# FLIP2M: Flexible Intra-layer Parallelism and Inter-layer Pipelining for Multi-model AR/VR Workloads

#### **Gabriele Tombesi**

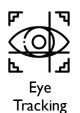
Je Yang Joseph Zuckerman Davide Giri William Baisi Luca Carloni





• AR/VR applications require multiple concurrent DNNs running on the same SoC:

- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc









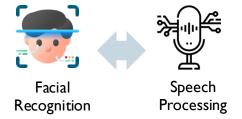
Facial Recognition



Speech Processing

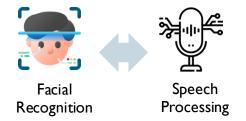
- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc
  - Complex intra- and inter-task dependencies (both static/dynamic)





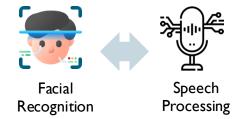
- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc
  - Complex intra- and inter-task dependencies (both static/dynamic)
  - Often coming with tight real-time constraints





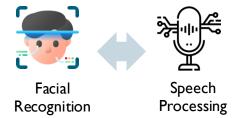
- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc
  - Complex intra- and inter-task dependencies (both static/dynamic)
  - Often coming with tight real-time constraints
- Main Execution Challenge → High heterogeneity in:





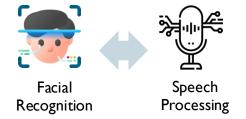
- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc
  - Complex intra- and inter-task dependencies (both static/dynamic)
  - Often coming with tight real-time constraints
- Main Execution Challenge → High heterogeneity in:
  - Layer shapes / operations





- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc
  - Complex intra- and inter-task dependencies (both static/dynamic)
  - Often coming with tight real-time constraints
- Main Execution Challenge → High heterogeneity in:
  - Layer shapes / operations
  - Intra-model /

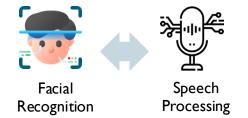




Model	Per-Layer Activation-Weight Size Ratio	Layer Operations
MobileNet-v2 [39]	Min: 0.0057 Median: 2.04, Max 784	Conv2D, PW-Conv, DW-Conv, Residual

- AR/VR applications require multiple concurrent DNNs running on the same SoC:
  - Object Detection, Hand Tracking, Segmentation, etc
  - Complex intra- and inter-task dependencies (both static/dynamic)
  - Often coming with tight real-time constraints
- Main Execution Challenge → High heterogeneity in:
  - Layer shapes / operations
  - Intra-model / Inter-model





Model	Per-Layer Activation-Weight Size Ratio	Layer Operations
MobileNet-v2 [39]	Min: 0.0057 Median: 2.04, Max: 784	Conv2D, PW-Conv, DW-Conv, Residual
ResNet-50 [19]	Min: 0.0106, Median: 0.68, Max: 49	Conv2D, FullyConnected, Residual
Vgg16 [40]	Min: 0.043, Median: 1.36, Max: 87.111	Conv2D, FullyConnected
SqueezeNet [20]	Min: 0.073, Median: 2.64, Max: 189.06	Conv2D, FullyConnected, Residual
Unet [38]	Min:0.13, Median: 8.88, Max: 568	Conv2D, FullyConnected, UpConv, Concat
MobileBert [43]	Min: 0.125, Median: 0.5, Max 8	Matrix-matrix multiplication, Residual

#### Background – Tiled Architectures

**AR/VR** Heterogeneity →

- Over-specialized hardware: latency/energy inefficiency across diverse DNN layers.
- Model concurrency: amplified range of operational requirements / resource demands.

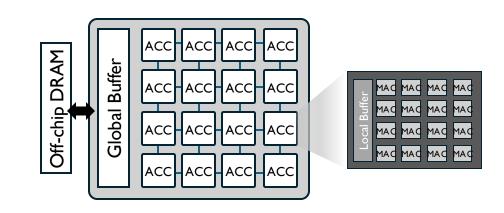
#### Background – Tiled Architectures

**AR/VR** Heterogeneity →

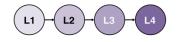
- Over-specialized hardware: latency/energy inefficiency across diverse DNN layers.
- Model concurrency: amplified range of operational requirements / resource demands.

Recent advancements in SoC design from monolithic accelerators towards **tiled architectures**:

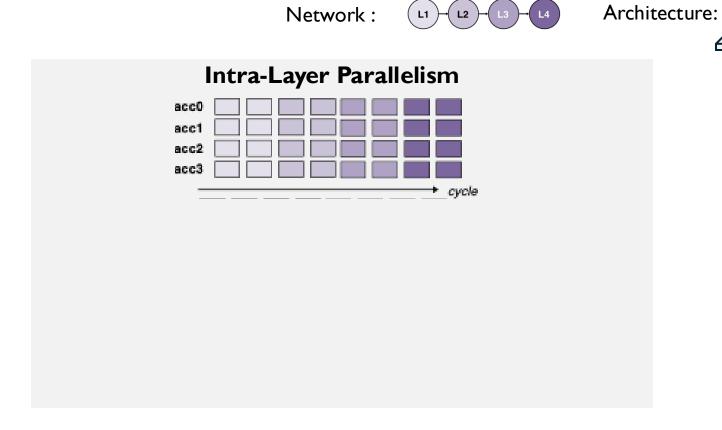
- Multiple accelerator tiles + on-chip global buffer
- Network-on-chip (NoC) based interconnect
- Multiple accelerator tiles cooperate to execute a task



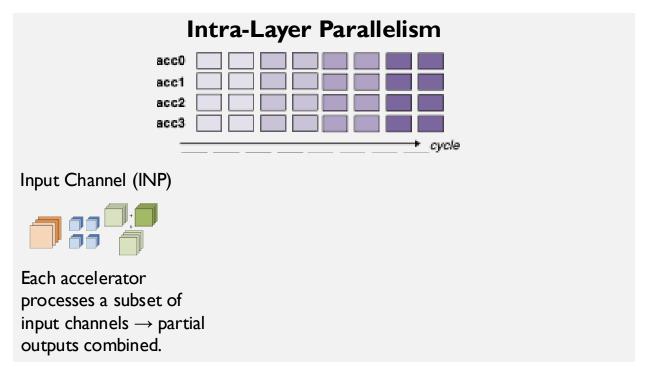
Network:







Network: L1 L2 L3 L4 Architecture: acc2 acc2 acc3



Architecture:

Intra-Layer Parallelism

acc0
acc1
acc2
acc3

Input Channel (INP)

Output Channel (OUTP)

Filters split across
processes a subset of
input channels → partial

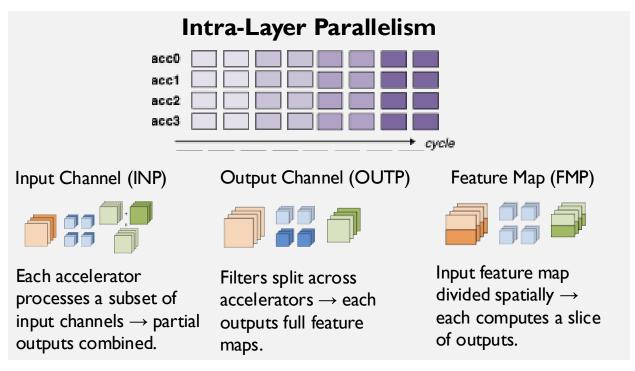
outputs full feature

maps.

