



# KAPLA: Scalable NN Accelerator Dataflow Design Space Structuring and Fast Exploring

## Zhiyao Li and Mingyu Gao

Tsinghua University Shanghai Artificial Intelligence Lab Shanghai Qi Zhi Institute



ASP-DAC 2025

# **Background & Motivation**

### Domain-specific accelerators (DSAs) for neural networks (NNs)

- Abundant parallelism of PEs
- Localized memory accesses
- Reduced control overheads



## **Background & Motivation**

### DSAs' hardware architectures and dataflow schedulers



# A Large Design Space of Dataflow Scheduling

### Complex hardware structures

- Diverse PE array structures: 1D, 2D, adder tree, ...
- Deep hierarchical buffer storage: register file, global buffer, main memory, ...
- Numerous constraints: PE dataflow, buffer capacity, ...

### • Various algorithm designs

- Complex model topologies
- 3D/4D tensor dimensions
- $_{\odot}$  Abundant layer types

### □ Flexible use scenarios

- Offline: traditional compilers
- o Online: MLaaS, NAS, ...



scheduling can no longer be ignored!

# **Hierarchical NN Dataflow Taxonomy**

#### **Segment slicing**

- Model DAG is sliced into multiple segments where each contains one or more layers
- Each segment is scheduled on hardware one after one

#### Node parallelism

- Each layer is parallelized across multiple hardware processing nodes
- Data parallelism, tensor parallelism, ...

#### Loop blocking

 For each (sub)layer in one hardware node, optimize off-chip data accesses using nested loop transformation techniques, including loop blocking, reordering, unrolling, etc.



#### Layer pipelining

- Layers in a segment are processed in a pipelined style
- Intermediate data are forwarded on-chip without spilled to off-chip memory

### **PE mapping**

• For each (sub)layer in one hardware node, optimize on-chip dataflow on the PE array, such as systolic, row-stationary, etc.

# **Prior Dataflow Techniques**



# **Prior Scheduling Methods**

### Prior schedulers are either too slow or inaccurate

Method	Representative	Problems
Exhaustive search	TANGRAM [ASPLOS'19]	Long search time
Random sampling	Timeloop [ISPASS'19]	Low scheduling quality
ML-guided search	AutoTVM [NeurIPS'18]	Long search time; no quality guarantees
Polyhedral models	AKG [PLDI'21]	Restricted search space

Design goal:

A generic, optimized, and fast scheduler is needed!

## **Our Key Ideas**



## **Tensor-Centric Dataflow Directives**

tensor(dim=size, ...[, shr])

 $\rightarrow$  (Sub)tensors in each buffer

stack(dim+=shift, ..., repl)

 $\rightarrow$  Spatial parallelism, in multiple buffers at one level

update(dim+=step, ...)

 $\rightarrow$  Temporal accesses, across multiple buffer levels

### Advantages

- Tensor as first-class citizen
  - Fast validity checking of tensor sizes in buffers
- Bottom-up dataflow construction
  - Fast memory hierarchy traffic calculatation
- Inter-operator pipeline native
  - Easy to express layer pipelining optimizations

Directive examples for CONV and DWCONV layers, with **row-stationary** PE mapping, output + batch **hybrid node parallelization**, and **layer pipelining**.

```
1 CONV :
  REGE :
    tensor{0}(N=1, C=2, Xi=5, Yi=1)
    tensor{w1}(C=2, K=3, R=5, S=1)
   tensor{1}(N=1, K=3, Xo=1, Yo=1)
   stack(Yi+=1, Yo+=1, 8) % PE columns
   stack(S+=1, Yi+=1, 5) % PE rows
   update(Xi+=1, Xo+=1)
                            % 1D conv
    update(Yi+=8, Yo+=8)
                            % folding
    update(N+=1)
    update(C+=2)
   update(K+=3)
13
  GBUF :
    tensor{0}(N=4, C=4, Xi=19, Yi=19, shr=4)
   tensor{w1}(C=4, K=6, R=5, S=5)
   tensor{1}(N=4, K=6, Xo=15, Yo=15)
   stack(K+=6, 4)
                   % output node parallel
17
18
   stack(N+=4, 16) % batch node parallel
   update(C+=4)
19
   update(K+=24)
20
   update(N+=64)
21
22 DWCONV :
  REGE :
23
   % . . .
24
  GBUF :
25
   % DWCONV input is the same as CONV output
   tensor{1}(N=4, C=4, Xi=9, Yi=15)
27
   tensor{w2}(C=4, R=3, S=3)
28
   tensor{2}(N=4, C=4, Xo=4, Yo=7)
29
   stack(C+=4, 6)
30
                    % channel node parallel
   stack(N+=4, 16) % batch node parallel
31
   stack(Xo+=4, 2) % output width node parallel
33
    update(Yo+=7)
34
    update(C+=24)
35
   update(N+=64)
```

