

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

Yannan Nellie Wu<sup>1</sup> , Joel S. Emer<sup>1,2</sup> , Vivienne Sze<sup>1</sup>

<sup>1</sup> MIT

<sup>2</sup> NVIDIA



# Accelergy Overview

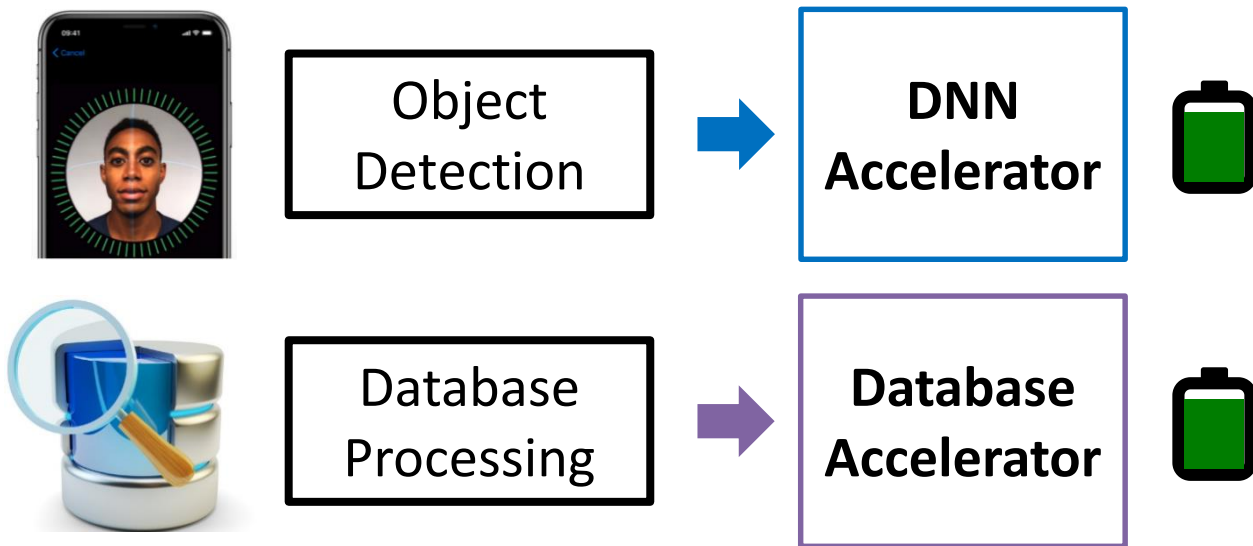
---

- **An architecture-level energy estimator**
- **Flexibly characterizes various basic building blocks of different technologies**
- **Succinctly models diverse and complicated designs**
- **Improves estimation accuracy via fine-grained classification of operations**
- **Achieves 95% accuracy in evaluating a deep neural network (DNN) accelerator – Eyeriss [ISSCC 2016]**

# Energy Consumption Concerns

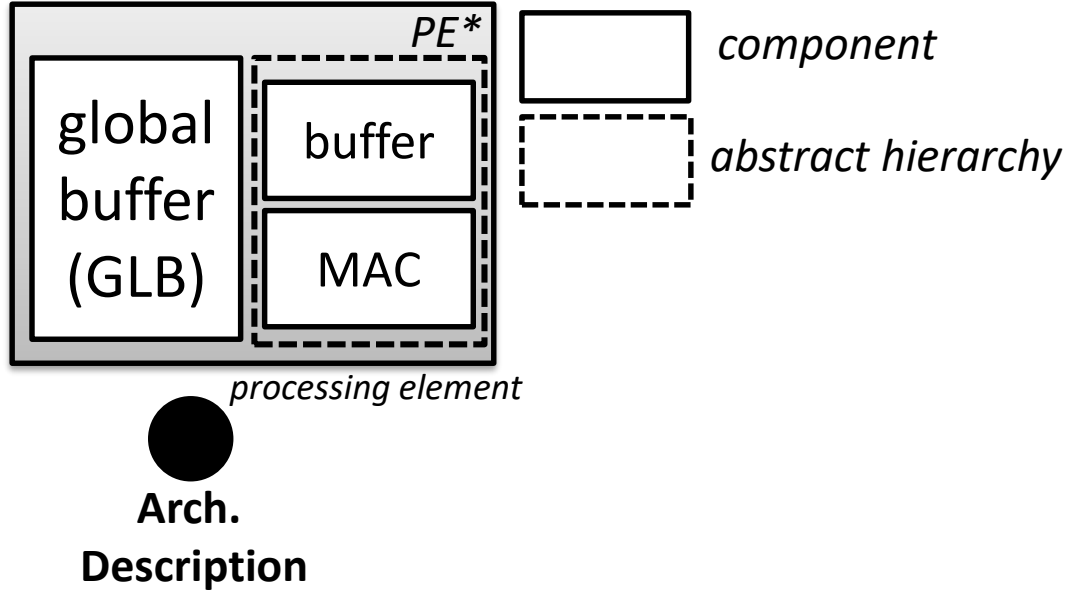
---

Data and computation-intensive applications are power hungry



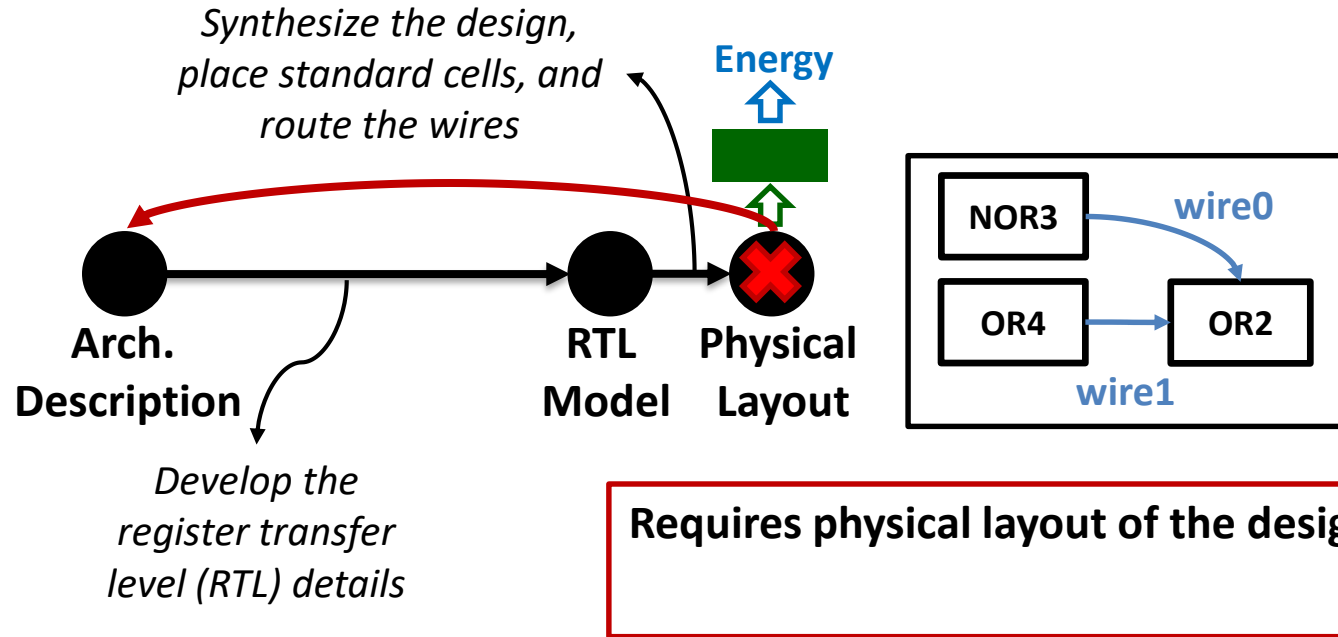
**We must quickly evaluate energy efficiency of arbitrary potential designs in the large design space**

