Accelergy: An Architecture-Level Energy Estimation Methodology for Accelerator Designs

MICRO Tutorial

October 12, 2019

Website: http://accelergy.mit.edu/tutorial.html



Massachusetts Institute of Technology





Accelergy Overview

- An architecture-level energy estimator that can be easily integrated into other design exploration tools
- Accurately and systematically captures the energy consumption of basic building blocks in the design
- Succinctly models diverse and complicated designs
- Achieves 95% accuracy in evaluating a well-know deep neural network accelerator design – Eyeriss [2016 ISSCC]



Accelergy Overview

- An architecture-level energy estimator that can be easily integrated into other design exploration tools
- Accurately and systematically captures the energy consumption of basic building blocks in the design
- Succinctly models diverse and complicated designs
- Achieves 95% accuracy in evaluating a well-know deep neural network accelerator design – Eyeriss [2016 ISSCC]



Importance of Energy Estimation

Energy efficiency is the general concern for accelerator designs



We must quickly evaluate energy efficiency of arbitrary potential designs in the <u>large design space</u>

Phir

Physical-level energy estimators



• Physical-level energy estimators





• Physical-level energy estimators



Physical-Level Energy Estimator



• Physical-level energy estimators



• Physical-level energy estimators



Physical-level energy estimators





• Architecture-Level Energy Estimators



• Architecture-Level Energy Estimators



Description

Architecture -Level Energy Estimator



Architecture-Level Energy Estimators



Action Counts



Architecture-Level Energy Estimators





Importance of Energy Estimations





Importance of Energy Estimations



• Design-Specific Accelerator Estimators: Aladdin[ISCA2014], fixed-cost[Asilomar2017]





Design-Specific Accelerator Estimators •



Architecture Description



Design-Specific Accelerator Estimators



Energy Estimator Energy/

action

100pJ

10pJ

5pJ



• Design-Specific Accelerator Estimators

Comp. name	Action name	Energy/ action
1024B GLB	access	100pJ
128B Buffer	access	10pJ
16b MAC	compute	5pJ

Design-Specific ERT





• Design-Specific Accelerator Estimators

Comp.
nameAction
nameEnergy/
action1024B GLBaccess100pJ128B Bufferaccess10pJ16b MACcompute5pJ

Design-Specific ERT

Coarse-grained estimations reduce estimation accuracies

Energy-Per-Actions of a Register File (normalized to idle) 4.7 Write the same data to ~5x different addresses 2.4 2.1 1.8 1.0 Random Repeated Random Repeated Constant Read Read Write Write Data Write **Repetitive access to the** same memory location



Design-Specific Accelerator Estimators





Energy/

action

100pJ

10pJ

5pJ



Design-Specific Accelerator Estimators







Design-Specific Accelerator Estimators



Design-Specific Accelerator Estimators



Action Counts

14115

Design-Specific Accelerator Estimators



Accelergy Overview

- An architecture-level energy estimator that can be easily integrated into other design exploration tools
- Accurately and systematically captures the energy consumption of basic building blocks in the design
- Succinctly models diverse and complicated designs
- Achieves 95% accuracy in evaluating a well-know deep neural network accelerator design – Eyeriss [2016 ISSCC].



- Accurate estimation for primitive components with a <u>primitive</u> <u>component library*</u> that defines
 - Necessary hardware attributes to accurately characterize each component
 - Necessary fine-grained actions that improve energy estimation accuracy



Fine-grained memory action types



Random Mult Reused Mult Gated Mult

Fine-grained multiplier action types

*open-source at http://acelergy.mit.edu 28







Use energy estimation plug-ins to characterize primitive components



*available at http://accelergy.mit.edu












Accelergy: Primitive Component Energy Estimation

Systematically generates ERTs for primitive component in various designs



Accelergy: Primitive Component Energy Estimation

Systematically generates ERTs for primitive component in various designs



How to use Accelergy?









