



# CoSA: Scheduling by Constrained Optimization for Spatial Accelerators

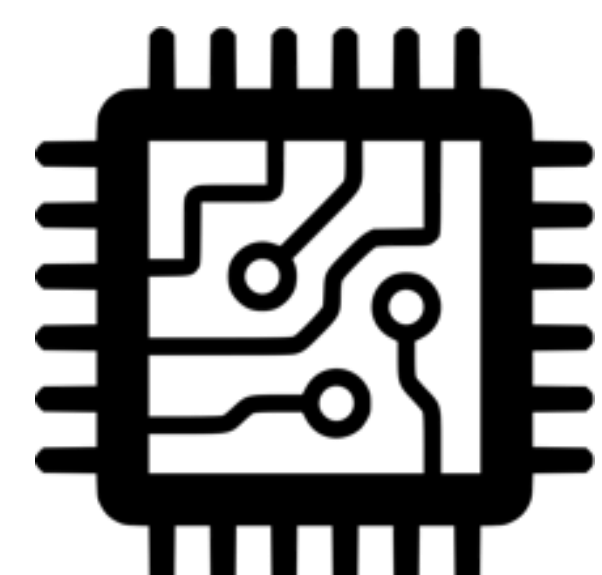
Qijing Huang, Minwoo Kang, Grace Dinh, Thomas Norell, Aravind Kalaiah,  
James Demmel, John Wawrzynek, Yakun Sophia Shao

Email: [jennyhuang@nvidia.com](mailto:jennyhuang@nvidia.com) Git repo: <https://github.com/ucb-bar/cosa>

## Scheduling is a big challenge



Intractable  
scheduling  
space



Exponentially growing  
algorithm complexity

- $\sim 10^{13}$  possible mappings for a typical ResNet layer on a 3-level architecture

Rapidly increasing  
hardware capacity

## Accelerator-oriented scheduling

Key DNN accelerator properties to leverage:

