Timeloop

Accelergy

Angshuman Parashar Yannan Nellie Wu Po-An Tsai Vivienne Sze Joel S. Emer NVIDIA MIT NVIDIA MIT NVIDIA, MIT

ISCA Tutorial

Hands-on session

May 2020







Resources

- Tutorial Website: http://accelergy.mit.edu/isca20_tutorial.html
- Tutorial Docker: https://github.com/Accelergy-Project/timeloop-accelergy-tutorial
 - Various exercises and example designs <u>and</u> environment setup for the tools



MOTIVATION

Infrastructure Download Instructions: http://accelergy.mit.edu/isca20_tutorial.html

EXPLOITING REUSE



MAPPING CHOICES

Energy-efficiency of peak-perf mappings of a single problem



480,000 mappings shown

Spread: 19x in energy efficiency

Only 1 is optimal, 9 others within 1%

A model needs a mapper to evaluate a DNN workload on an architecture

6,582 mappings have min. DRAM accesses but vary 11x in energy efficiency

A mapper needs a good cost model to find an optimal mapping



Infrastructure Download Instructions: http://accelergy.mit.edu/isca20_tutorial.html

TIMELOOP / ACCELERGY

Tools for Evaluation and Architectural Design-Space Exploration of DNN Accelerators



WHY TIMELOOP/ACCELERGY?

Microarchitectural model (Timeloop/Accelergy)

- Expressive: generic, template based hardware model
- Fast: faster than native execution on host CPUs
- Accurate: validated vs. design-specific models

Technology model (Accelergy)

- Allows user-defined complex architectural components
- Plugins for various technology models, e.g., Cacti, Aladdin, proprietary databases

Built-in Mapper (Timeloop)

• Addresses the hard problem of optimizing data reuse, which is required for faithful evaluation of a workload on an architecture



FUN WITH TIMELOOP

THE MODEL

Infrastructure Download Instructions: http://accelergy.mit.edu/isca20_tutorial.html



Architecture

Infrastructure Download Instructions: http://accelergy.mit.edu/isca20_tutorial.html

EXERCISE 0: PROBLEM



Infrastructure Download Instructions: http://accelergy.mit.edu/isca20_tutorial.html

EXERCISE 0: ARCHITECTURE

1-Level Temporal





Infrastructure Download Instructions: <u>http://accelergy.mit.edu/isca20_tutorial.html</u>

EXERCISE 0: MAPPING

1-Level Temporal





mapping:

target: Buffer
 type: temporal
 factors: R=3 P=16
 permutation: RP



EXERCISE 0

Follow the instructions in the README.



Infrastructure Download Instructions: http://accelergy.mit.edu/isca20_tutorial.html

EXERCISE 0

Run Timeloop model:

>> timeloop-model arch.yaml problem.yaml map.yaml

Output:

```
timeloop-model.map.txt
```

```
Buffer [ Weights:3 Inputs:18 Outputs:16 ]
| for P in [0:16)
| for R in [0:3)
```

<pre>timeloop-model.stats.txt</pre>	
•••••	
Summary Stats	
Utilization: 1.00 Cycles: 48 Energy: 0.00 uJ Area: 0.00 mm^2	
MACCs = 48 pJ/MACC MACC Buffer	= 0.60 = 1.54
Total	= 2.14

Follow the instructions in the exercise's README



EXERCISE 1: ARCHITECTURE

2-Level Temporal





EXERCISE 1: MAPPING

Weight Stationary





EXERCISE 1: MAPPING

Output Stationary



Phi 📀

EXERCISE 1

Follow the directions in the README.



EXERCISE 2: PROBLEM

Conv1D + Output Channels To represent this... Think about: And write: problem: for k = [0:K]shape: Operation for r = [0:R): name: Conv1D for p = [0:P): **Outputs** dimensions: [K, R, P] Output[k][p] += Weight[k][r] * Input[p+r]; data-spaces: - name: Weights Weights projection: Weights - [[K]] 0 Space - [[R]] Cti - name: Inputs R Projec Projection projection: R - [[P], [R]] - name: Outputs Inputs projection: - [[K]] Inputs W=P+R-1 - [[P]] read-write: True W=P+R-1 **Outputs** Data Spaces instance: K: 32 Κ R: 3 P: 16 Ρ

EXERCISE 2: MAPPINGS

Untiled vs. K-tiled





EXERCISE 2

Follow the directions in the README.



