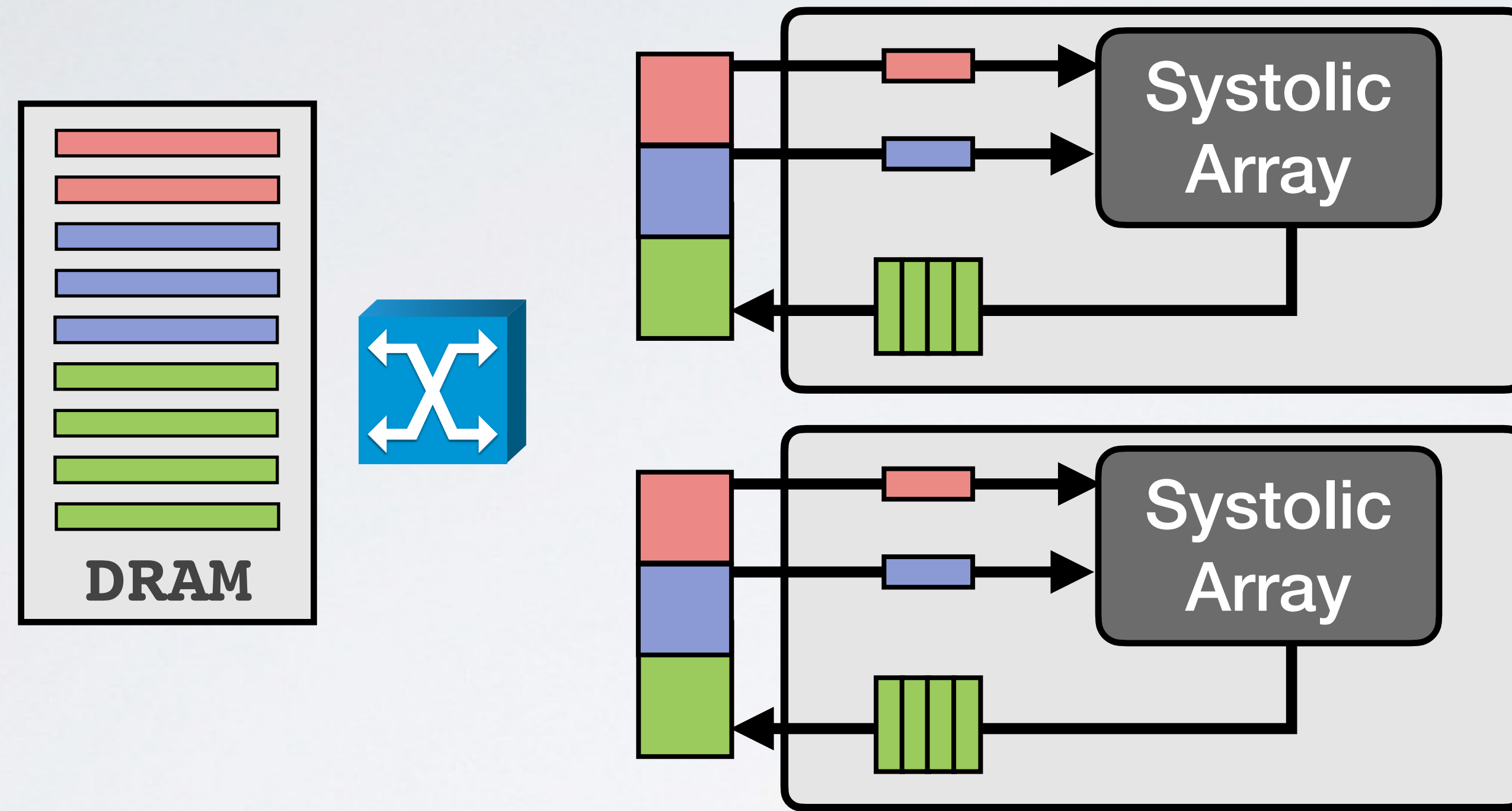
- Define custom memory hierarchy — Buffer, FIFO
- Tiling asynchronous task blocks

μ IR Structural Graph : Stage 3 — Pipelining



- Define custom memory hierarchy — Buffer, FIFO
- Tiling asynchronous task blocks
- Start pipelining the operands