# Energy Estimation and Design Exploration



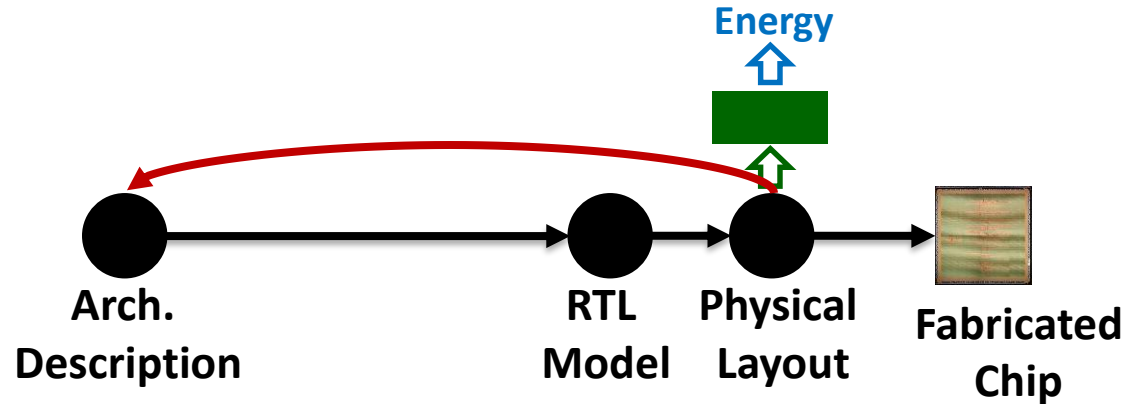
# Energy Estimation and Design Exploration

- Physical-Level Energy Estimator (Synopsys Prime Power, Cadence Joules)



# Energy Estimation and Design Exploration

- Physical-Level Energy Estimator (Synopsys Prime Power, Cadence Joules)



Requires physical layout of the design  
Slow design space exploration

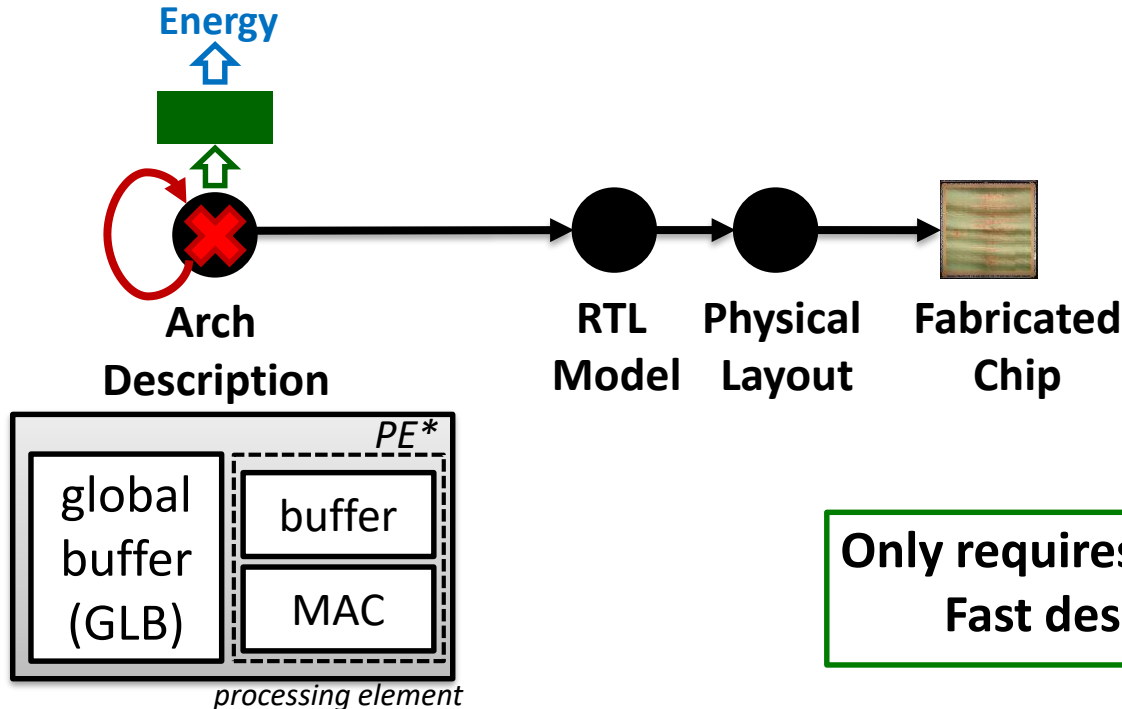
# Accelergy Overview

---

- **An architecture-level energy estimator**
- Flexibly characterizes various basic building blocks in the design
- Succinctly models diverse and complicated designs
- Improves estimation accuracy via fine-grained classification of operations
- Achieves 95% accuracy in evaluating a deep neural network (DNN) accelerator – Eyeriss [ISSCC 2016]