YAML Syntax

- Contains various types of primitive component classes
- Each primitive component class contains attribute names and action names that can be shared by its instances

Object-Oriented Approach

Components are **instances** of component classes

version: 0.2 # vcs
classes:
- name: bitwise
attributes:
technology
datawidth
actions:
- name: process
- name: intadder
attributes:
technology
datawidth
actions:
- name: add



YAML Syntax

- Contains various types of primitive component classes
- Each primitive component class contains attribute names and action names that can be shared by its instances
- Default attributes are defined

Object-Oriented Approach Components are instances of component classes

```
version: 0.2 \# vcs
classes:
  name: bitwise
  attributes:
     technology: 40nm
     datawidth: 1
 actions:
    - name: process ...
  name: intadder
  attributes:
     technology: 40nm
     datawidth: 16 ...
  actions:
     name: add ...
```

- Contains various types of primitive component classes
- Each primitive component class contains attribute names and action names that can be shared by its instances
- Default attributes are defined

Object-Oriented Approach Components are instances of component classes

YAML Syntax



Representation of fine-grained action types

14111



 Fine-grained actions with runtime arguments

Original action type	New Action type	
Repeated read	road	
Random read	reau	
Repeated write		
Random write	write	
Repeated data write		

YAML Syntax

- name: regfile
 - attributes: ...

actions:

- name: read

- name: write



• Fine-grained actions with runtime arguments

Original action type Action type	New	Argument	
	data∆	address∆	
Repeated read	road	0	0
Random read	reau	1	1
Repeated write		0	0
Random write	write	1	1
Repeated data write		0	1

YAML Syntax

name: regfile attributes: ... actions: - name: read 0≤ data_delta ≤1 arguments: data delta: 0..1 address delta: 0..1 name: write arguments: data delta: 0..1 address delta: 0..1





• Accelergy shares similar architecture description syntax with Timeloop



YAML Syntax









• Accelergy shares similar architecture description syntax with Timeloop



Mit 💿



YAML Syntax



• Accelergy shares similar architecture description syntax with Timeloop



Mit 💿



YAML Syntax





Accelergy: Energy Reference Table (ERT)

 Contains energy/action values for all the components in the designs YAML Syntax







Accelergy: Energy Reference Table (ERT)

 Contains energy/action values for all the components in the designs YAML Syntax



```
tables:
  name:Design.GLB
  read:
    - arguments:
        data delta: 1
        address delta: 1
      energy: 100
  ...
  name: Design.PE.buffer
  read: ...
  name: Design.PE.MAC
 mac random:
    energy: 5
  ...
```



Exercise 1: Generate ERT for Simple Architecture

- Instructions:
 - cd exercises/accelergy/01_primitive_architecute_ERT
 - accelergy input/*.yaml -o output/ --enable_flattened_arch 1
 - ERT.yaml, ERT_summary.yaml, and flattened_architecture will be generated in current directory
- Questions:
 - Each YAML file contains a top-key that indicates its content. What top-level keys are used to identify the input file content?
 - What information in ERT is used for Timeloop?
 - Examine ERT_summary.yaml, why it's bad to have coarse-grained action type?





 Action counts are generated using external performance models, such as a cycle-level simulator or Timeloop.
 <u>Subtree</u>:



 Action counts are generated using external performance models, such <u>YAML Syntax</u>
 as a cycle-level simulator or Timeloop.
 whtreet







Exercise 2: Generate Energy Estimation for Simple Architecture

- Instructions:
 - cd .../02_primitive_architecute_estimation
 - accelergy *.yaml --enable_flattened_arch 1 -o output/
 - estimation.yaml will be generated in current directory
- Questions:
 - Did Accelergy run ERT generator part? Why?
 - Is energy calculation faster than ERT generation?
 - Would adding the architecture description in this directory confuse Accelergy to rerun ERT generator?
 - Why is separating the two parts good?



- Practical architecture designs involve much more details
 - Example: storage units with local address generators (AG)*



- AG_buffer_SRAM is an abstract hierarchy
- Buffer is an instance of SRAM primitive class
- AGs are instances of counter primitive class

64



- Practical architecture designs involve much more details
 - Example: storage units with local address generators (AG)



- AG_buffer_RF is an abstract hierarchy
- AGs are instances of counter primitive class
- Buffer is an instance of regfile primitive class



- Practical architecture designs involve much more details
 - Example: storage units with local address generators (AG)



- Practical architecture designs involve much more details
 - Example: storage units with local address generators (AG)



- Practical architecture designs involve much more details
 - Example: storage units with local address generators (AG)



Design













- Existing work that aims to succinctly model complicated architectures
 - Wacttch, McPAT, PowerTrain





Use a fixed set of <u>compound components</u> to represent model the architecture Components that can be decomposed into lower level components


Modeling Complicated Designs

• Existing work that aims to succinctly model complicated architectures



The fixed set of compound components is not sufficient to describe arbitrary accelerator designs

Accelergy Overview

- An architecture-level energy estimator that can be easily integrated into other design exploration tools
- Accurately and systematically captures the energy consumption of basic building blocks in the design
- Succinctly models diverse and complicated designs
- Achieves 95% accuracy in evaluating a well-know deep neural network accelerator design – Eyeriss [2016 ISSCC].