EXERCISE 2: O.S. DATAFLOW VARIANTS



Mit 💿

R

R

R

R

22 📀 nvidia

EXERCISE 3: ARCHITECTURE

3-Level Temporal





EXERCISE 3B: BYPASSING LEVELS

3-Level Temporal with Level Bypassing



EXERCISE 3B: BYPASSING

Bypassing

- Avoids energy cost of reading and writing buffers
- May result in additional accesses to outer buffers
- Does not change energy cost of moving data over network wires

For brevity in expressing mappings, Timeloop's evaluator assumes each datatype is stored at each level.

• We will see later that Timeloop's *mapper* makes no such assumption

Follow the directions in the README.

Challenge

• Experiment with bypass strategies to find out if there's any benefit in bypassing for this problem.



EXERCISE 4: SPATIAL INSTANCES

3-Level with multiple PEs



EXERCISE 4: MAPPING

Spatial levels need loops too





EXERCISE 4

Follow the directions in the README.



EXERCISE 4: SPATIAL INSTANCES

Spatial levels need to be mapped.

By convention, a block of spatial_for loops representing a spatial fanout from storage level *Outer* to storage level *Inner* are described as a spatial mapping directive targeted at level *Outer*.

Specifying complete mappings manually is beginning to get tedious. Space of choices and consequences is getting larger. Moving to realistic problem shapes and hardware topologies, we get a combinatorial explosion.

Fortunately, Timeloop's mapper was built exactly for this.



FUN WITH TIMELOOP

THE MAPPER

INVOKING THE MAPPER



To understand how the mapper works, let's go back to a simpler hardware architecture.



Arch: 3-Level, Problem: 1D + Output Channels



Recall:

mapping:

- target: MainMemory
 type: temporal
 factors: R=1 P=16 K=4
 permutation: RPK
- target: GlobalBuffer
 type: temporal
 factors: R=3 P=1 K=2
 permutation: RPK
 - target: RegisterFile type: temporal factors: R=1 P=1 K=4 permutation: RPK

Mapper constructs a mapping template:

mapping:

- target: MainMemory type: temporal factors: R=_ P=_ K=_ permutation: _ _ _
- target: GlobalBuffer
 type: temporal
 factors: R=_ P=_ K=_
 permutation: _ _ _
- target: RegisterFile
 type: temporal
 factors: R=_ P=_ K=_
 permutation: ____

Arch: 3-Level, Problem: 1D + Output Channels



Arch: 3-Level, Problem: 1D + Output Channels



Mapspace: An enumeration of ways to fill in these _ red blanks:

- Factors
- Permutations
- Dataspace Bypass

Mapspaces can be constrained by the user.

- Architecture constraints
- Mapspace constraints

Mapper constructs a mapping template:

mapping:

- target: MainMemory type: temporal factors: R=_ P=_ K=_ permutation: _ _ _
- target: GlobalBuffer type: temporal factors: R=_ P=_ K=_ permutation: ____
 target: RegisterFile type: temporal factors: R=_ P=1 K=1

permutation: **R**

Arch: 3-Level, Problem: 1D + Output Channels



Mapspace: An enumeration of ways to fill in these _ red blanks:

- Factors
- Permutations
- Dataspace Bypass

Mapspaces can be **constrained** by the user.