# Energy Estimation and Design Exploration

- Architecture-Level Energy Estimators

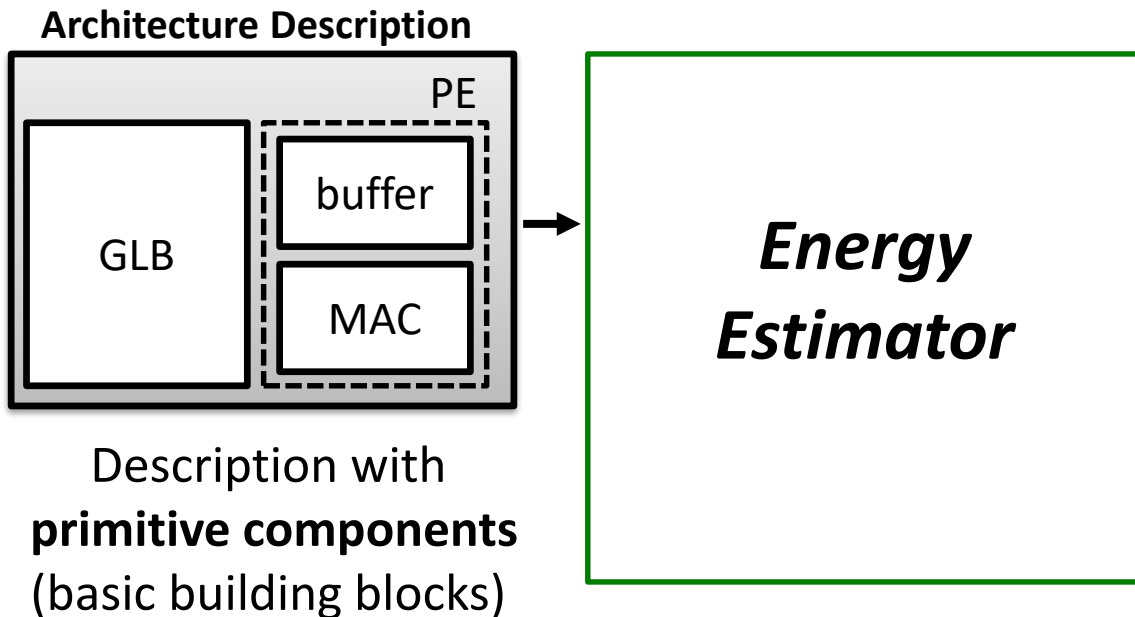


Only requires architecture-level design  
Fast design space exploration



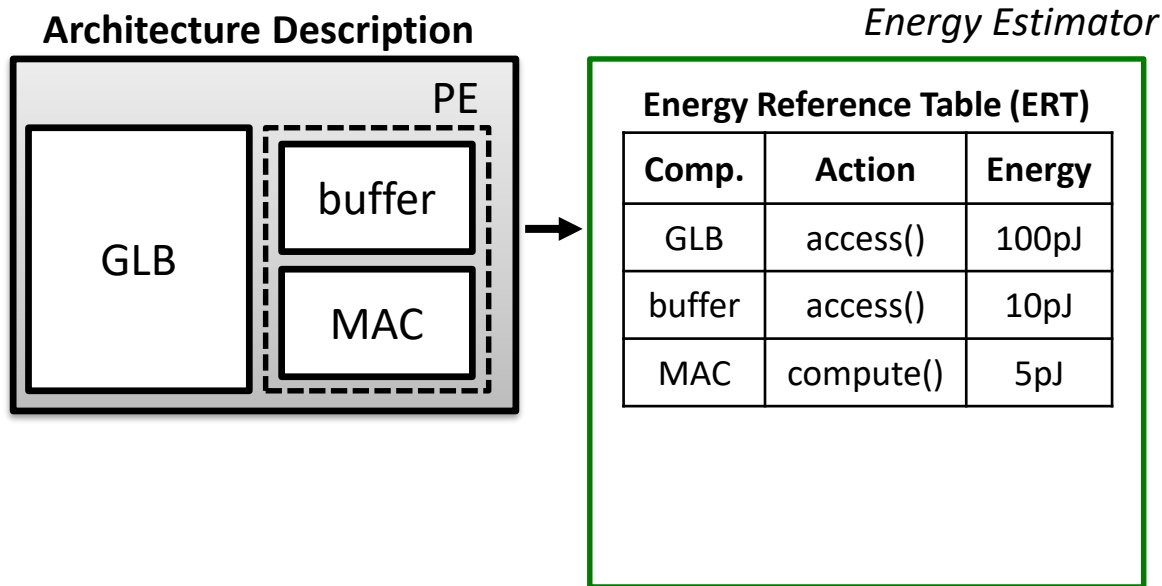
# Existing Architecture-Level Energy Estimators

- **Design-Specific Accelerator Estimators:** Aladdin<sup>[ISCA2014]</sup>, fixed-cost<sup>[Asilomar2017]</sup>



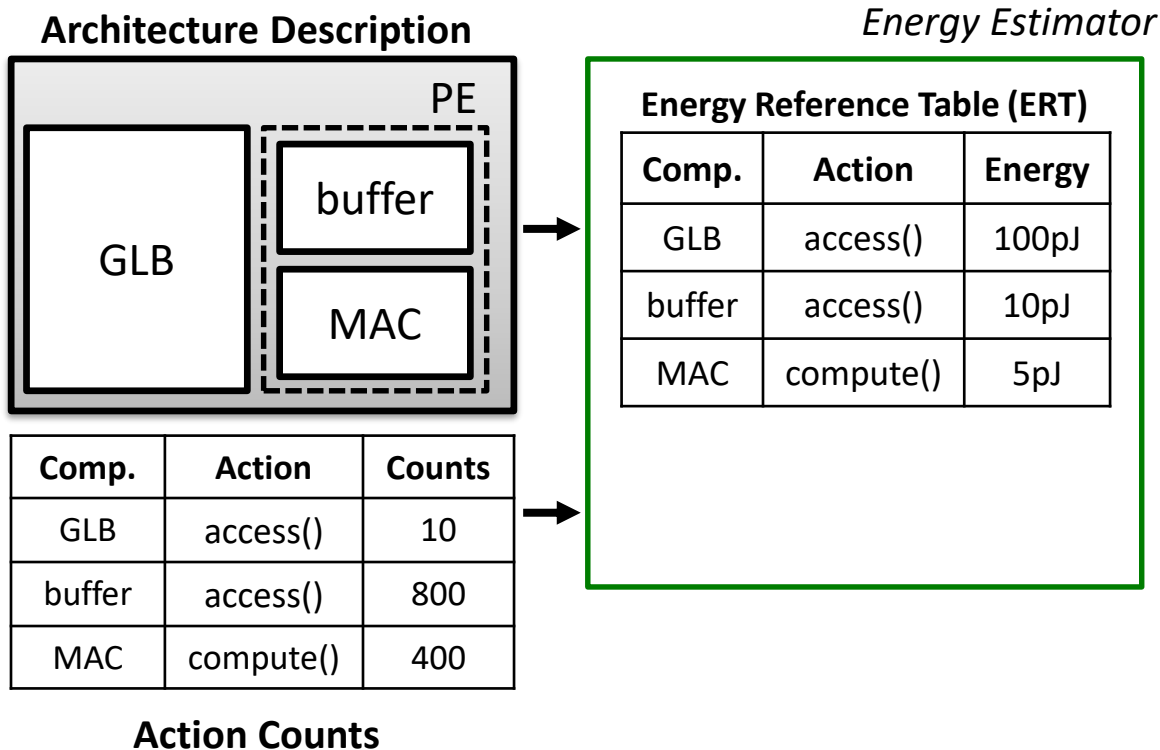
# Existing Architecture-Level Energy Estimators

- **Design-Specific Accelerator Estimators: Aladdin**<sup>[ISCA2014]</sup>, **fixed-cost**<sup>[Asilomar2017]</sup>