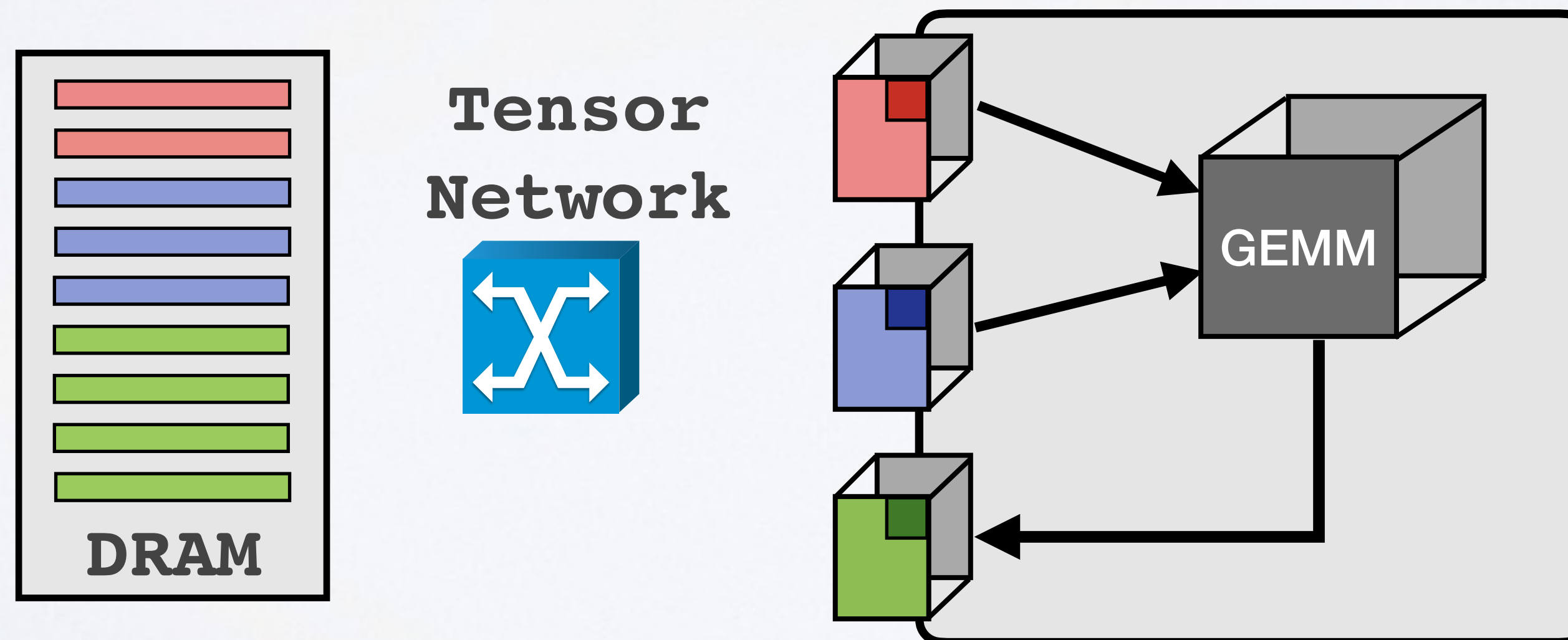
μ IR Structural Graph : Stage 3 — Pipelining



- Define custom memory hierarchy — Buffer, FIFO
- Tiling asynchronous task blocks
- Start pipelining the operands
- Define custom operator

μ IR Structural Graph : Stage 4 Higer-Order Ops

- Extensible IR
 - Introduce new ops and new types
- Existing stages/transformations have to work
- μ IR components are generic: Dataflow nodes, buffers and memory network



- Motivation
- μ IR behaviour graph
- μ IR structural graph
- **Evaluation**
- Summary

μIR Targets

- **FPGA**

- Arria 10 Soc
- Scripted process



<https://github.com/sfu-arch/uir-fpga>

- **Cycle accurate simulation (e.g., gem5)**

- C++ driver and a Python binding for ease of use
- Cache-based memory interface



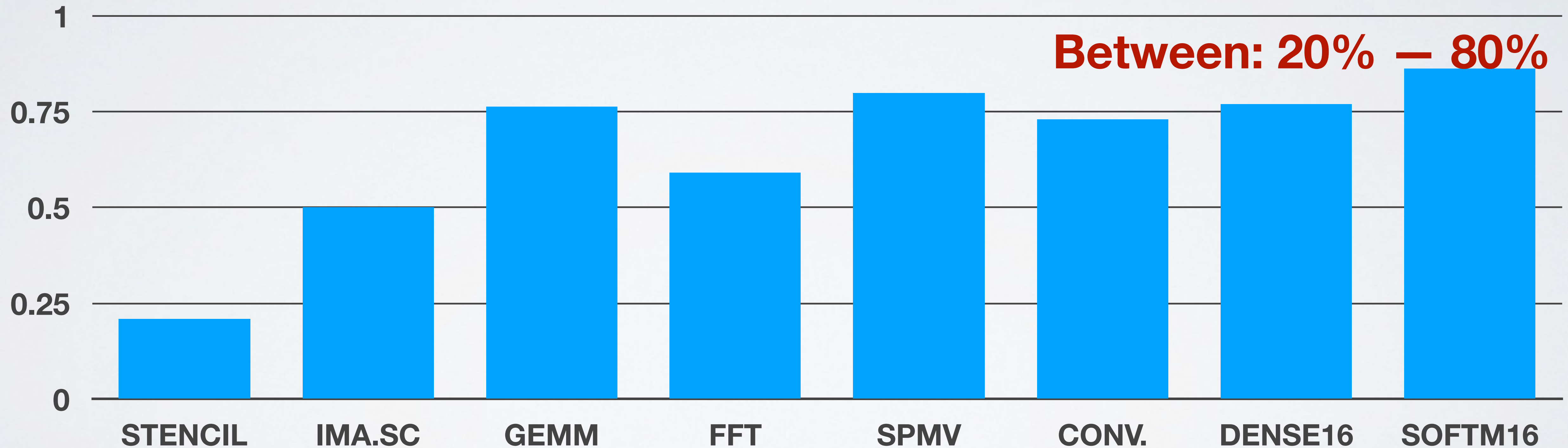
<https://github.com/sfu-arch/uir-sim>

μ IR Iterative Optimizations — Stacking

- Normalized to the baseline accelerator

Banking, Localization, Op-Fusion, Tiling

Between: 20% — 80%



μIR Productivity

	#Changes μIR	#Changes Firrtl	Ratio
Saxpy	9	73	9.3
Stencil	12	142	12.4
Image Sca.	10	84	8.4

Available Now!



<https://github.com/sfu-arch/uir>



<https://github.com/sfu-arch/uir-sim>



<https://github.com/sfu-arch/uir-fpga>



<https://github.com/sfu-arch/uir-docker>

Thanks to Open-Source community!