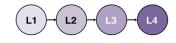
outputs combined.

Network:

Network: L1 L2 L3 L4 Architecture: acc0 acc1 acc2 acc3

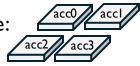


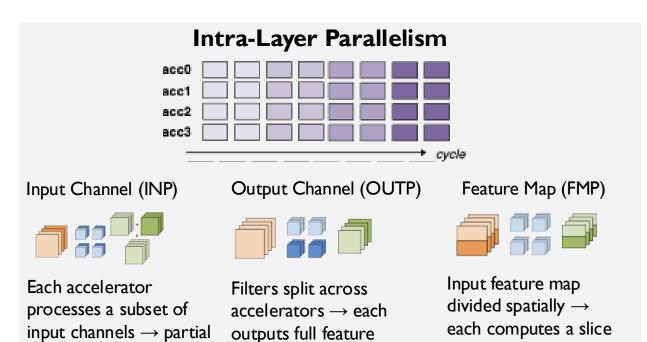
Network:



of outputs.

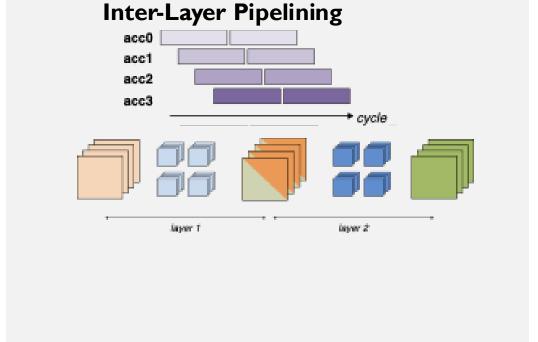




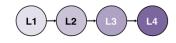


maps.

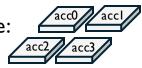
outputs combined.



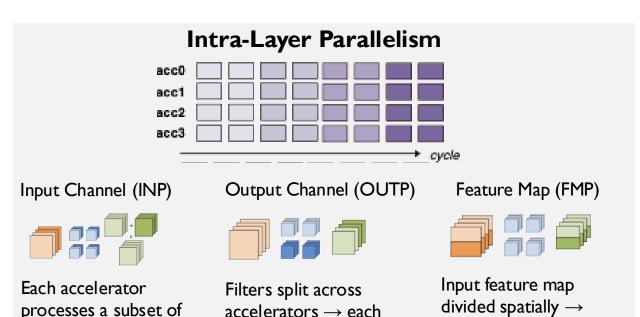
Network:



Architecture:



activations on-chip.



outputs full feature

maps.

input channels  $\rightarrow$  partial

outputs combined.

Inter-Layer Pipelining

acc0
acc1
acc2
acc3

Layer 1

Adjacent layer fused into segments, forwarding intermediate

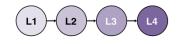
- Benefits: off-chip accesses reduction
- Cost: pipeline fill/flush overhead and larger on-chip buffers.

CASES'25

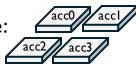
each computes a slice

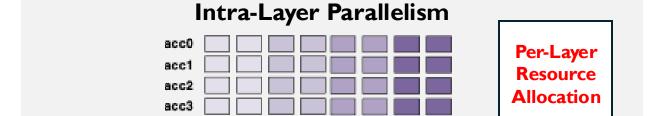
of outputs.

Network:



Architecture:





Input Channel (INP)



Each accelerator processes a subset of input channels → partial outputs combined.

Output Channel (OUTP)

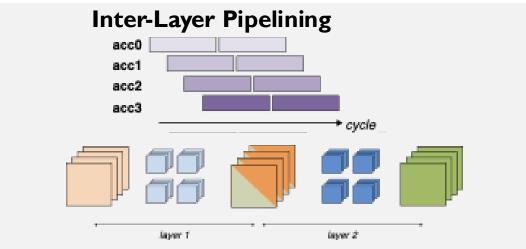


Filters split across accelerators → each outputs full feature maps.

Feature Map (FMP)



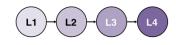
Input feature map divided spatially → each computes a slice of outputs.



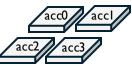
Adjacent layer fused into **segments**, forwarding intermediate activations on-chip.

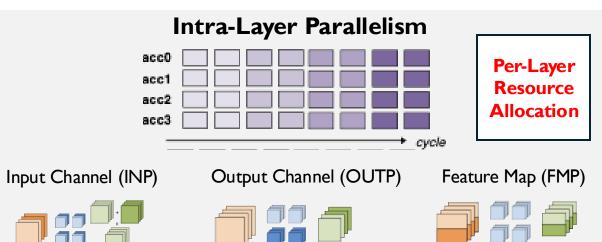
- Benefits: off-chip accesses reduction
- Cost: pipeline fill/flush overhead and larger on-chip buffers.

Network:



Architecture:



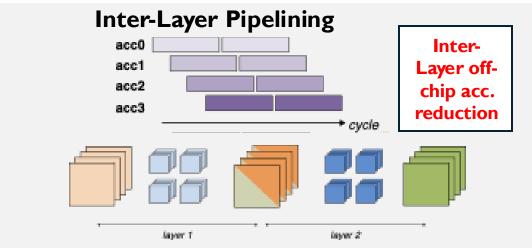


Each accelerator processes a subset of input channels  $\rightarrow$  partial outputs combined.

Filters split across  $accelerators \rightarrow each$ outputs full feature maps.



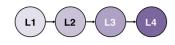
Input feature map divided spatially  $\rightarrow$ each computes a slice of outputs.



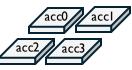
Adjacent layer fused into segments, forwarding intermediate activations on-chip.

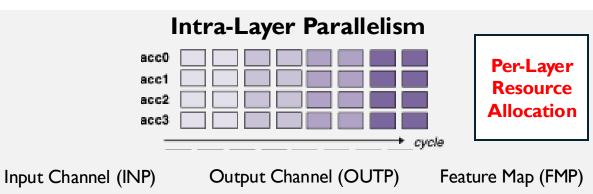
- **Benefits**: off-chip accesses reduction
- **Cost**: pipeline fill/flush overhead and larger on-chip buffers.