Workload  
Regularity

Known HW  
Constraints

Explicit Data  
Management

## CoSA: a one-shot approach

**Brute-force**

Timeloop  
dMazeRunner  
Triton Interstellar  
Marvel

- Costly
- Sample invalid space
- Hard to generalize

**Feedback-based**

AutoTVM Halide  
FlexFlow Gamma  
MindMapping

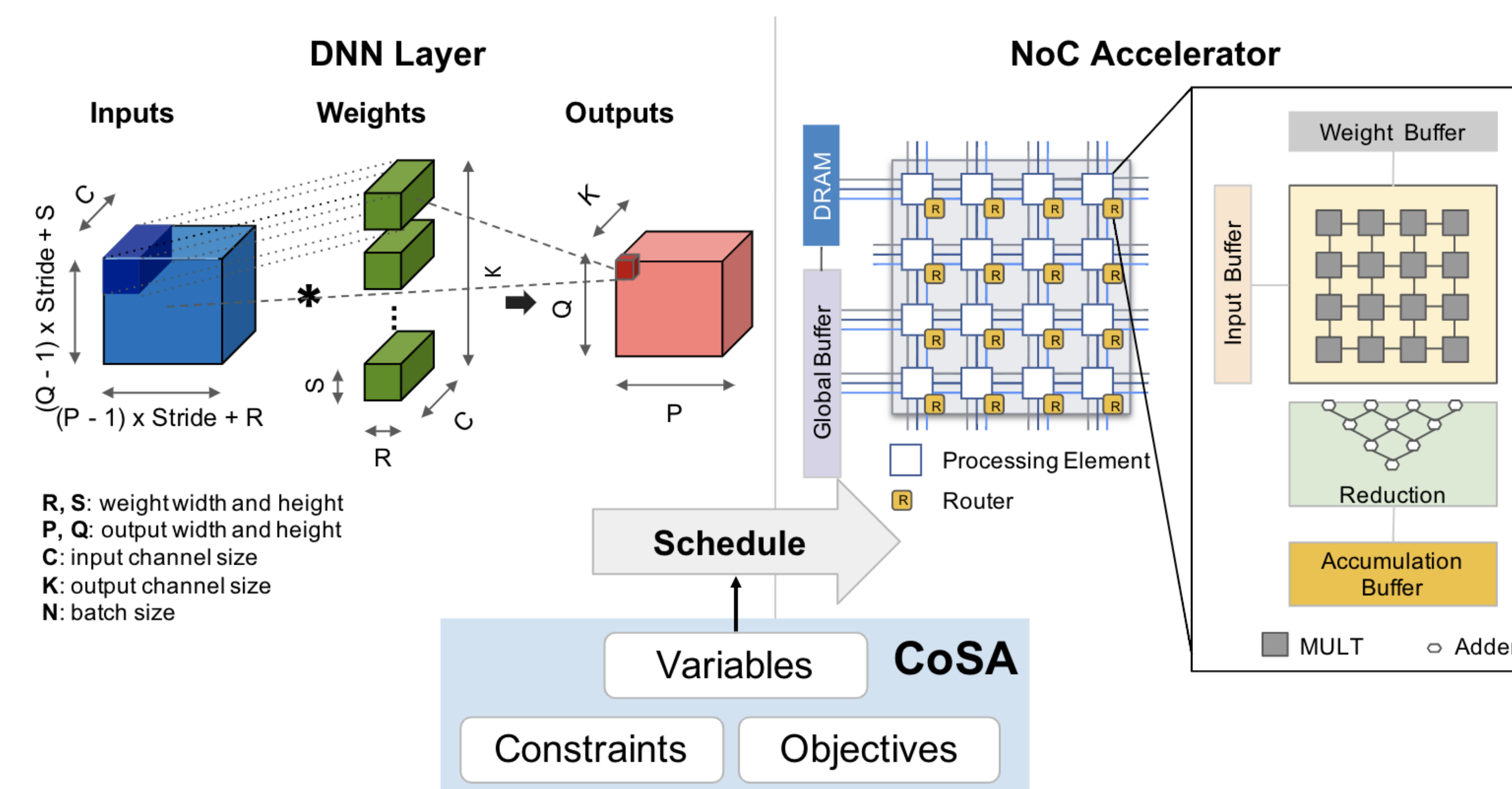
**Constrained  
Optimization**

Polly+Pluto TC  
Tiramisu

- Unable to determine  
tiling factor sizes

**CoSA - addresses key  
scheduling decisions**

## DNN scheduling formulation with CoSA



## Three scheduling decisions

**DRAM level**

for q2 = [0 : 2) :

**Global Buffer level**

for q1 = [0 : 7) :

for n0 = [0 : 3) :

*spatial* for r0 = [0 : 3) :

*spatial* for k1 = [0 : 2) :

**Input Buffer level**

for c1 = [0 : 2) :

for p1 = [0 : 2) :

**Weight Buffer level**

for p0 = [0 : 2) :

*spatial* for k0 = [0 : 2) :

**1. Tiling Factors**

How do we break up our tensors?

**2. Spatial / Temporal**

Which loops do we parallelize?