Previous: Abstract Hierarchies



• Allows users to describe the design-specific compound components



Now: User-Defined Compound Components





Previous: Architecture described with primitive components

Ш

Allows Architecture to be described with compound components





How to use Accelergy?

- 1. Estimate architectures with primitive components
- 2. Estimate architectures with compound components







Exercise 3a

- Instructions:
 - cd ../03_compound_architecute_estimation
- Questions:
 - Does the architecture description with compound components follow the previous general format?
 - Does the action counts follow the previous general format?



Exercise 3a

- Instructions:
 - cd ../03_compound_architecute_estimation
- Questions:
 - Does the architecture description with compound components follow the previous general format?
 - Does the action counts follow the previous general format?

Architecture description and Action counts follow the same general format











Accelergy: Compound Component Description



Accelergy: Compound Component Description



Accelergy: Compound Component Description







Exercise 3b

- Instruction:
 - accelergy *.yaml components/*.yaml -o output/
- Questions:
 - Are the ERTs in terms of compound components or primitive components?
 Why?
 - Compare the current ERT with the previous ERT, what is different?
 - If we change the hardware implementation of a compound component,
 e.g. add a FIFO to AG_buffer_SRAM, what needs to be changed?



Accelergy Overview

- An architecture-level energy estimator that can be easily integrated into other design exploration tools
- Accurately and systematically captures the energy consumption of basic building blocks in the design
- Succinctly models diverse and complicated designs
- Achieves 95% accuracy in evaluating a well-know deep neural network accelerator design – Eyeriss [2016 ISSCC].



Energy Evaluations on Eyeriss

- **Experimental Setup:**
 - Workload: Alexnet weights & ImageNet input feature maps
 - Ground Truth: Energy obtained from post-layout simulations



93

Energy Evaluations on Eyeriss

- Experimental Setup:
 - Workload: Alexnet weights & ImageNet input feature maps
 - Ground Truth: Energy obtained from post-layout simulations



- Total energy estimation is 95% accurate of the post-layout energy.
- Estimated relative breakdown of the important units in the design is within 8% of the post-layout energy.





- Comparisons with existing work: Aladdin and fixed-cost
- Aladdin: [Shao et al., ISCA2014]
 - lack of fine-grained action types, e.g., gated-read
- fixed-cost: [Yang et al., Asilomar2017]
 - <u>Designed specifically for Eyeriss analysis</u>
 - lack of fine-grained component characterizations, e.g., all local buffers are characterized as identical components



Design-specific

energy estimators

• Comparisons with existing work: Aladdin and fixed-cost





• Comparisons with existing work: Aladdin and fixed-cost





PEs that process data of different sparsity



• Comparisons with existing work: Aladdin and fixed-cost

Not sensitive to sparsity as the zero-gating optimization is not recognized



Energy Breakdown of PEs across the Array

PEs that process data of different sparsity

• Comparisons with existing work: Aladdin and fixed-cost

Missing components characterizations



Energy Breakdown of PEs across the Array

PEs that process data of different sparsity



• Comparisons with existing work: Aladdin and fixed-cost

Energy Breakdown of components inside a PE





• Comparisons with existing work: Aladdin and fixed-cost

Energy Breakdown of components inside a PE



Exercise 4

- Instructions
 - cd ../04_eyeriss_like
 - Run accelergy *.yaml components/*.yaml –o output/
- Questions:
 - What compound components are involved in the architecture?
 - What are the subcomponents of each compound component?
 - What actions in the ERT can help better characterize data-gating?
 - Which components in the PE are sensitive to data sparsity? How do you know?



Free Lab Session

Timeloop + Accelergy



Have fun exploring both tools

- Go to the timeloop + accelergy exercise folder
 - Readme contains instructions for various experiments!
 - Run an mapping exploration with an eyeriss-like architecture that is described by compound components
 - Run an mapping exploration with an floating-point eyeriss-like architecture. Which component's energy increases the most?
 - Try different DRAM types. Modify the architecture description to exploration the energy impact brought by DRAM technologies.