Network:



Architecture:





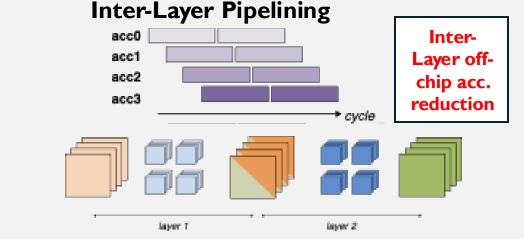
Each accelerator processes a subset of input channels → partial outputs combined.



Filters split across accelerators → each outputs full feature maps.



Input feature map divided spatially → each computes a slice of outputs.

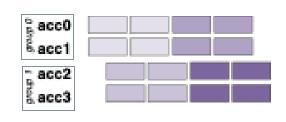


Adjacent layer fused into **segments**, forwarding intermediate activations on-chip.

- Benefits: off-chip accesses reduction
- Cost: pipeline fill/flush overhead and larger on-chip buffers.

→ Combination can strike the right balance but:

- Further increases the mapping space
- Requires advanced architectural support



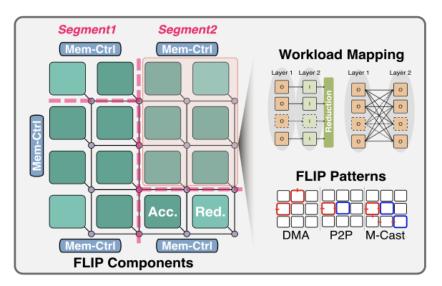
- Intra/Inter model heterogeneity (layers vary in intensity and size)
- Multi-model concurrency (contention for compute and off-chip bandwidth)
- Large mapping space and architectural challenges in tiled architectures (Intra/Inter-layer mapping)

- Intra/Inter model heterogeneity (layers vary in intensity and size)
- Multi-model **concurrency** (contention for compute and off-chip bandwidth)
- Large mapping space and architectural challenges in tiled architectures (Intra/Inter-layer mapping)

→ We propose **FLIP2M**: a holistic solution combining intra- and inter-layer optimization for multi-model AR/VR workloads on tiled architecture

- Intra/Inter model heterogeneity (layers vary in intensity and size)
- Multi-model concurrency (contention for compute and off-chip bandwidth)
- Large mapping space and architectural challenges in tiled architectures (Intra/Inter-layer mapping)

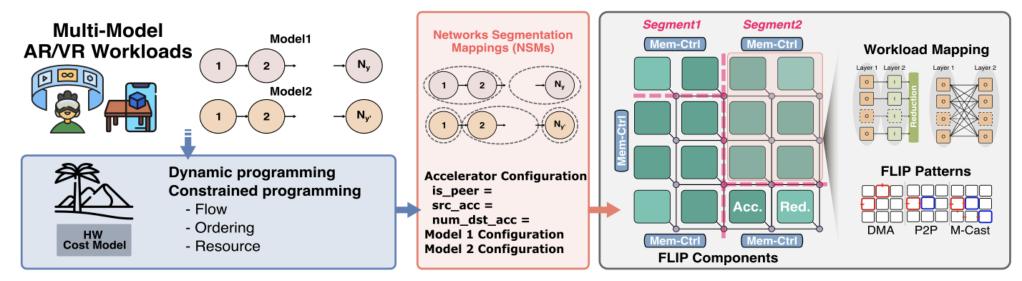
→ We propose **FLIP2M**: a holistic solution combining intra- and inter-layer optimization for multi-model AR/VR workloads on tiled architecture



O FLIP Acceleration Fabric

- Intra/Inter model **heterogeneity** (layers vary in intensity and size)
- Multi-model concurrency (contention for compute and off-chip bandwidth)
- Large mapping space and architectural challenges in tiled architectures (Intra/Inter-layer mapping)

→ We propose **FLIP2M**: a holistic solution combining intra- and inter-layer optimization for multi-model AR/VR workloads on tiled architecture

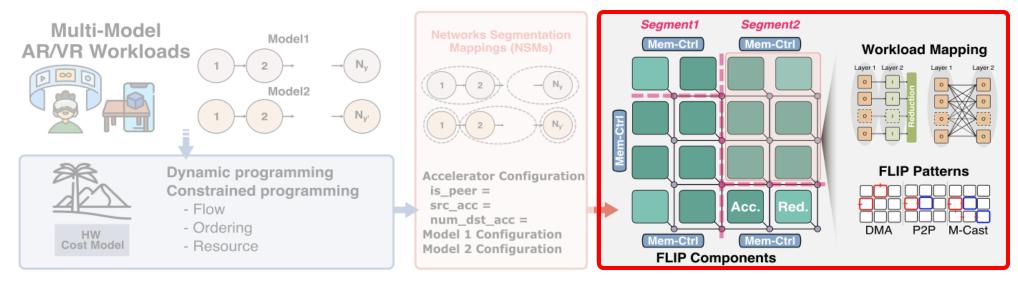


**@ OASIS Framework** 

• FLIP Acceleration Fabric

- Intra/Inter model heterogeneity (layers vary in intensity and size)
- Multi-model concurrency (contention for compute and off-chip bandwidth)
- Large mapping space and architectural challenges in tiled architectures (Intra/Inter-layer mapping)

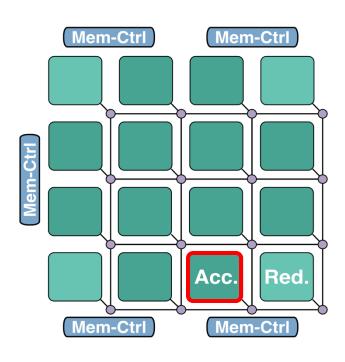
→ We propose **FLIP2M**: a holistic solution combining intra- and inter-layer optimization for multi-model AR/VR workloads on tiled architectures



**@ OASIS Framework** 

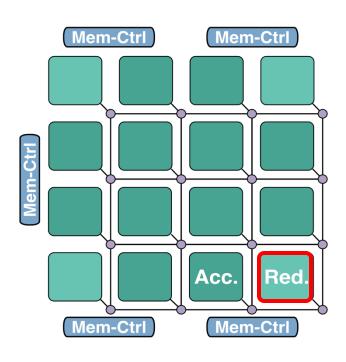
O FLIP Acceleration Fabric

#### • FLIP Architecture - Overview



- Accelerator Tiles:
  - Coarse-grained engines  $\rightarrow$  execute full layers or large portions.
  - Include private buffers, issue long load/store bursts.
  - Flexible design (vector lanes, systolic, etc.), to support DNN kernels

#### • FLIP Architecture - Overview



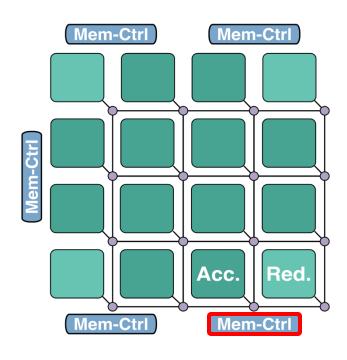
#### Accelerator Tiles:

- Coarse-grained engines  $\rightarrow$  execute full layers or large portions.
- Include private buffers, issue long load/store bursts.
- Flexible design (vector lanes, systolic, etc.), to support DNN kernels

#### Reduction Tiles:

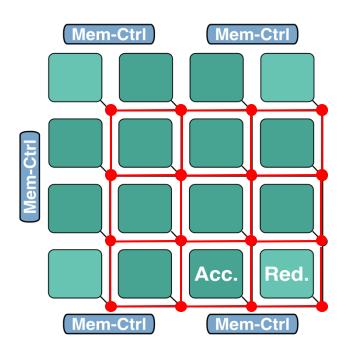
- Specialized for tensor addition (e.g., INP partial sums, residual merges).
- Hardware support avoids software bottlenecks.