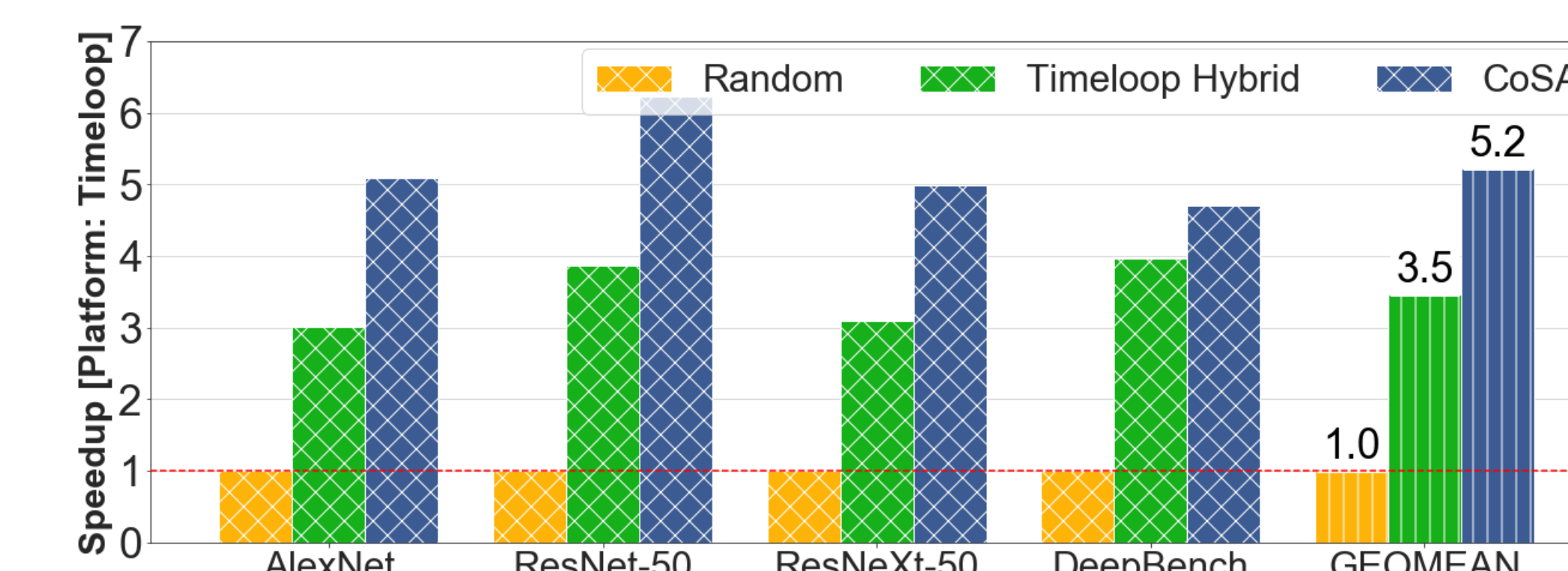
**3. Loop Permutation**

How do we order the loops?  
What dataflow do we use"