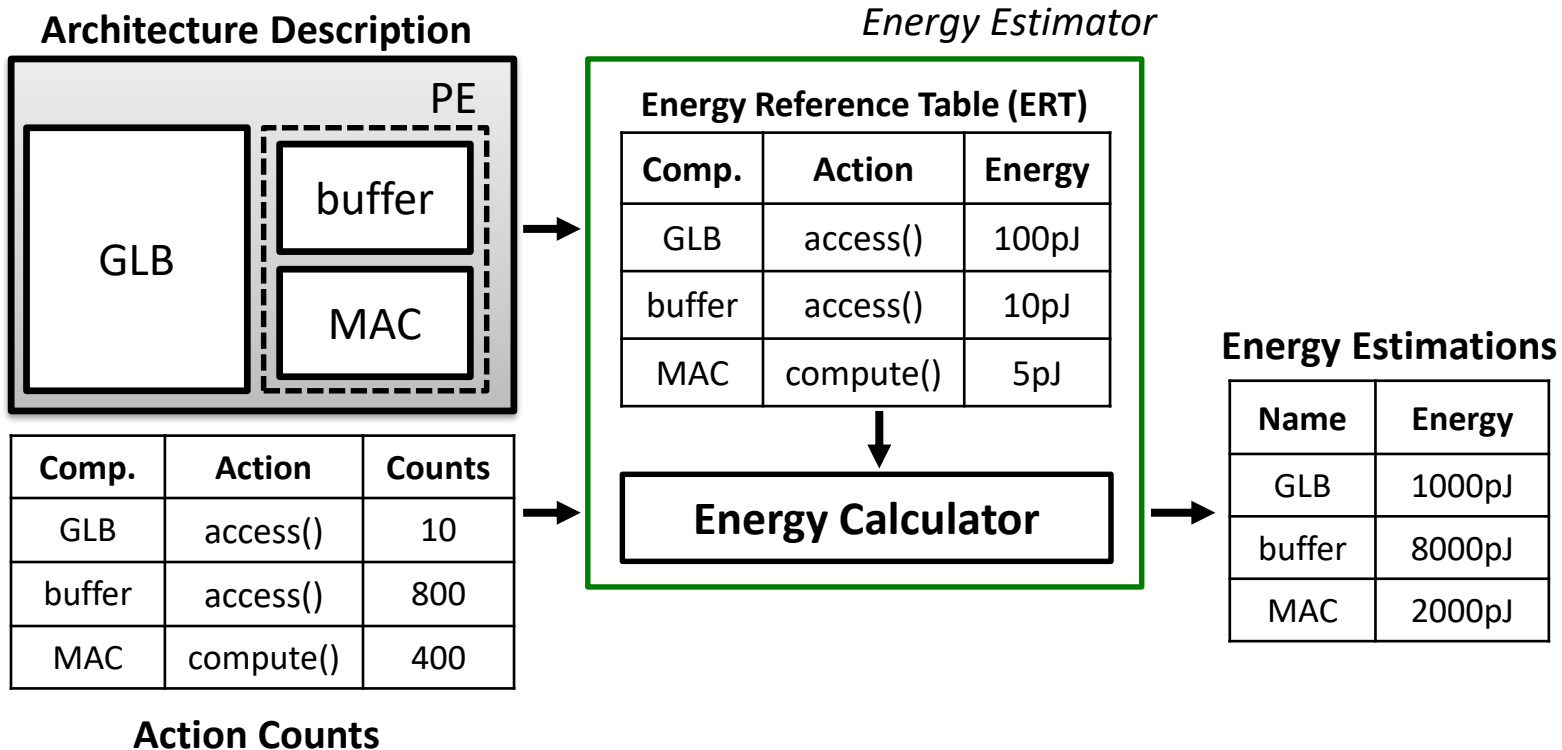
# Existing Architecture-Level Energy Estimators

- Design-Specific Accelerator Estimators: **Aladdin**<sup>[ISCA2014]</sup>, **fixed-cost**<sup>[Asilomar2017]</sup>



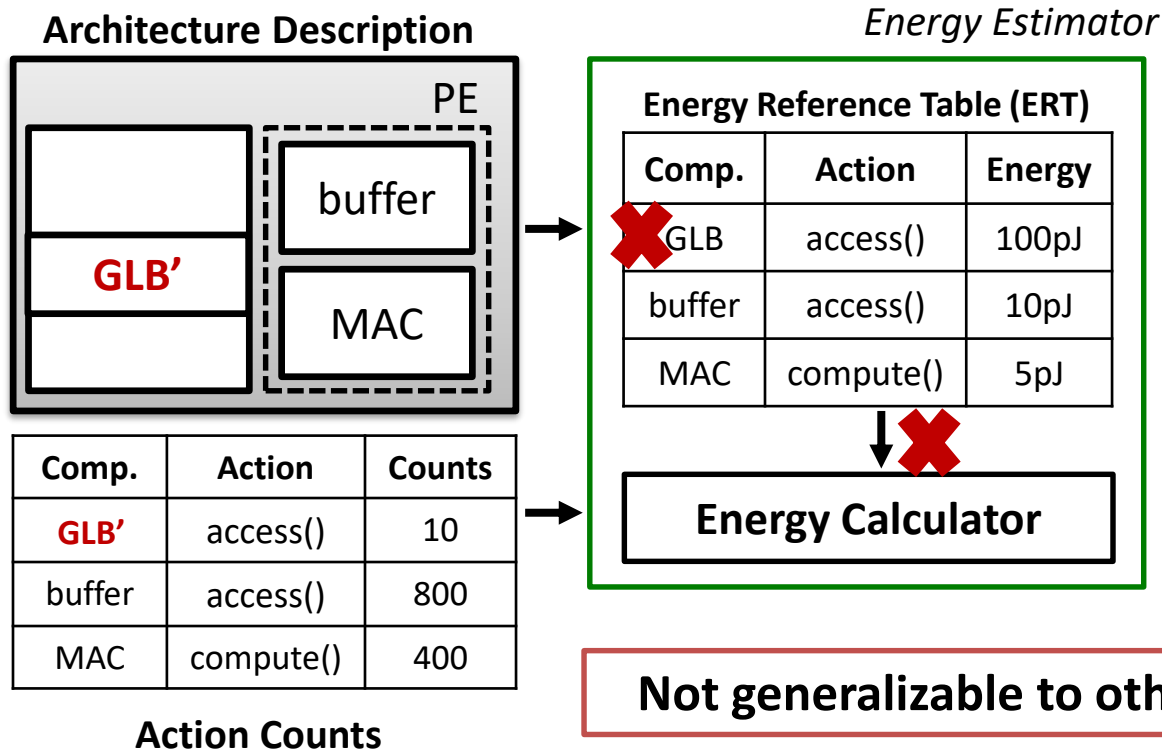
# Existing Architecture-Level Energy Estimators

- Design-Specific Accelerator Estimators: **Aladdin**<sup>[ISCA2014]</sup>, **fixed-cost**<sup>[Asilomar2017]</sup>



# Existing Architecture-Level Energy Estimators

- Design-Specific Accelerator Estimators: **Aladdin**<sup>[ISCA2014]</sup>, **fixed-cost**<sup>[Asilomar2017]</sup>

