PLEASE INSTALL THE DOCKER FOR OUR FUN EXERCISES

Instructions on our tutorial website: <u>http://accelergy.mit.edu/tutorial.html</u>

- 0. If you do not have the docker app installed already, please install docker (community edition) (Windows users please manually turn on virtualization via BIOS settings.)
- 1. Go to our tutorial repo and get a copy of the single file docker-compose.yaml
- 2. Make an empty directory and place the docker-compose.yaml inside. Change the UID or GID as instructed in the file
- 3. To pull the newest docker image, run command docker-compose pull
- 4. Run command docker-compose run --rm exercises
- 5. Once invoked, you can cd ./exercises/timeloop/01-conv1d-2level and timeloop-model *.yaml to run a test example.

Updated test: cd ./exercises/timeloop/00-conv1d-1level; timeloop-model */*.yaml

(if needed)

TIMELOOP / ACCELERGY TUTORIAL



Joel S. Emer Vivienne Sze Angshuman Parashar Po-An Tsai Yakun Sophia Shao Yannan Nellie Wu NVIDIA, MIT MIT NVIDIA NVIDIA UC Berkley MIT



SCHEDULE

- 1:00 2:30PM Timeloop lecture + exercises
- 2:30 3:00PM Timeloop free lab time
- 3:00 3:30PM Coffee break
- 3:30 4:30PM Accelergy lecture + exercises
- 4:30 5:00PM Accelergy + Timeloop free lab time

MOTIVATION

DNN ACCELERATORS

Design Considerations



DNN ACCELERATORS

Design Considerations



DATA MOVEMENT

Why it's important

Energy costs	
8-bit Integer Multiply	0.2 pJ
Fetch two 8-bit operands from DRAM	128 pJ
Fetch two 8-bit operands from large SRAM	2 pJ

VGG16 conv 3_2		
Multiply Add Ops	1.85 Billion	
Weights	590 K	Re-use
Inputs	803 K	ine use
Outputs	803 K	

Fortunately...

EXPLOITING REUSE



Flexible architectures may allow millions of alternative mappings of a single workload

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

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

TIMELOOP VALIDATION

VALIDATION: EYERISS

Vs. ISCA 2016 Eyeriss Energy Model



VALIDATION: NVDLA-DERIVED PE (ENERGY)



DeepBench Workload

Within 8% error across all workloads

VALIDATION: NVDLA-DERIVED PE (PERFORMANCE)



CASE STUDIES

CASE STUDY: TECHNOLOGY MODEL



CASE STUDY: MEM HIERARCHY



CASE STUDY: EXISTING DESIGNS



CASE STUDY: EXISTING DESIGNS



FUN WITH TIMELOOP

THE MODEL





^{23 📀} nvidia.

EXERCISE 0: ARCHITECTURE

1-Level Temporal



EXERCISE 0

Follow the instructions in the README.

EXERCISE 1: ARCHITECTURE

2-Level Temporal



26 📀 nvidia.

EXERCISE 1: MAPPING

Weight Stationary



EXERCISE 1: MAPPING

Output Stationary



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 Ŭ projection Proj 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



33 📀 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.

FUN WITH TIMELOOP

THE MAPPER

THE MAPPER



EXERCISE 5: ARCHITECTURE

3-Level Temporal



EXERCISE 5: MAPSPACE CONSTRAINTS

Retain old "mapping" directives, now treat them as *mapspace constraints*.

Does this set of constraints result in only 1 legal mapping?

- No! Bypass options haven't been set. If you fire the model directly on the mapping, it assumes no bypass. But if you fire the *mapper* on these constraints, it recognizes that bypass has been left unspecified and explores all options.
- Each of 3 dataspaces may either be kept or bypassed at each of the 2 inner levels (RegisterFile and GlobalBuffer) => (2²)³ = 64 choices!

Run Timeloop. Does it find a better bypassing strategy?

Comment out all the factors and permutations, comment out the num-threads parameter and re-run Timeloop. Does it find a better mapping?

EXERCISE 5

Follow the directions in the README.

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 neurses-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

EXERCISE 6: TUNING THE MAPPER'S SEARCH

Termination conditions

- 1. Mapspace exhausted
- 2. #Valid mappings encountered >= "search-size"
- 3. #Consecutive invalid mappings encountered >= "timeout"
- 4. #Consecutive sub-optimal valid mappings encountered >= "victory-condition"
- 5. Ctrl+C

EXERCISE 6

Follow the directions in the README.

DEEP DIVE: UNDERSTANDING THE MAPPER

CANONICAL MAPPING REPRESENTATION

Mapping: way in which the operation space and the associated data spaces are split into tiles at each level of the architecture's hierarchy and among multiple instances at each level



MAPSPACES

Set of legal mappings of a workload onto an architecture

```
// === 1D Convolution Workload ===
for r=[0:R):
  for p=[0:P):
    Output[p] += Weight[r] * Input[r+p];
// === Mapping ===
// DRAM
for r3=[0:R3):
  for p3=[0:P3):
    // GBuf
    for r2=[0:R2):
      for p2=[0:P2):
      // Spatial: GBuf->RFile
      parallel for r1=[0:R1):
        parallel for p1=[0:P1):
        // RFile
        for r0=[0:R0):
          for p0=[0:P0):
            p = p3*P2*P1*P0 + p2*P1*P0 + p1*P0 + p0;
            r = r3*R2*R1*R0 + r2*R1*R0 + r1*R0 + r0;
            Output[p] += Weight[r] * Input[r+p];
```