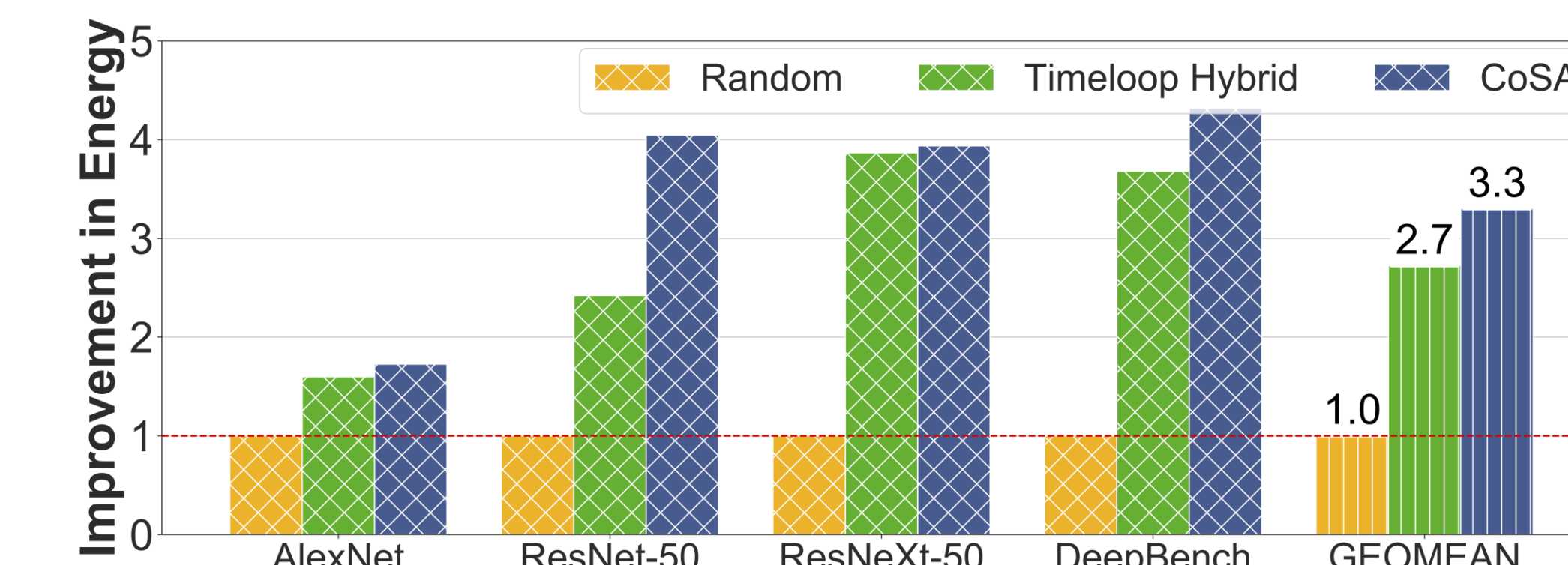
## CoSA Evaluation

- Baselines:
  - Random (best out of 5 valid schedules)
  - Timeloop Hybrid (best out of 16K valid schedules)
- Platforms:
  - Timeloop Simulator

### 1.5x latency speedup



### 1.2x better energy efficiency



### 90x faster time-to-solution

	CoSA	Random	Timeloop Hybrid
Runtime / Layer	4.2s	4.6s (1.1x)	379.9s (90.5x)
Samples / Layer	1	20K	67M
Evaluations/ Layer	1	5	16K