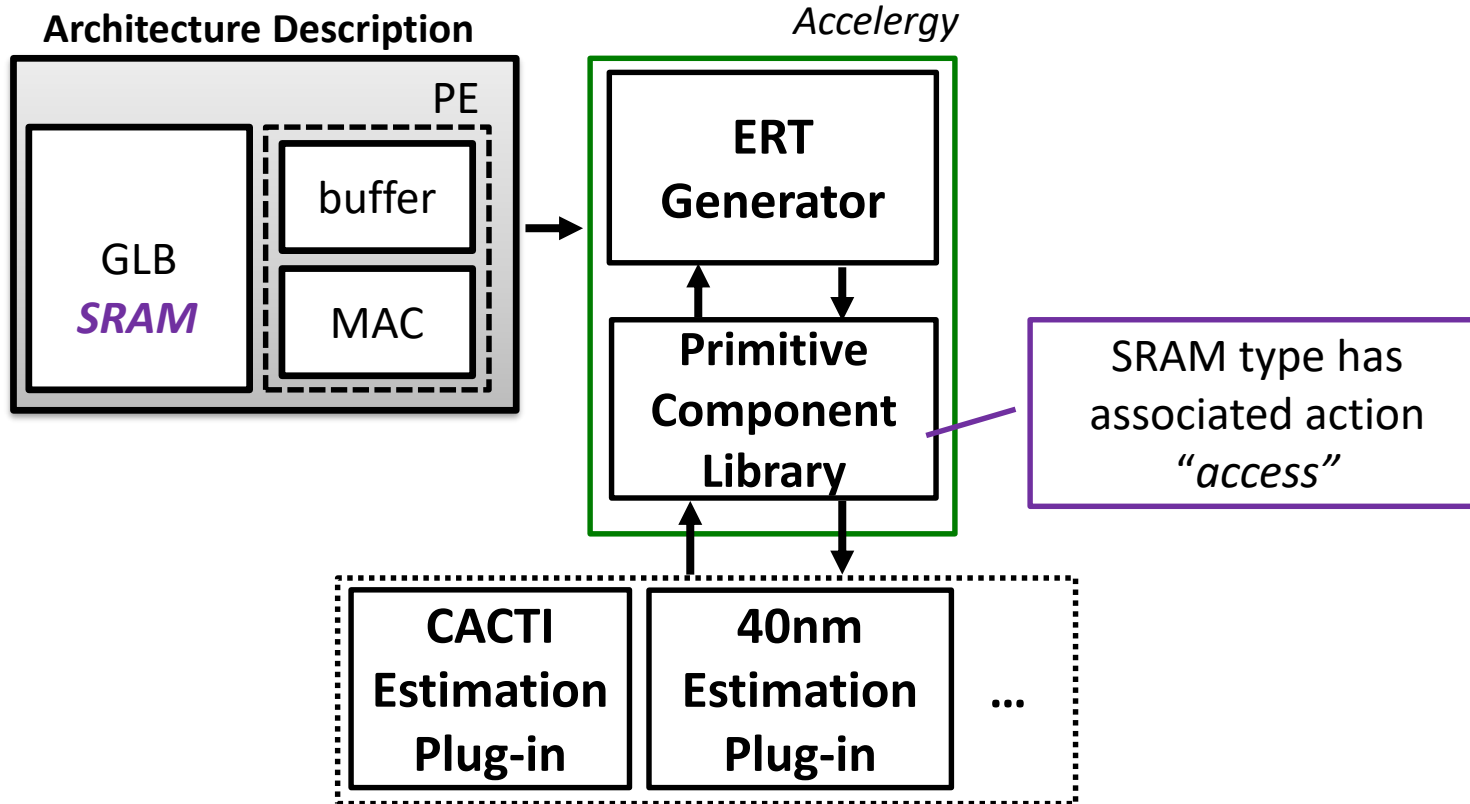


# Accelergy Overview

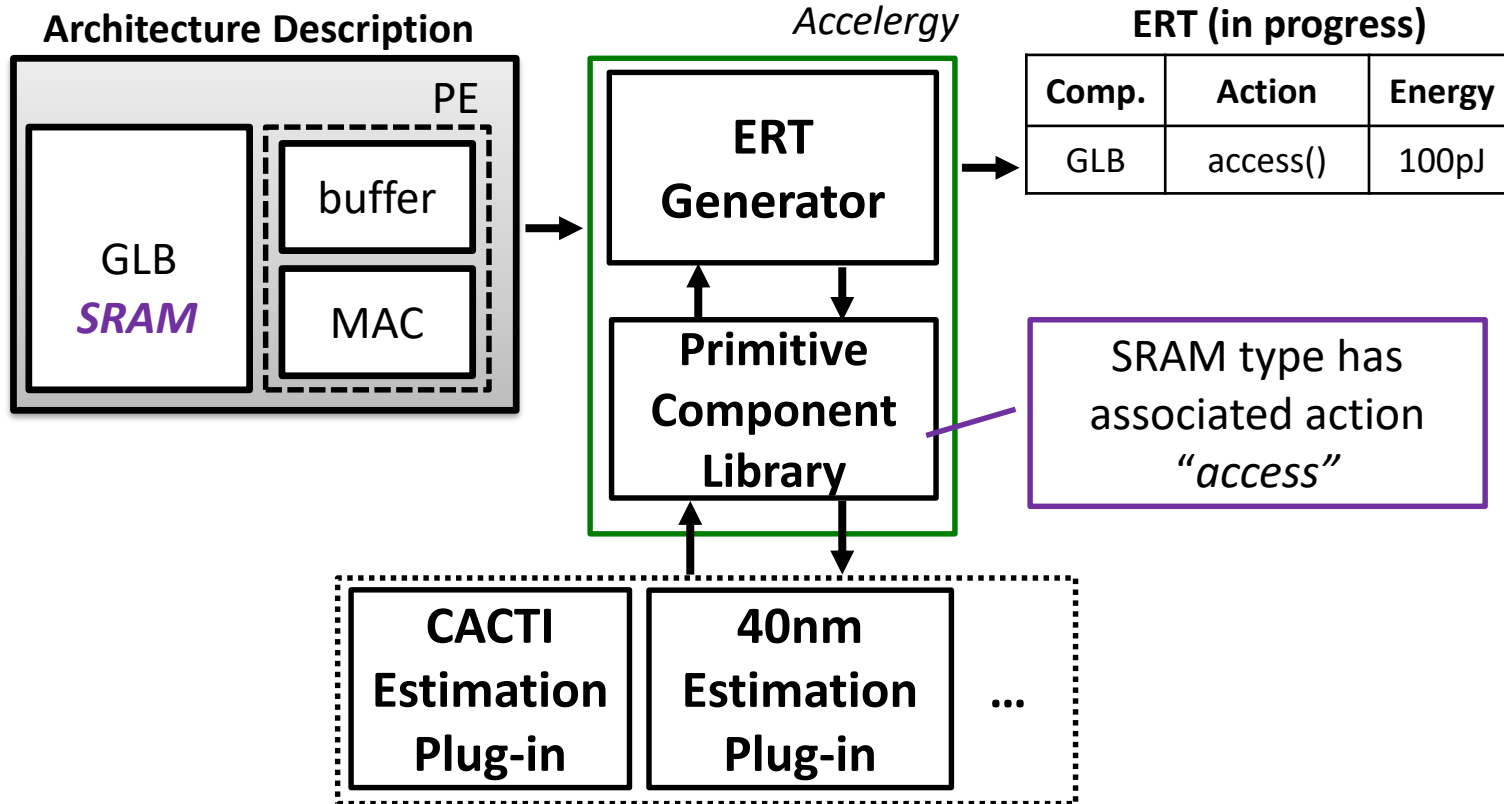
---

- An architecture-level energy estimator
- **Flexibly characterizes various primitive components of different technologies**
- Succinctly models diverse and complicated designs
- Improves estimation accuracy via fine-grained classification of operations
- Achieves 95% accuracy in evaluating a deep neural network (DNN) accelerator – Eyeriss [ISSCC 2016]