#### FLIP Architecture - Overview



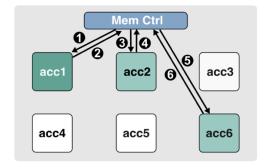
- Accelerator Tiles:
  - Coarse-grained engines  $\rightarrow$  execute full layers or large portions.
  - Include private buffers, issue long load/store bursts.
  - Flexible design (vector lanes, systolic, etc.), to support DNN kernels
- Reduction Tiles:
  - Specialized for tensor addition (e.g., INP partial sums, residual merges).
  - Hardware support avoids software bottlenecks.
- Memory Controllers:
  - Multiple DRAM controllers scale bandwidth and reduce contention.
     Partitioned address space → better isolation between segments.

#### • FLIP Architecture - Overview



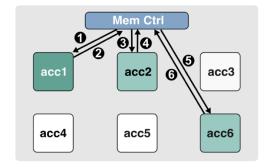
- Accelerator Tiles:
  - Coarse-grained engines  $\rightarrow$  execute full layers or large portions.
  - Include private buffers, issue long load/store bursts.
  - Flexible design (vector lanes, systolic, etc.), to support DNN kernels
- Reduction Tiles:
  - Specialized for tensor addition (e.g., INP partial sums, residual merges).
  - Hardware support avoids software bottlenecks.
- Memory Controllers:
  - Multiple DRAM controllers scale bandwidth and reduce contention.
     Partitioned address space → better isolation between segments.
- NoC-based Interconnect supporting a variety of communication patterns

Direct Memory Access (DMA)



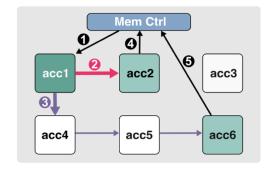
- Moves data between
   DRAM and accelerator
   buffers.
- Sustains high bandwidth for large transfers (weights, features).

#### Direct Memory Access (DMA)



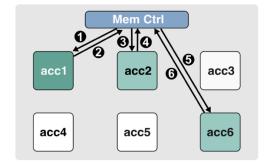
- Moves data between
   DRAM and accelerator
   buffers.
- Sustains high
   bandwidth for large
   transfers (weights,
   features).

#### Peer to Peer (P2P)



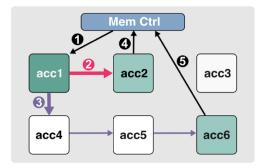
- Direct transfer between compute tiles. Key for interlayer pipelining → avoids round-trip to DRAM.
- Requires synchronization to prevent deadlock.

Direct Memory Access (DMA)



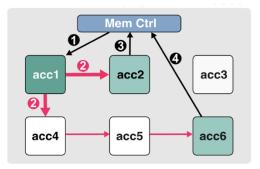
- Moves data between
   DRAM and accelerator
   buffers.
- Sustains high
   bandwidth for large
   transfers (weights,
   features).

Peer to Peer (P2P)



- Direct transfer between compute tiles. Key for interlayer pipelining → avoids round-trip to DRAM.
- Requires synchronization to prevent deadlock.

Multicast



- One tile sends
   output to multiple
   downstream tiles
   in parallel.
- Crucial for intra + inter-layer optimizations

#### • FLIP Architecture – Workload Mapping

• Single-layer segments: OUTP / INP

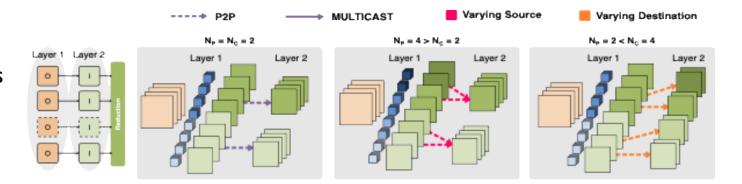
- Single-layer segments: OUTP / INP
- Multi-layer segments: Layers forward outputs directly to next layer's accelerators (pipelined)

- Single-layer segments: OUTP / INP
- Multi-layer segments: Layers forward outputs directly to next layer's accelerators (pipelined)
  - INP→OUTP bottlenecked by reduction tile → excluded.
  - Too many pipelined layers increase fill/flush overhead → capped at 3.

- Single-layer segments: OUTP / INP
- Multi-layer segments: Layers forward outputs directly to next layer's accelerators (pipelined)
  - INP→OUTP bottlenecked by reduction tile → excluded.
  - Too many pipelined layers increase fill/flush overhead → capped at 3.

#### $OUTP \rightarrow INP$

- Producers split output channels (OUTP).
- Consumers use INP → need only subsets of channels.
- Communication depends on producer/consumer ratio



- Single-layer segments: OUTP / INP
- Multi-layer segments: Layers forward outputs directly to next layer's accelerators (pipelined)
  - INP→OUTP bottlenecked by reduction tile → excluded.
  - Too many pipelined layers increase fill/flush overhead → capped at 3.

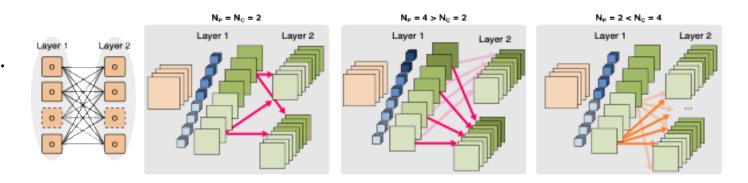
#### $OUTP \rightarrow INP$

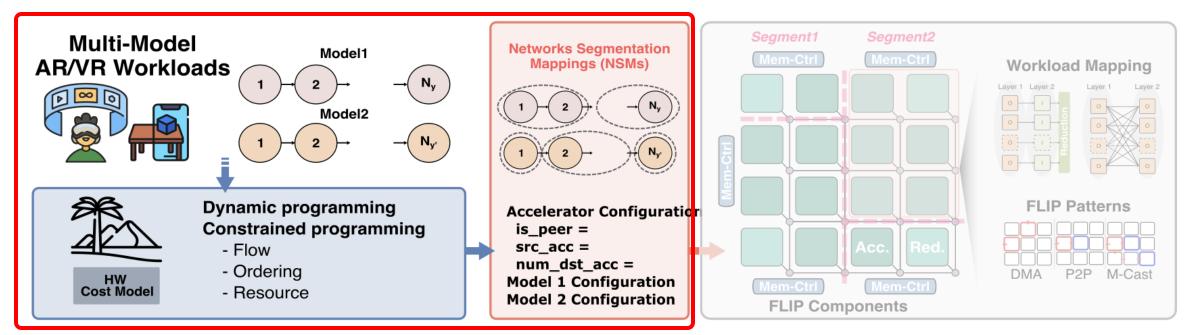
- Producers split output channels (OUTP).
- Consumers use INP → need only subsets of channels.
- Communication depends on producer/consumer ratio