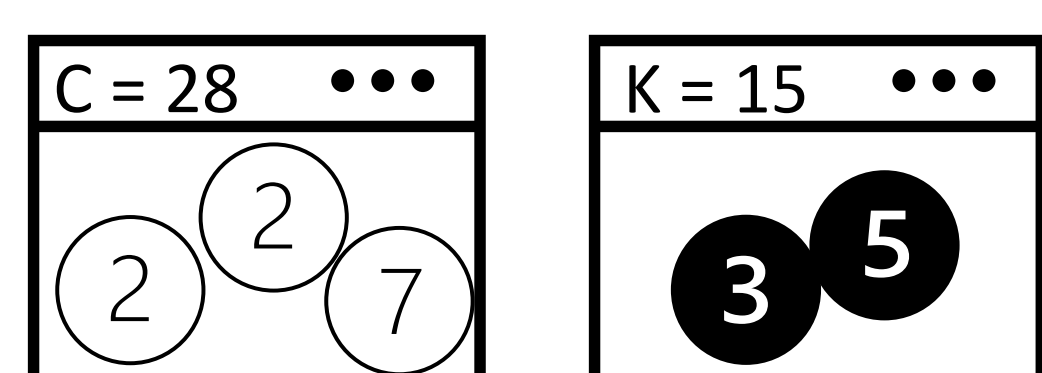
3

## Key idea: prime factor allocation

**Matrix-vector mult:**

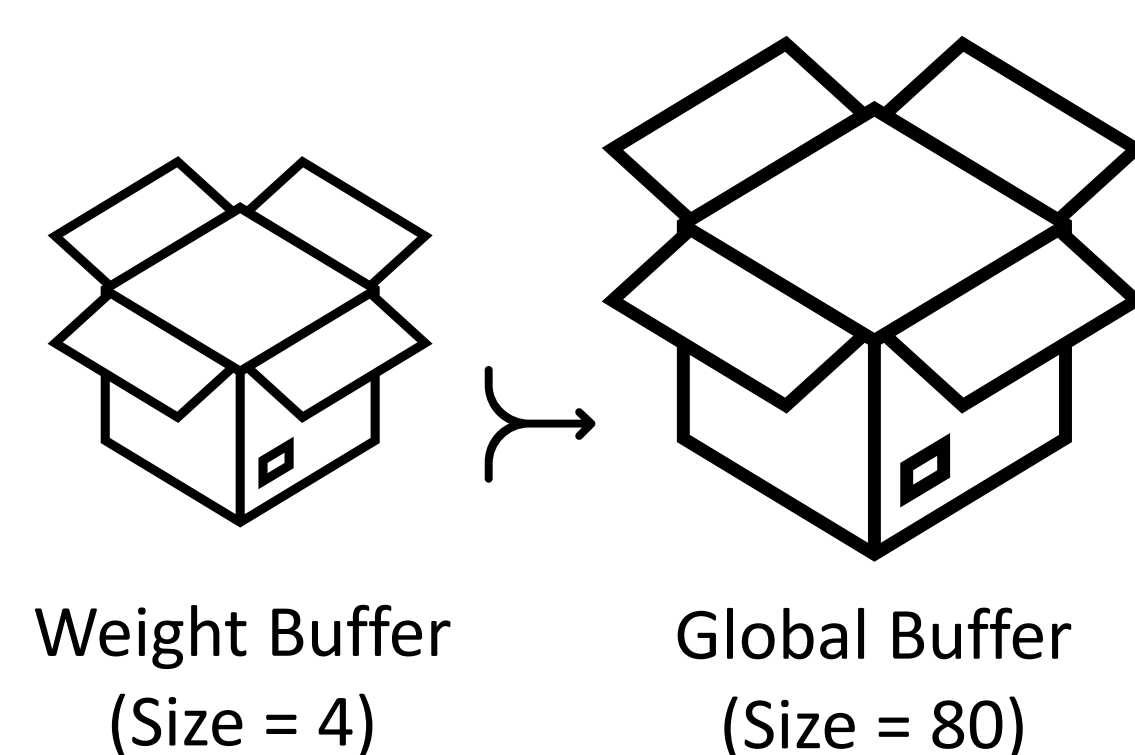
for c in [0:C) // C = 28  
for k in [0:K) // K = 15  
OA[k] += IA[c] x W[c,k]

**Prime factor items:**



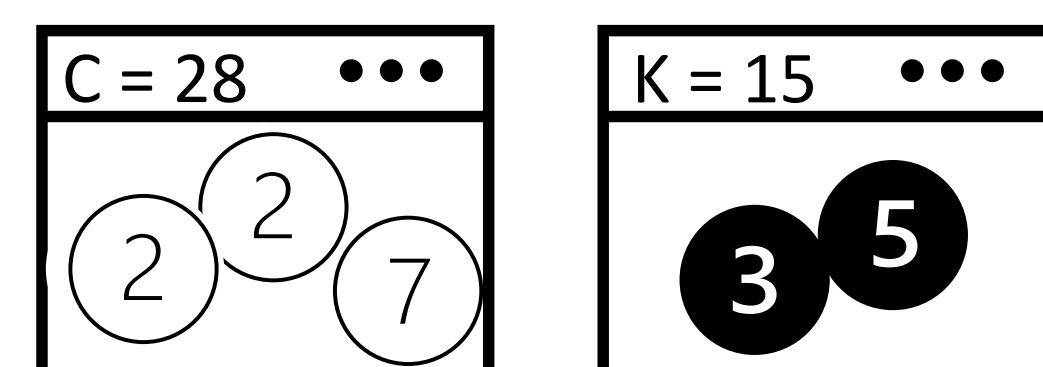
**Local buffers:**

- Weight buffer
- Global buffer



### CoSA Variable X

**Prime factor items :**



#### 1.Tiling Factors

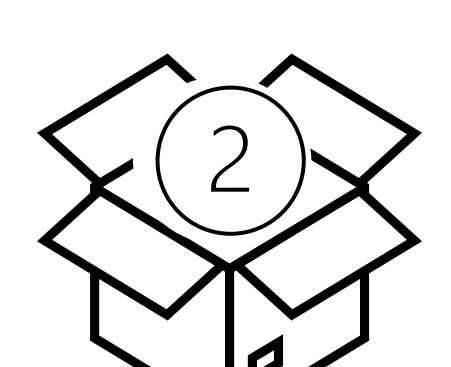
	C=28			K=15	
Prime Factors	2	2	7	3	5
WeightBuf	✓				
GlobalBuf		✓		✓	✓
DRAM			✓		