- Architecture constraints
- Mapspace constraints

Mapper runs a search heuristic over the constrained mapspace

Mapper constructs a mapping template:

mapping:

- target: MainMemory type: temporal factors: R=_ P=_ K=_ permutation: _ _ _
- target: GlobalBuffer
 type: temporal
 factors: R=_ P=_ K=_
 permutation: _ _ _
- target: RegisterFile
 type: temporal
 factors: R=_ P=1 K=1
 permutation: R

EXERCISE 5: MAPSPACE CONSTRAINTS

We provide 3 alternative sets of constraints:

- 1mapping: Constrain mapspace to the point that only 1 legal mapping remains in it!
- *freebypass:* Factors and permutations are forced, but bypass options are left unspecified.
 - Each of 3 dataspaces may either be kept or bypassed at each of the 2 inner levels (RegisterFile and GlobalBuffer) => (2²)³ = 64 choices!
 - Does Timeloop find a better bypassing strategy?
- *null*: Fully unconstrained.
 - How large is the mapspace?
 - Does Timeloop find a better mapping?



EXERCISE 6: PROBLEM

Convolutional Network Layer





EXERCISE 6: ARCHITECTURE

Eyeriss-256



EXERCISE 6: CNN LAYER ON EYERISS-256

Mapper is multi-threaded.

- Mapspace is split between each mapper thread.
- Default number of threads = number of logical CPUs on host machine.

For long mapper runs, you can use the interactive ncurses-based status tracker by setting mapper.live-status = True

- Tracks various statistics for each mapper thread:
 - Best energy-efficiency/performance seen so far
 - Number of legal/illegal/total mappings examined so far
 - Number of consecutive illegal mappings
 - Number of consecutive legal sub-optimal mappings



TUNING THE MAPPER'S SEARCH

Search heuristics (as of this recording)

- Linear
- Random
- Hybrid

Optimization criteria: prioritized list of statistics emitted by the model, e.g.,

- [cycles, energy]
- [last-level-accesses]

Termination conditions

- Mapspace exhausted
- #Valid mappings encountered >= "search-size"
- #Consecutive invalid mappings encountered >= "timeout"
- #Consecutive sub-optimal valid mappings encountered >= "victory-condition"
- Ctrl+C



HARDWARE X/Y DIMENSIONS



Architecture

name: GlobalBuffer
 class: SRAM
 attributes:

name: RegFile[0..11]
 class: regfile
 attributes:

.... meshX: 4

Mapping (also applies to Constraints)

mapping: target: GlobalBuffer type: spatial factors: C=4 K=3 R=1 S=1 P=1 Q=1 N=1 permutation: C_K_R_S_P_Q_N split: 1 X Y 0 1 2 3 4 5 6 7 41

41 💿 NVIDIA

HARDWARE X/Y DIMENSIONS



PARTITIONED BUFFERS



This is also a temporary workaround. Partitioned buffers will be supported natively in future.



EXERCISE 6

Follow the directions in the README.

Complete the exercise and enjoy!



TIMELOOP



Timeloop aims to serve as a vehicle for quality research on flexible DNN accelerator architectures. The infrastructure is released at https://github.com/NVlabs/timeloop under a BSD license.

Please join us in making Timeloop better and more useful for research opportunities across the community.



Resources

- Tutorial Related
 - Tutorial Website: <u>http://accelergy.mit.edu/isca20_tutorial.html</u>
 - Tutorial Docker: <u>https://github.com/Accelergy-Project/timeloop-accelergy-tutorial</u>
 - Various exercises and example designs <u>and</u> environment setup for the tools
- Other
 - Infrastructure Docker: <u>https://github.com/Accelergy-Project/accelergy-timeloop-infrastructure</u>
 - Pure environment setup for the tools without exercises and example designs
 - Open Source Tools
 - Accelergy: <u>http://accelergy.mit.edu/</u>
 - Timeloop: <u>https://github.com/NVlabs/timeloop</u>
 - Papers:
 - A. Parashar, et al. "Timeloop: A systematic approach to DNN accelerator evaluation," ISPASS, 2019.
 - Y. N. Wu, V. Sze, J. S. Emer, "An Architecture-Level Energy and Area Estimator for Processing-In-Memory Accelerator Designs," ISPASS, 2020.
 - Y. N. Wu, J. S. Emer, V. Sze, "Accelergy: An Architecture-Level Energy Estimation Methodology for Accelerator Designs," ICCAD, 2019.