## **KAPLA Dataflow Solver**

	Validity check	<b>Efficiency estimation</b>
Inter-layer	Conservative pruning	
Intra-layer		

Estimate the minimum required buffer capacity of a layer, with optimistic intra-layer schemes

	<b>Total Schemes</b>	After Pruning	% Pruned	
AlexNet	700	2	99.7%	
MobileNet	331	42	87.3%	
VGGNet	7	1	85.7%	
GoogLeNet	560	1	99.8%	
ResNet	63	4	93.6%	
MLP	784	18	97.7%	
LSTM	70	3	95.7%	



## **KAPLA Dataflow Solver**





	Validity check	Efficiency estimation	
Inter-layer	Conservative pruning	Prioritization based on estimated cost	
Intra-layer	Bottom-up tensor construction		

Bottom-up tensor construction ensures schemes are always valid





	Validity check	<b>Efficiency estimation</b>	
Inter-layer	Conservative pruning	Prioritization based on estimated cost	
Intra-layer	Bottom-up tensor construction	Cost descending search	

#### Iteratively apply stacking and caching with the best loop orders

Algorithm 1 Greedy cost descending algorithm

- 1: **Input:** NN layer spec LayerSpec, hardware spec HWSpec, and constraints Constrs
- 2: Output: Optimized dataflow scheme OptDf
- 3: MemoryHier = [RegFile, GBUF, DRAM]
- 4: Df, LoopCnt = initUnitMapping(LayerSpec, HWSpec)
- 5: for each m in MemoryHier do
- 6: **for** order **in** genLoopOrder(m) **do**
- 7: Df, LoopCnt = stack(Df, m, LoopCnt, order)
- 8: Df, LoopCnt = cache(Df, m, LoopCnt, order)
- 9: **if checkValid**(Constrs, LoopCnt) is False **then**
- 10: continue
- 11: **if** evalCost(Df) > evalCost(OptDf) then
- 12: OptDf = Df



# **Evaluation Methodology**

### Evaluated workloads

- Inference & training on various NNs
- $_{\odot}$  Batch sizes 64 & 1

### Evaluated architectures

- 16 x 16 multi-array cloud accelerator
- $_{\circ}$  single-array edge accelerator

### Baselines

- o (B) nn-dataflow [ASPLOS'19]: exhaustive search
- (S) KAPLA directives: exhaustive search
- o (R) Timeloop [ISPASS'19]: random search
- (M) AutoTVM [NeurlPS'18]: XGBoost

# **Evaluation Results: Scheduling Speed**

### 518x search speedup

	В	S	R	М	К
AlexNet	8.7 h	12.7 h	4.6 min	3.2 h	32 s
MobileNet	17.6 h	33.8 h	16.6 min	10.1 h	53 s
VGGNet	6.8 h	12.2 h	5.5 min	3.6 h	35 s
GoogLeNet	64.2 h	111.0 h	31.4 min	37.1 h	262 s
ResNet	78.6 h	114.8 h	42.4 min	28.6 h	118 s
MLP	0.2 h	0.7 h	1.0 min	0.2 h	13 s
LSTM	0.2 h	0.9 h	7.4 min	0.1 h	6 s

Comparison of scheduling time for NN **training** on **multi-array** accelerators. Measured on an Intel Xeon Gold 5120 processor with 8 parallel processes.

# **Evaluation Results: Scheduling Quality**

### Only 2.2% quality (energy) loss on average



Comparison of energy for training (left) and inference (right) on multi-array Eyeriss-like accelerators with batch 64.

## Summary

- A hierarchical taxonomy that highlights the tight coupling across different levels of the dataflow space as the major difficulty for fast design exploration
- A set of tensor-centric directives that can accurately express various inter-layer and intra-layer schemes, and allow for quickly determining their validity and efficiency
- A generic, optimized, and fast dataflow solver, KAPLA, which makes use of the proposed directives to explore the design space, and achieves similar scheduling quality with much faster scheduling speed