#### 2.Spatial/Temporal Factors

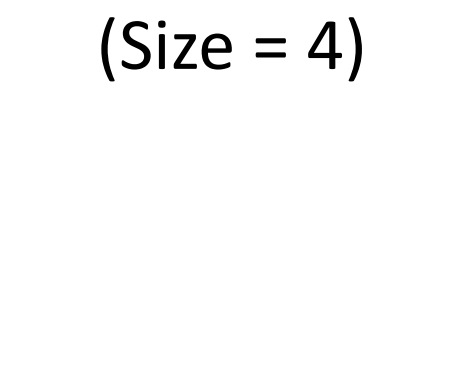
	C=28			K=15	
Prime Factors	2	2	7	3	5
Spatial				✓	
Temporal	✓				✓

**Local buffers:**

Utilized: 2



Weight Buffer (Size = 4)

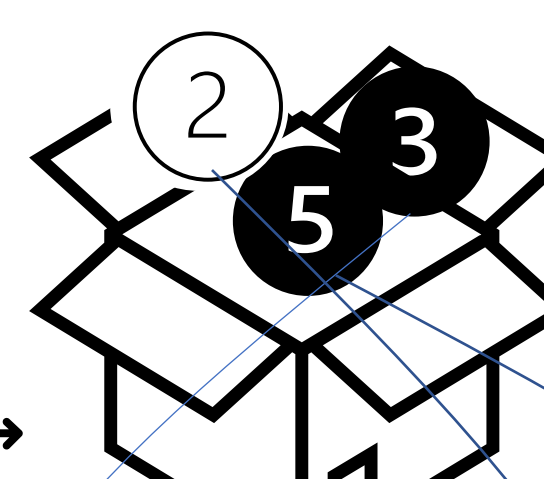


Global Buffer (Size = 80)

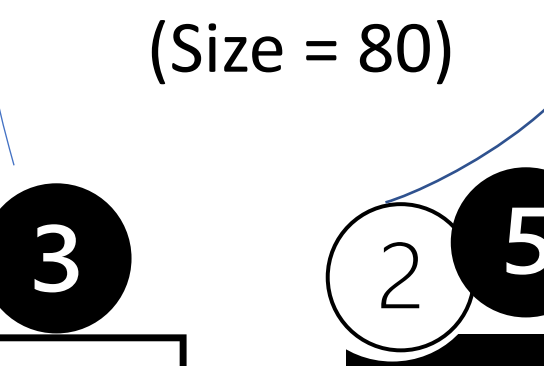
Spatial Part (Limit = 4)

Temporal Part

Utilized: (2x3x5)x(2)=60



Weight Buffer (Size = 4)



Global Buffer (Size = 80)

Spatial Part (Limit = 4)

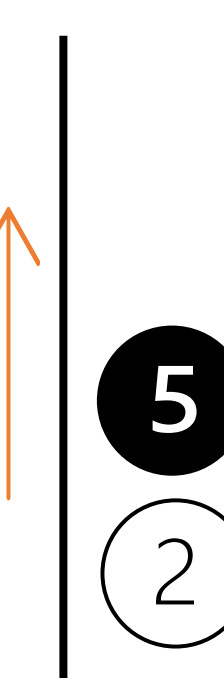
Temporal Part

### 3.Loop Permutation

	C=28			K=15	
Prime Factors	2	2	7	3	5
rank0	✓				
rank1					✓
rank2					
rank3					
rank4					