# Accelergy: Flexibly Model Various Primitive Components

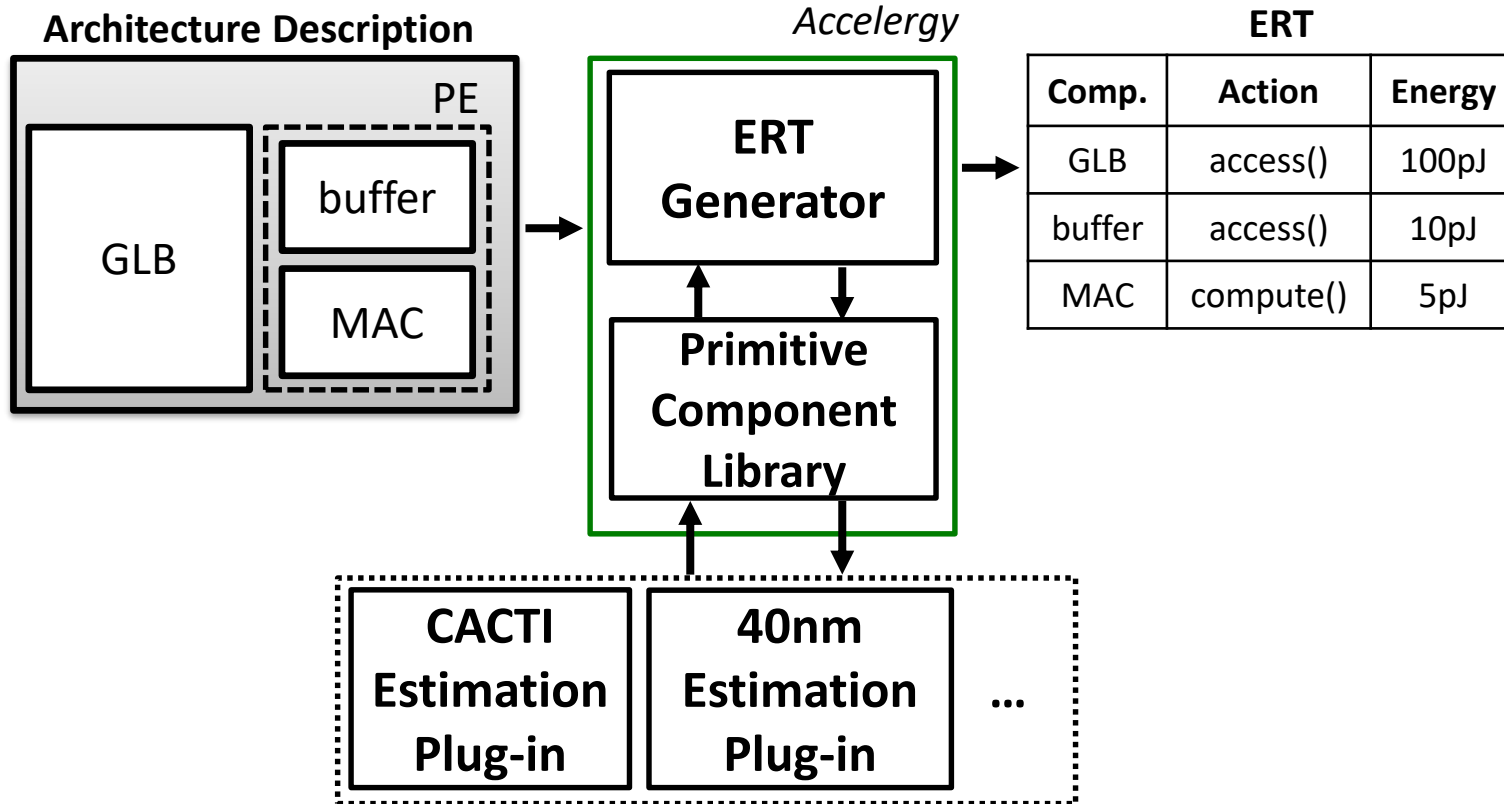


# Accelergy: Flexibly Model Various Primitive Components

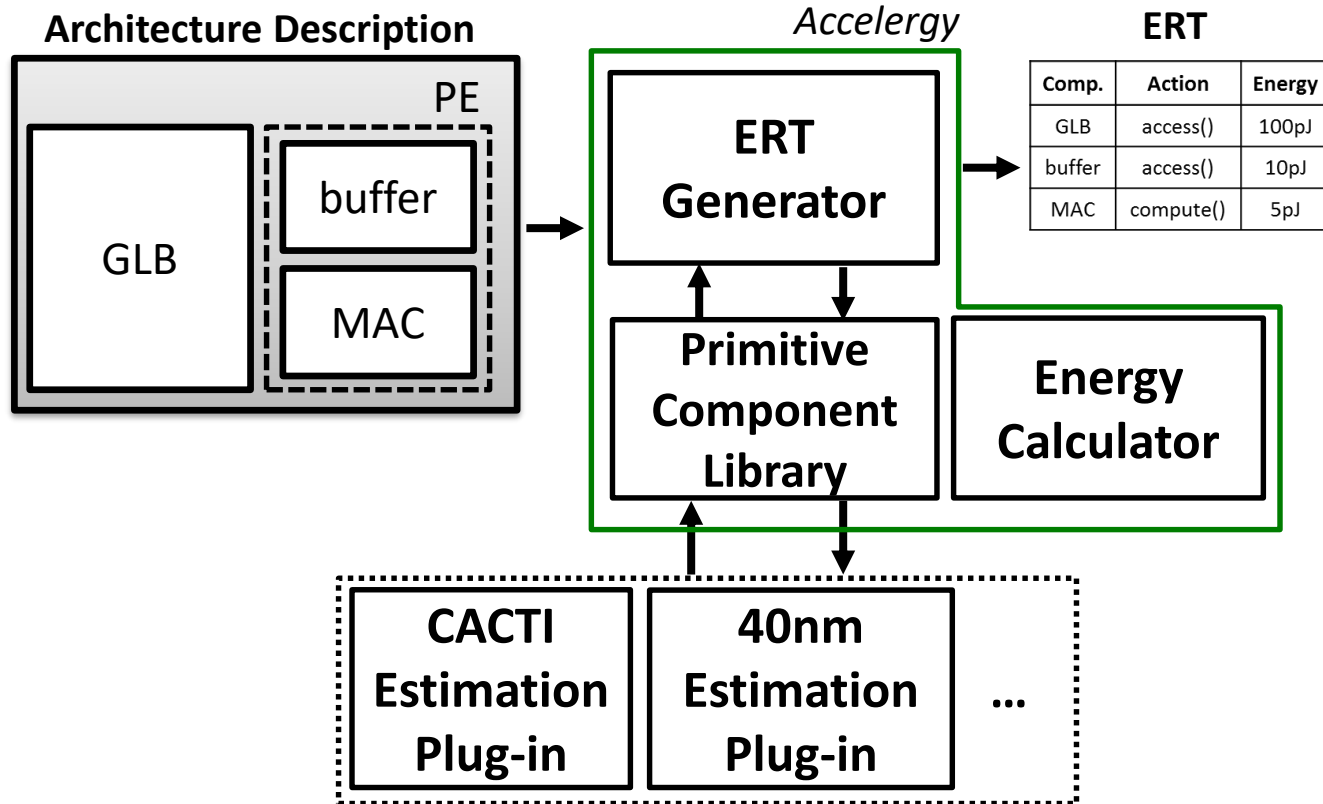




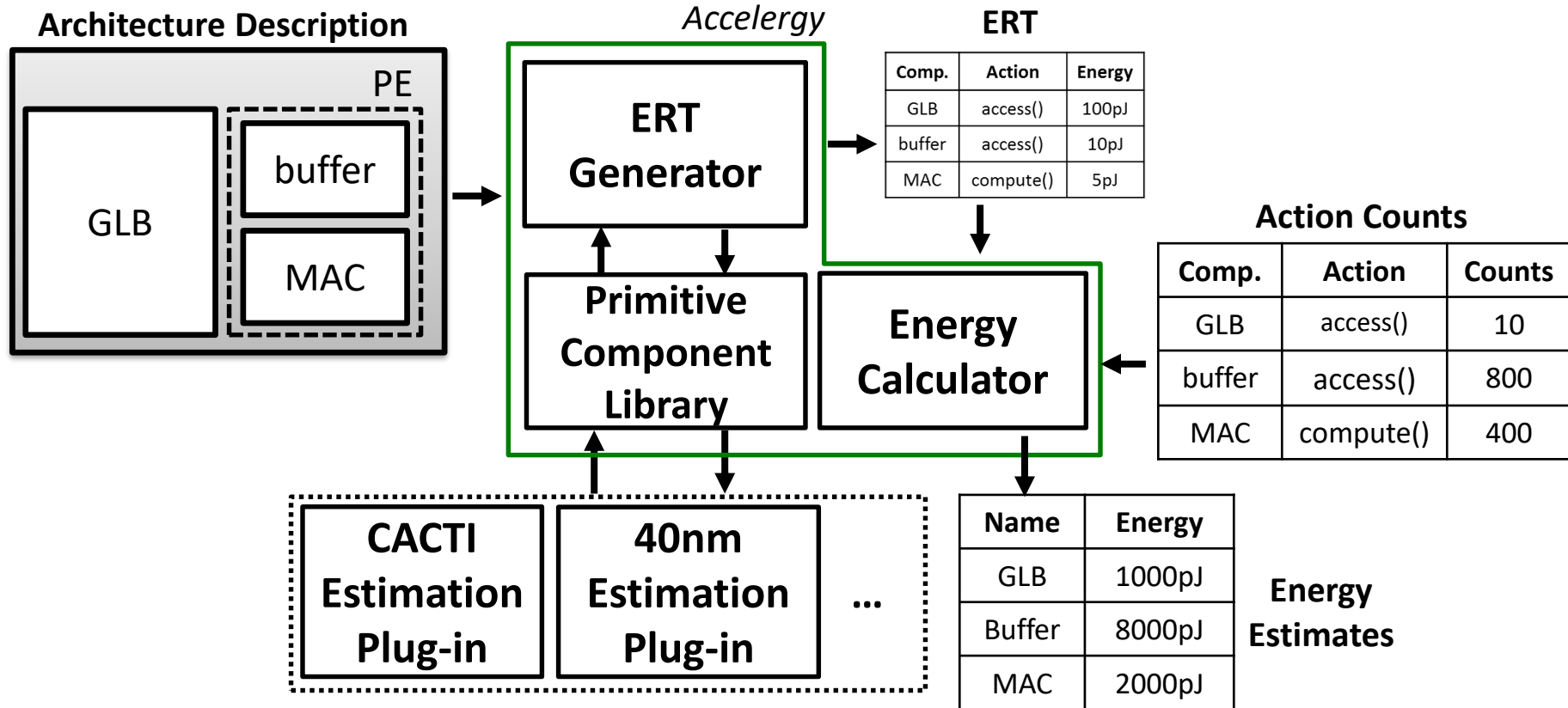
# Accelergy: Flexibly Model Various Primitive Components



# Accelergy: Flexibly Model Various Primitive Components



# Accelergy: Flexibly Model Various Primitive Components



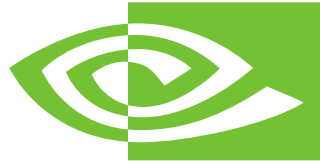
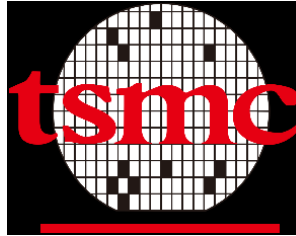
# Accelergy: Flexibly Model Various Primitive Components

Use energy estimation plug-ins to characterize primitive components

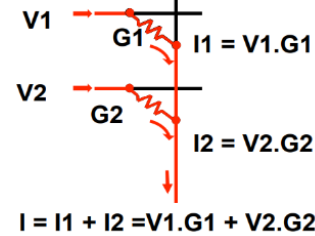
CACTI  
Estimation  
Plug-in

40nm  
Estimation  
Plug-in