## 

#### $OUTP \rightarrow OUTP$

- Both layers parallelized by filters (OUTP).
- Each consumer needs *all* input channels.
- Requires multicast: producers send to all consumers, each consumer gathers from all producers.



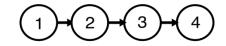


**OASIS** Framework

**O** FLIP Acceleration Fabric

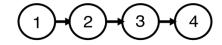
• **Goal:** Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.

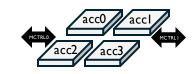
 Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.



- Inputs:
  - DNN topologies

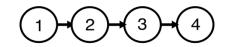
- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype

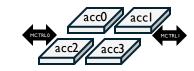




- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:





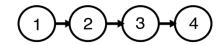
• **Goal:** Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.

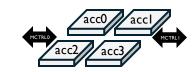


- DNN topologies
- FLIP prototype
- Segment:

Execution mode from layer i:

- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)

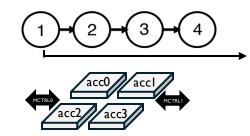




- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

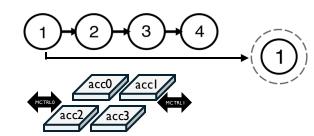
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

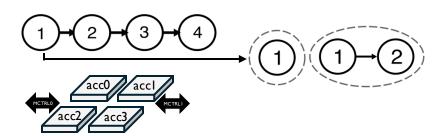
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- **Goal:** Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

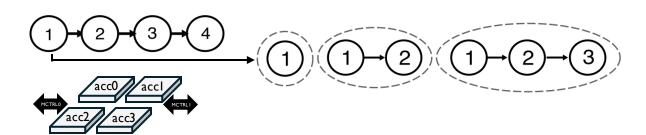
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

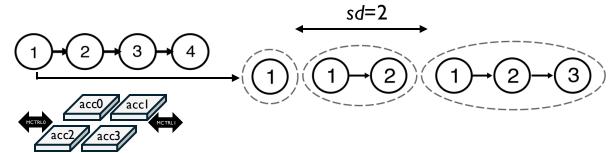
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

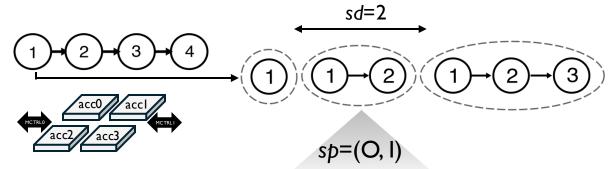
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

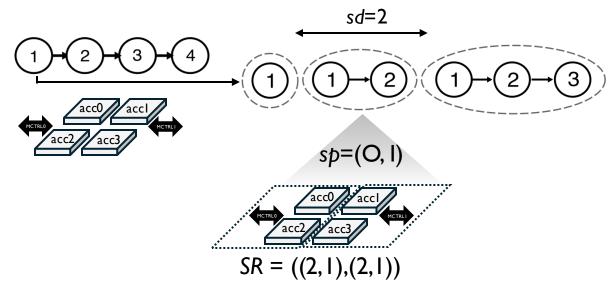
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

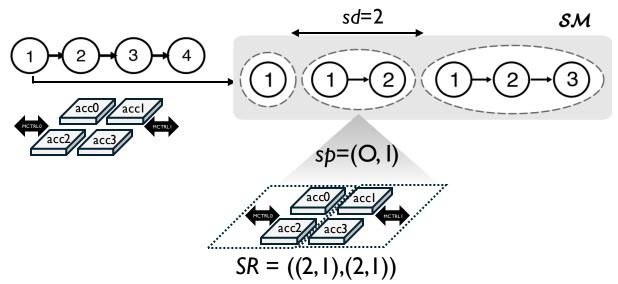
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

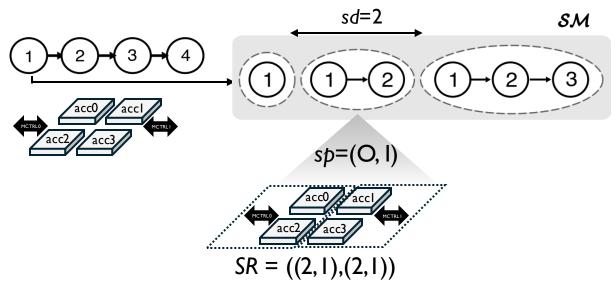
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)
- Mapping Space:
  - Segment mapping space (SM) = all possible segments from a layer



- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)
- Mapping Space:
  - Segment mapping space (SM) = all possible segments from a layer
  - Network mapping = sequence of segment mappings covering all layers.

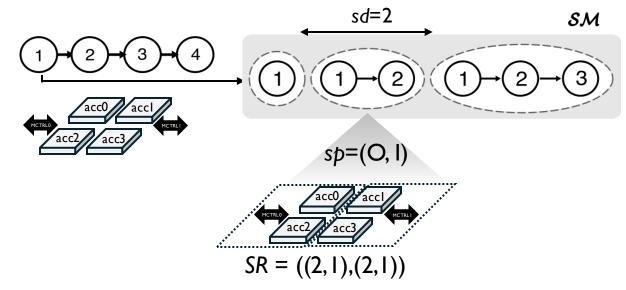


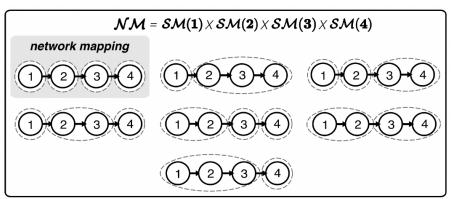


- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

Execution mode from layer i:

- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)
- Mapping Space:
  - Segment mapping space (SM) = all possible segments from a layer
  - Network mapping = sequence of segment mappings covering all layers.
  - Network mapping space (NM) = all possible network mappings