MAPSPACE CONSTRAINTS

Index Factorization and Loop Permutation

```
// === 1D Convolution Workload ===
for r=[0:R):
 for p=[0:P):
   Output[p] += Weight[r] * Input[r+p];
// === Mapping ===
                                                                               target = "RFile";
// DRAM
for r3=[0:R3):
                                                                               factors = "R=3";
 for p3=[0:P3):
                                                                               permutation = "PR";
    // GBuf
    for r2=[0:R2):
     for p2=[0:P2):
                                                                   Note that P is a free variable. Timeloop will
     // Spatial: GBuf->RFile
                                                                   automatically choose an optimal P and will know not
     parallel for r1=[0:R1):
                                                                   to exceed the capacity of the buffer at this level.
       parallel_for p1=[0:P1):
       // RFile
                                                                   Inequalities are allowed, e.g., P<=16.
       for r0=[0:3):
         for p0=[0:P0):
                                                                   Permutations can be partially specified.
           p = p3*P2*P1*P0 + p2*P1*P0 + p1*P0 + p0;
           r = r3*R2*R1*R0 + r2*R1*R0 + r1*R0 + r0;
           Output[p] += Weight[r] * Input[r+p];
```

MAPSPACE CONSTRAINTS

Buffer Level Bypassing



CONSTRAINING A MAPSPACE



PRUNING A MAPSPACE



PRUNING A MAPSPACE



THE MAPPER



MULTI-THREADING



DEEP DIVE: UNDERSTANDING THE MODEL

THE MODEL



THE MODEL: TILE ANALYSIS

delta

Mapping



Determines tiles, spatial partitioning of tiles, and schedule Point sets: at each loop iteration (time step OR instance of a hardware unit)

Deltas: set-difference between point sets

tile *i*+1

tile i

- Temporal: Indicates stationarity, slidingwindow behavior, etc.
- Spatial: Indicates overlaps/halos between adjacent units, multicasting/forwarding opportunities

<u>Tile Analysis</u>: Measure and record deltas over all space and time.

Naïve but robust approach: *simulate* execution of entire loop nest.

Regular problems:

- Compute 1st, 2nd, and last iterations of each loop
- Point sets are Axis-Aligned Hyper Rectangles (AAHR)

THE MODEL: UARCH MODEL



ESTIMATING PERFORMANCE AND ENERGY

Performance: Throughput of rate-limiting step across:

- Multipliers
- Buffer read/write ports
- Networks

Assumption: Buffers are either double-buffered or use buffets*

Energy: summation of costs for various activities:

- Multiplier accesses
- Buffer accesses
- Network transfers
- Temporal and spatial accumulations
- Address-generator in



* Buffets: An Efficient and Composable Storage Idiom for Explicit Decoupled Data Orchestration; Michael Pellauer, Yakun Sophia Shao, Jason Clemons, Neal Crago, Kartik Hegde, Rangarajan Venkatesan, Stephen W. Keckler, Christopher W Fletcher, Joel Emer; ASPLOS 2019

FUTURE WORK

Search Heuristics

Workloads: Complete networks, with inter-layer optimization

Modeling fragmentation and load imbalance in architectures that exploit unstructured sparsity for performance



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.





DNN ACCELERATORS

Design Considerations



BUFFER PARAMETERS



CASE STUDY: BENCHMARK SWEEP



Problem Index

ARCHITECTURE SPECIFICATION: EXAMPLE


MAPPINGS

Mapping: way in which the operation space and the associated data spaces are split into tiles at each level of the architecture's hierarchy and among multiple instances at each level



MAPPING CONSTRAINTS

Factors, Permutation and LevelBypass

```
target = "RFile";
```

```
type = "temporal";
```

```
factors = "Q1 R1 S1 C1 K1 N1";
```

```
permutation = "PQRSCKN";
```

Note: P is a free variable. Timeloop
will automatically choose an optimal
P and will know not to exceed the
capacity of the buffer at this level.

```
target = "GBuf_RFile";
```

```
type = "spatial";
```

```
factors = "C16 K1 R1 S1 P1 Q1 N1";
```

```
target = "RFile";
type = "datatype";
keep = [ "Weights" ];
bypass = [ "Inputs", "Outputs" ];
```

BUFFER MODEL



Buffets: An Efficient and Composable Storage Idiom for Explicit Decoupled Data Orchestration; Michael Pellauer, Yakun Sophia Shao, Jason Clemons, Neal Crago, Kartik Hegde, Rangarajan Venkatesan, Stephen W. Keckler, Christopher W Fletcher, Joel Emer; ASPLOS 2019

CASE STUDY: AREA VS. ENERGY



PROBLEM SPECIFICATION

Deep Loop Nest



PROBLEM SPECIFICATION

Deep Loop Nest





ARCHITECTURE SPECIFICATION



MAPPINGS

Mapping: way in which the operation space and the associated data spaces are split into tiles at each level of the architecture's hierarchy and among multiple instances at each level



EXERCISE 6: PROBLEM

7D Convolutional Network Layer



MAPPINGS

Mapping: way in which the operation space and the associated data spaces are split into tiles at each level of the architecture's hierarchy and among multiple instances at each level



MULTIPLE TOOLS WITHIN TIMELOOP









MAPPER TERMINATION CONDITIONS

Mapper termination conditions

- 1. Mapspace exhausted
- 2. #Valid mappings encountered >= "search-size"
- 3. #Consecutive invalid mappings encountered >= "timeout"
- 4. #Consecutive sub-optimal valid mappings encountered >= "victory-condition"
- 5. Ctrl+C

ALTERNATE INTRO

EVALUATING ARCHITECTURES



EVALUATING ARCHITECTURES



Flexible accelerator architectures



EVALUATING ARCHITECTURES



Flexible accelerator architectures