Traditional open-source  
plug-ins\*



Proprietary plug-ins



NVSIM  
[TCAD 2012]

Emerging technology  
plug-ins

\*available at <http://accelergy.mit.edu>

# Accelergy: Flexibly Model Various Primitive Components

Use energy estimation plug-ins to characterize primitive components

CACTI  
Estimation  
Plug-in



**Detailed plug-in interface in open-source repo**

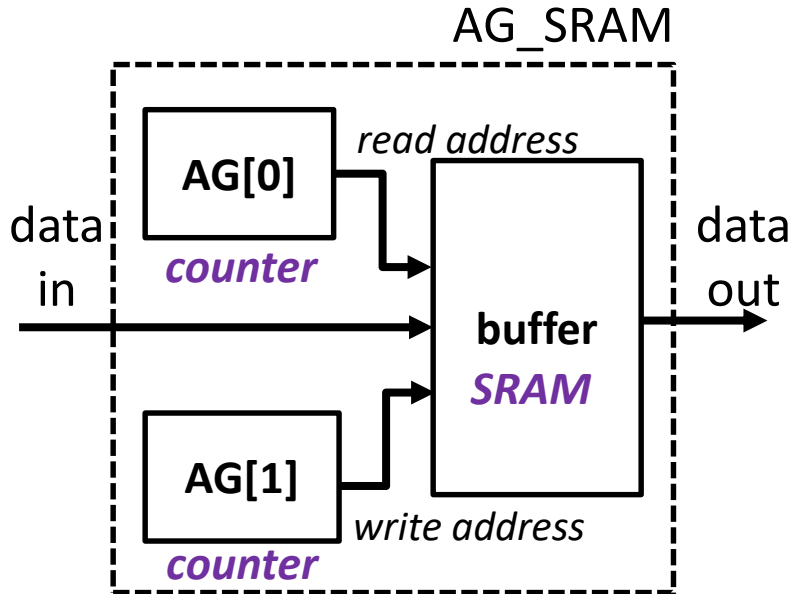
Traditional open source  
plug-ins\*

Accelergy  
plug-ins

\*available at <http://accelergy.mit.edu>

# Modeling Complicated Designs

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

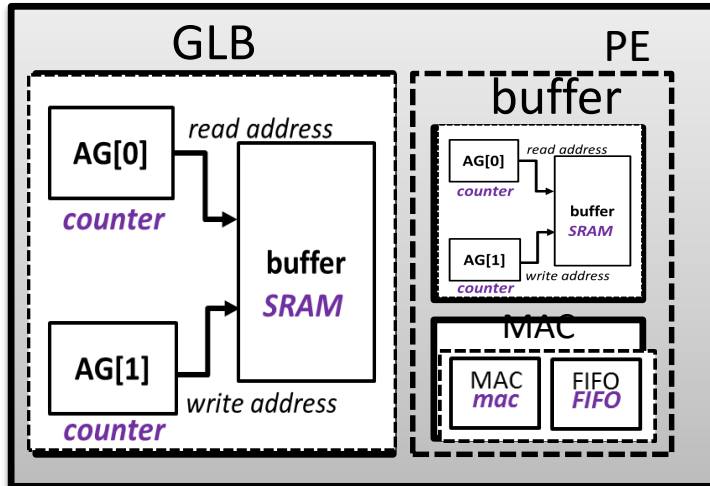


- AG\_SRAM is an abstract hierarchy
- Buffer is of **SRAM** type
- AGs is of **counter** type