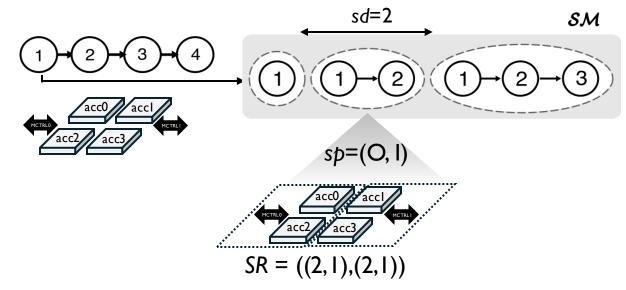


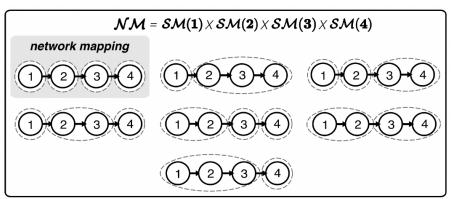


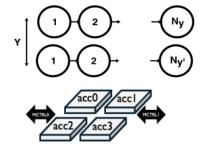
- Goal: Optimize segmentation, resource allocation, and scheduling for multi-model AR/VR on FLIP.
- Inputs:
  - DNN topologies
  - FLIP prototype
- Segment:

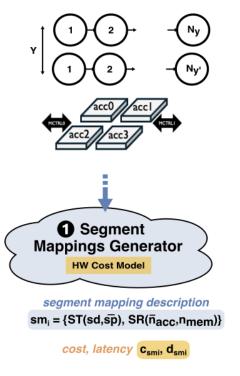
Execution mode from layer i:

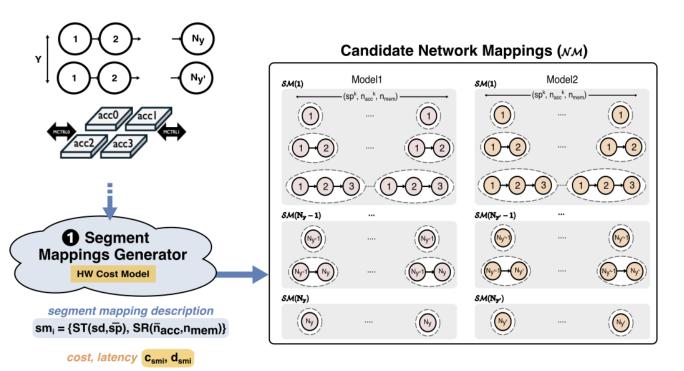
- Segment depth (sd)
- Segment parallelism (sp)
- Segment resource (SR)
- Mapping Space:
  - Segment mapping space (SM) = all possible segments from a layer
  - Network mapping = sequence of segment mappings covering all layers.
  - Network mapping space (NM) = all possible network mappings