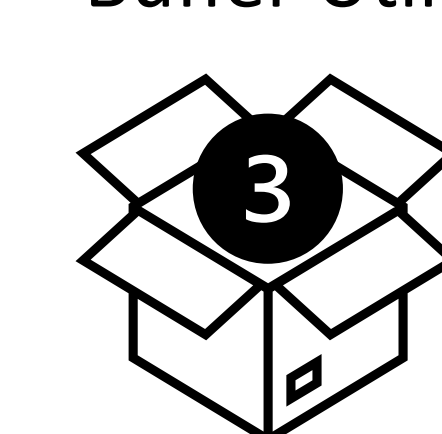
Binary allocation var X

Rank in Global Buffer:

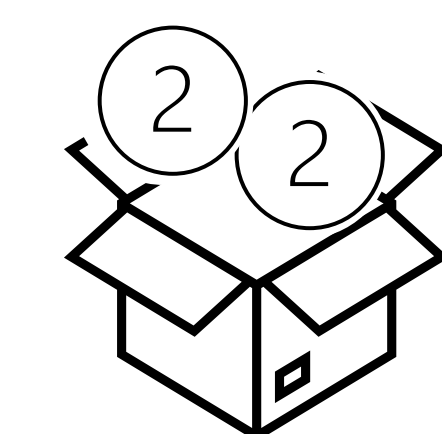


### CoSA Constraints

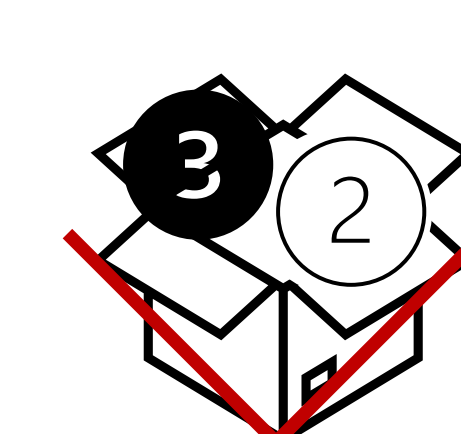
**Buffer Utilization:**



Weight Buffer (Size = 4)

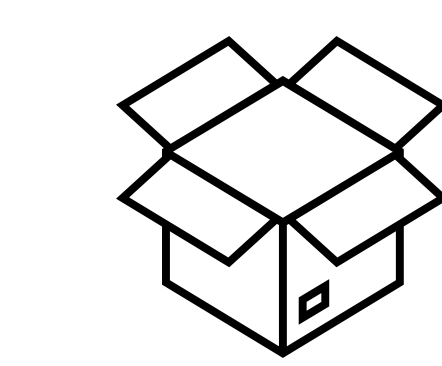


Weight Buffer (Size = 4)

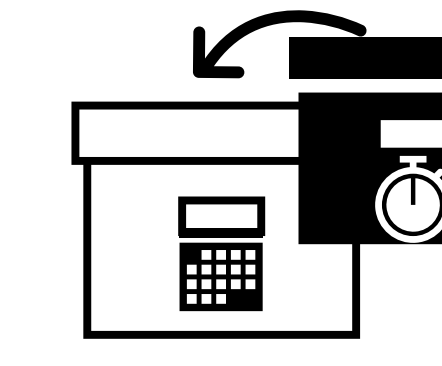


Weight Buffer (Size = 4)

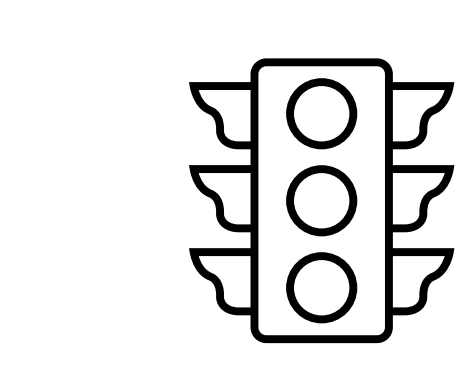
### CoSA Objectives



Utilization-driven



Compute-driven



Traffic-driven