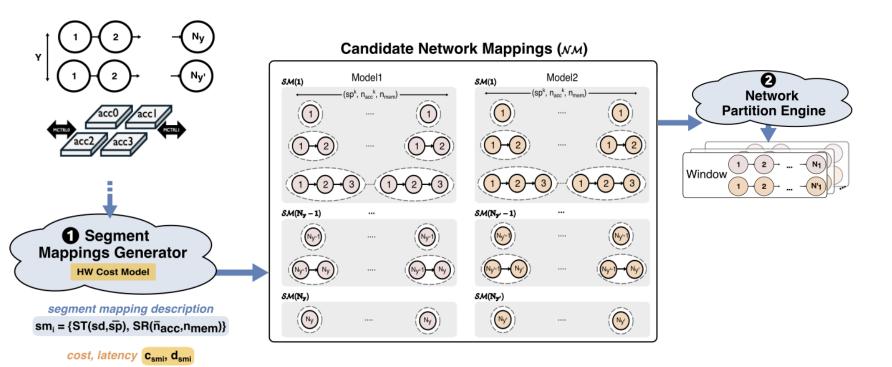


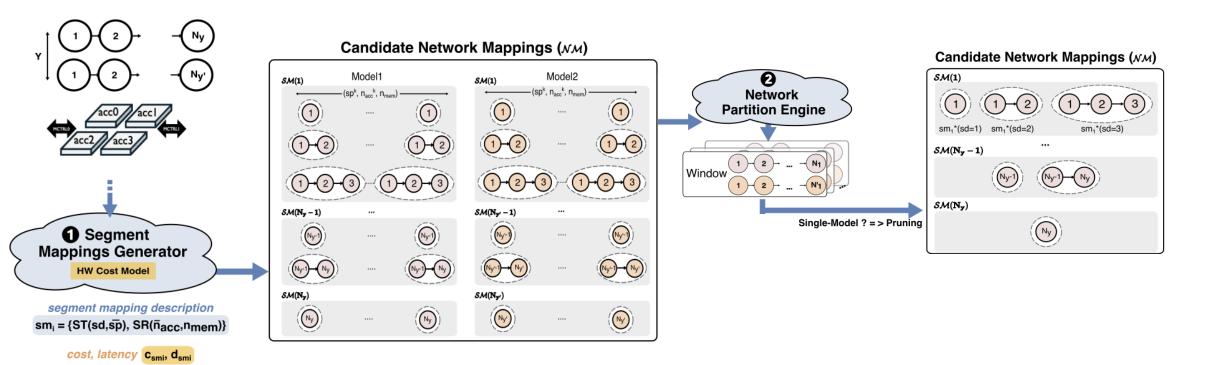


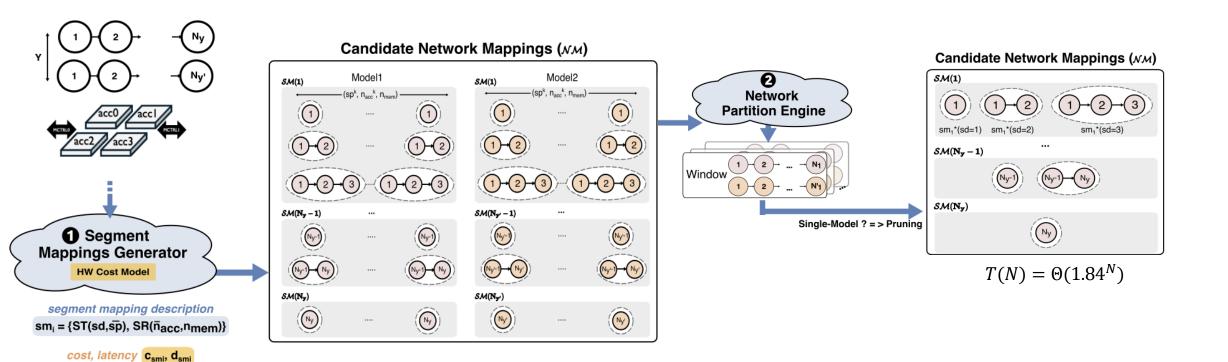


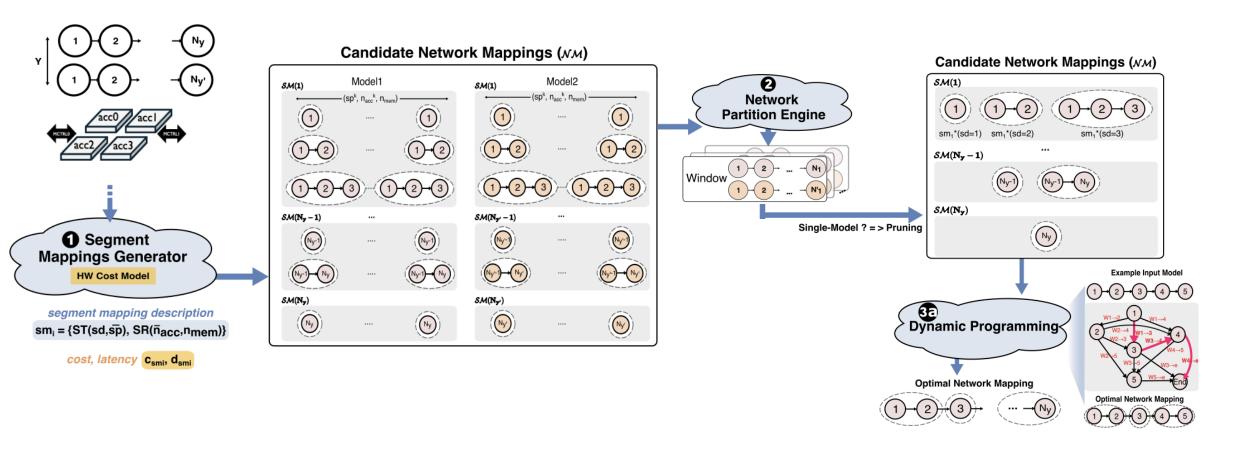


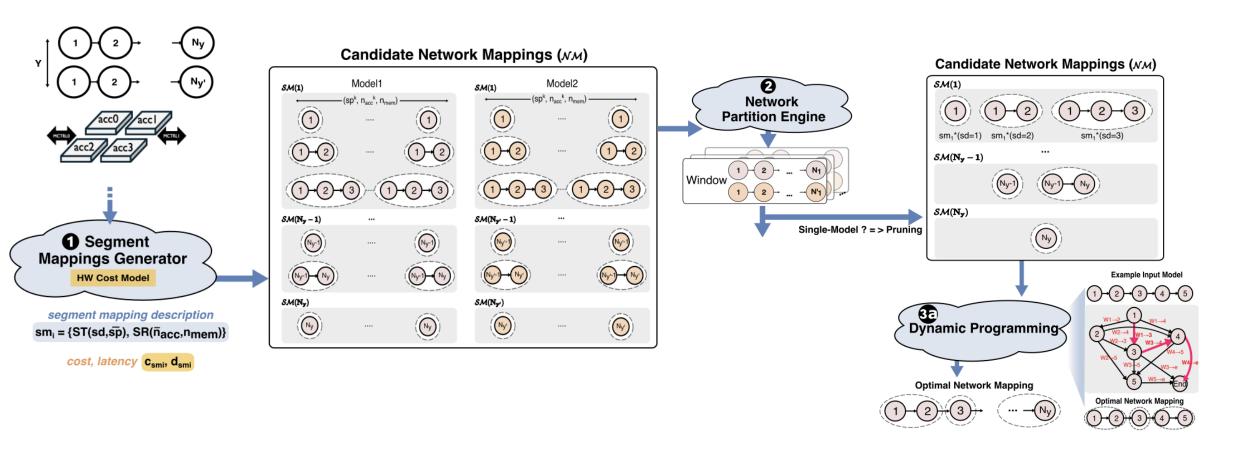


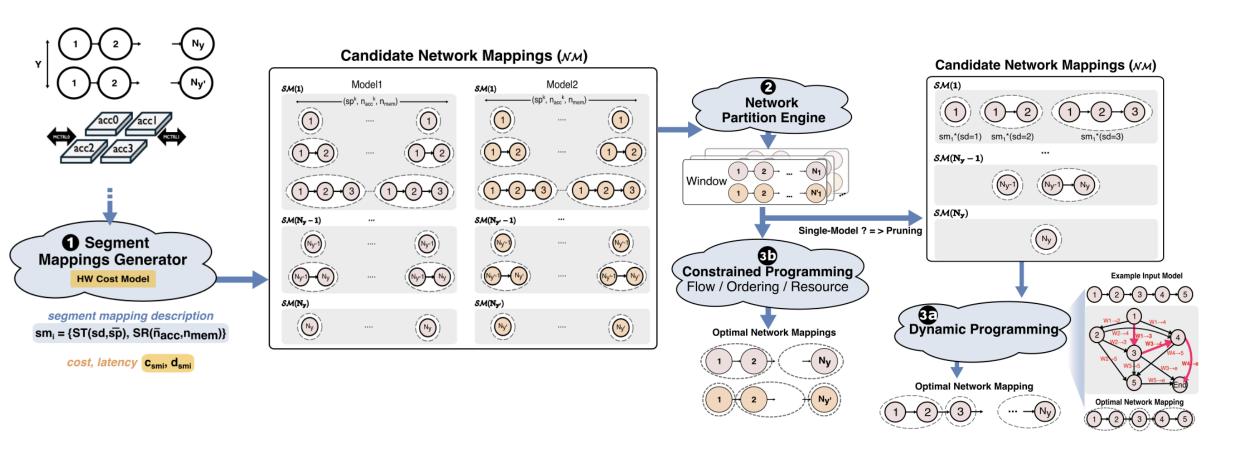












## Evaluation - FLIP Prototype

**Platform:** Built on open-source ESP:

esp.cs.columbia.edu



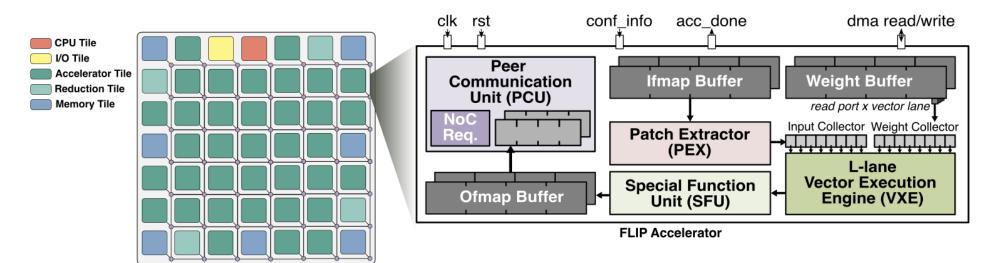
Implemented on Xilinx XCVU19P FPGA. **Tile architecture (7×7 grid = 49 tiles):**36 acc tiles / 4 red. Tiles / 7 memory tiles / 1 processor tile (CVA6 host) / I auxiliary tile (peripherals).

#### **Custom features:**

- Multi-level data reuse (weights & inputs) for acc tiles.
- Dynamic communication patterns: DMA, P2P, multicast.
- Modified ESP NoC for flexible producer/consumer access
   + multicast support.

**Resource utilization:** ~48% LUTs, ~21% registers, ~60% DSPs, ~74% BRAMs.

**Power:** Dominated by accelerators, reduction tiles, and NoC



## Evaluation – Experimental Setup

**Benchmarks:** Multi-model AR/VR workloads (from XRbench)  $\rightarrow$  3 scenarios combining different models & batch sizes:

AR/VR1			AR/VR2			AR/VR3		
Task	Model	Batch size	Task	Model	Batch size	Task	Model	Batch size
Gaze Estimation	ResNet-18	2	Hand Tracking	ResNet-34	4	Eye Segmentation	Vgg16	2
Keyword Detection	SqueezeNet	1	Plane Detection	Unet	2	Depth Refinement	MobileNet-v2	2
Object Detection	MobileNet-v2	2	Speech Recognition	MobileBert	1	Action Segmentation	ResNet-18	4
			Depth Estimation	ResNet-50	2	Speech Recognition	MobileBert	1

#### **Evaluation:**

- Single- vs. multi-model deployments
- Metrics include latency, energy, and EDP (from Vivado power + performance counters).

### Evaluation – Multi-Model Performance

#### **Configurations:**

- Baseline  $\rightarrow$  fixed per-model resources, no pipelining.
- FlexIntra I/2  $\rightarrow$  flexible per-model accelerator & memory BW allocation.
- Full  $\rightarrow$  + inter-layer pipelining.

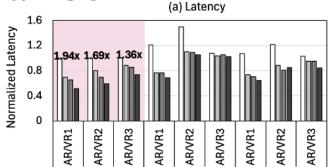
### Evaluation – Multi-Model Performance

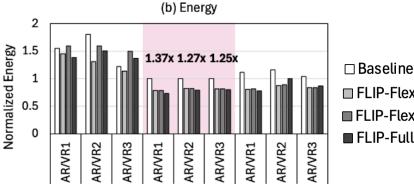
#### **Configurations:**

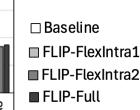
- Baseline  $\rightarrow$  fixed per-model resources, no pipelining.
- FlexIntra I/2 → flexible per-model accelerator & memory BW allocation.
- Full  $\rightarrow$  + inter-layer pipelining.

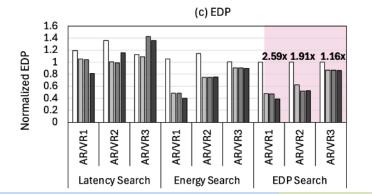
#### **Key Results:**

- **Latency:** Up to **1.9× speedup** (higher than single-model gains).
- **Energy:** FlexIntra variants give most savings ( $\leq 1.4\times$ ).
- **EDP:** Consistently improved by reduced idle cycles & better allocation









### Evaluation – Multi-Model Performance

#### **Configurations:**

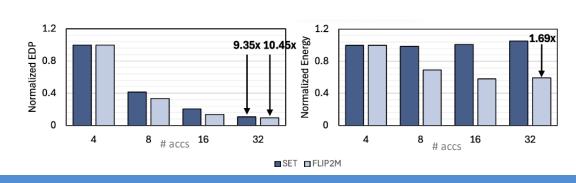
- Baseline  $\rightarrow$  fixed per-model resources, no pipelining.
- FlexIntra I/2 → flexible per-model accelerator & memory BW allocation.
- Full  $\rightarrow$  + inter-layer pipelining.

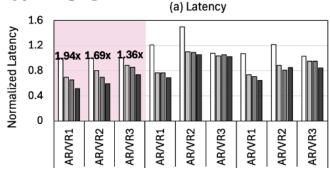
#### **Key Results:**

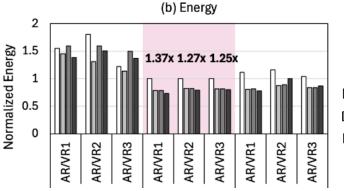
- Latency: Up to 1.9× speedup (higher than single-model gains).
- **Energy:** FlexIntra variants give most savings (≤1.4×).
- **EDP:** Consistently improved by reduced idle cycles & better allocation

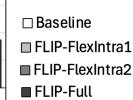
#### Comparison to SET [1] (state-of-the-art temporal multi-model):

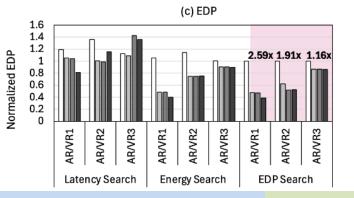
- FLIP2M achieves 10.5× EDP reduction vs. 9.4× at 32 tiles.
- Better energy scaling (up to 1.7× lower).
- Advantage: concurrent pipelines via spatial partitioning.











• **FLIP2M**: integrates intra-layer parallelism + inter-layer pipelining → efficient multi-model AR/VR on tiled accelerators.

- **FLIP2M**: integrates intra-layer parallelism + inter-layer pipelining → efficient multi-model AR/VR on tiled accelerators.
- **FLIP accelerator fabric**: flexible architecture enabling scalable parallelism & diverse on-chip communication.

- FLIP2M: integrates intra-layer parallelism + inter-layer pipelining → efficient multi-model AR/VR
  on tiled accelerators.
- **FLIP accelerator fabric**: flexible architecture enabling scalable parallelism & diverse on-chip communication.
- OASIS framework: navigates vast mapping space for single- & multi-model workloads.

- FLIP2M: integrates intra-layer parallelism + inter-layer pipelining → efficient multi-model AR/VR on tiled accelerators.
- **FLIP accelerator fabric**: flexible architecture enabling scalable parallelism & diverse on-chip communication.
- OASIS framework: navigates vast mapping space for single- & multi-model workloads.
- **Prototype**: 49-tile FPGA SoC on ESP platform.

- FLIP2M: integrates intra-layer parallelism + inter-layer pipelining → efficient multi-model AR/VR
  on tiled accelerators.
- **FLIP accelerator fabric**: flexible architecture enabling scalable parallelism & diverse on-chip communication.
- OASIS framework: navigates vast mapping space for single- & multi-model workloads.
- Prototype: 49-tile FPGA SoC on ESP platform.
- Results: up to 1.9× latency speedup, 1.4× energy reduction, 2.6× EDP improvement over baseline.

#### In Memory of Davide Giri (1990-2021)



#### Thank you!

# FLIP2M: Flexible Intra-layer Parallelism and Inter-layer Pipelining for Multi-model AR/VR Workloads

#### Gabriele Tombesi

Je Yang Joseph Zuckerman Davide Giri William Baisi Luca Carloni