# Modeling Complicated Designs

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

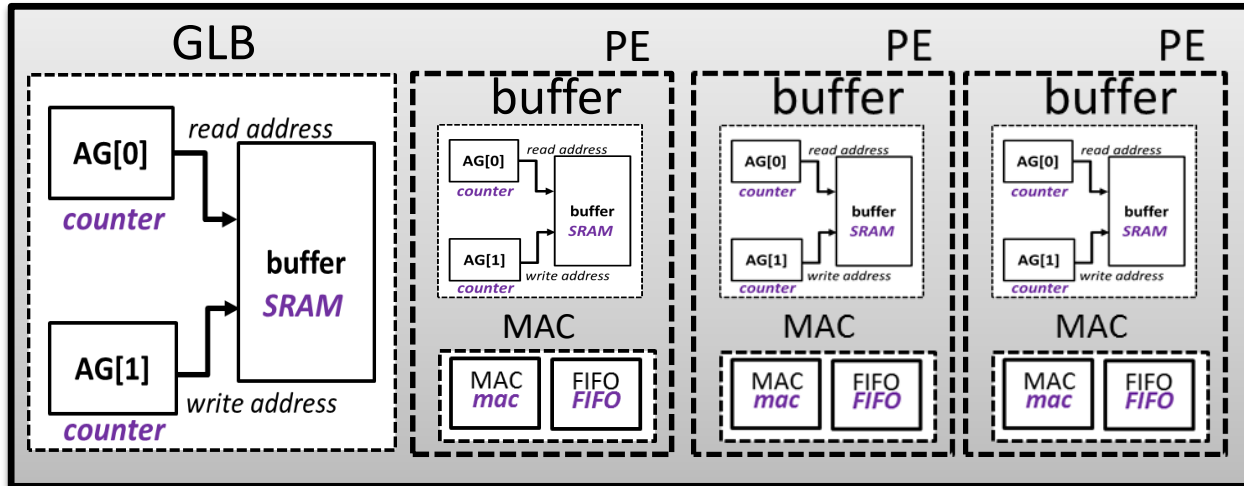
*Let's construct a more practical design!*



# Modeling Complicated Designs

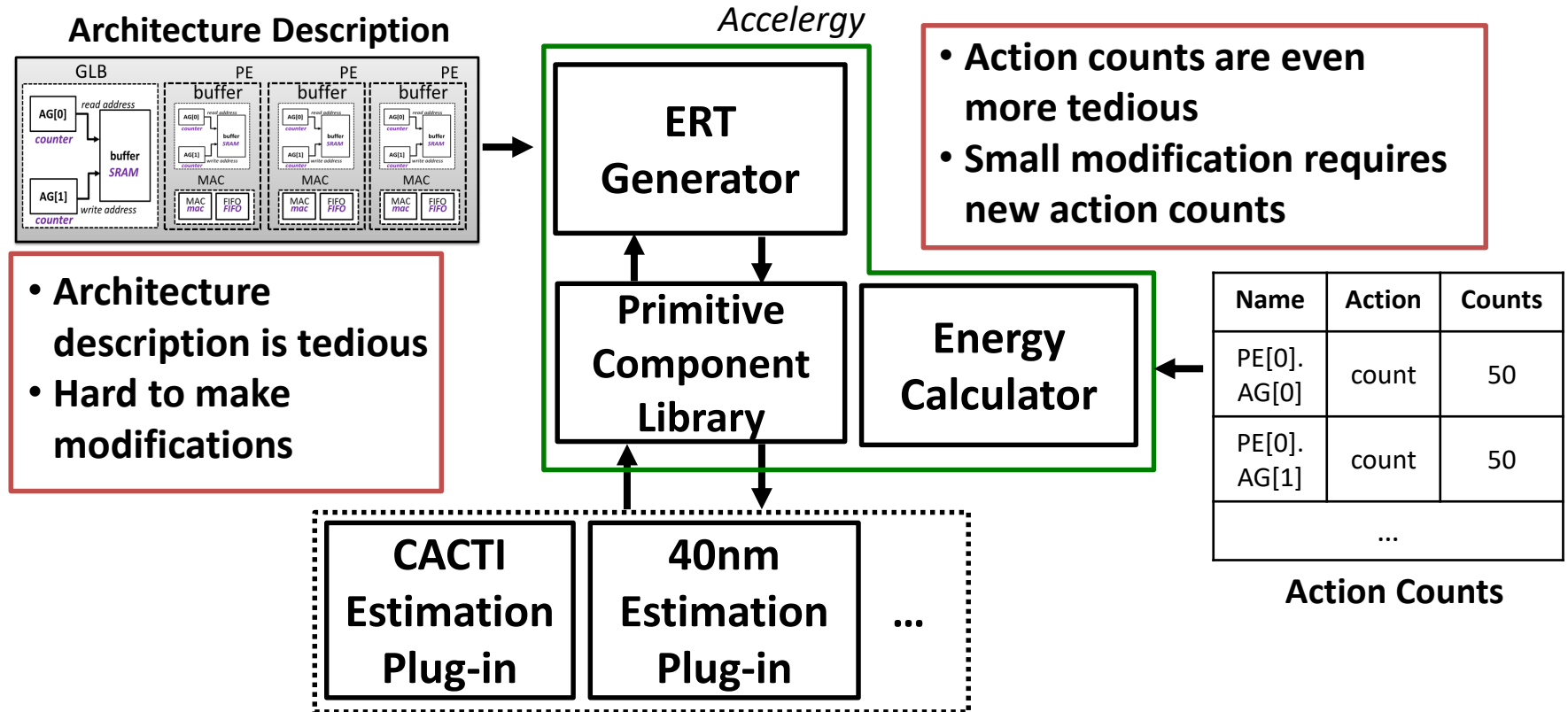
- Practical architecture designs involve much more details
  - Example: storage units with local address generators (AG)

*Let's construct a more practical design!*





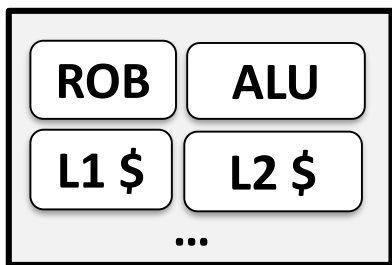
# Modeling Complicated Designs



# Existing Work - Modeling Complicated Designs

- Existing work that aims to succinctly model complicated architectures
  - **Wattch**[ISCA2000], **McPAT**[MICRO2009]

## CPU-Centric Architecture Model



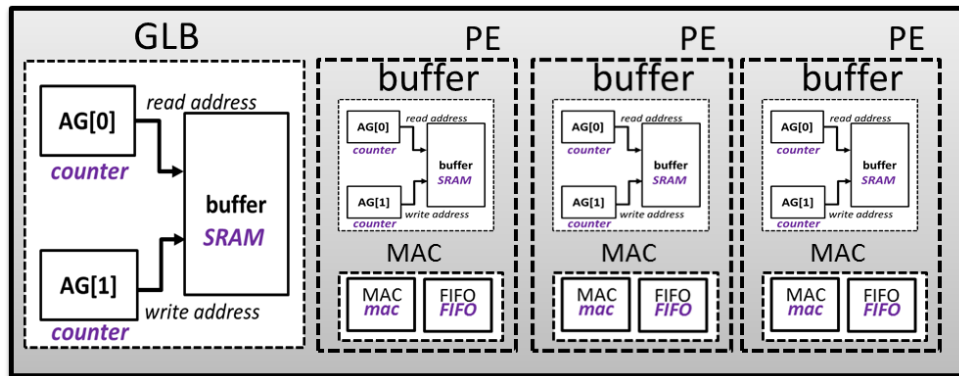
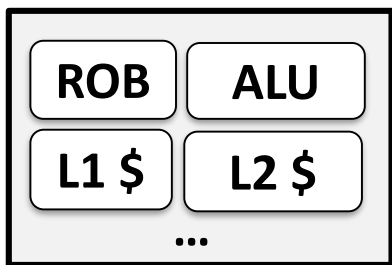
Use a fixed set of **compound components** to  
represent the architecture

*Components that can  
be decomposed into  
lower level components*

# Existing Work - Modeling Complicated Designs

- Existing work that aims to succinctly model complicated architectures
  - **Wattch**[ISCA2000], **McPAT**[MICRO2009]

## CPU-Centric Architecture Model



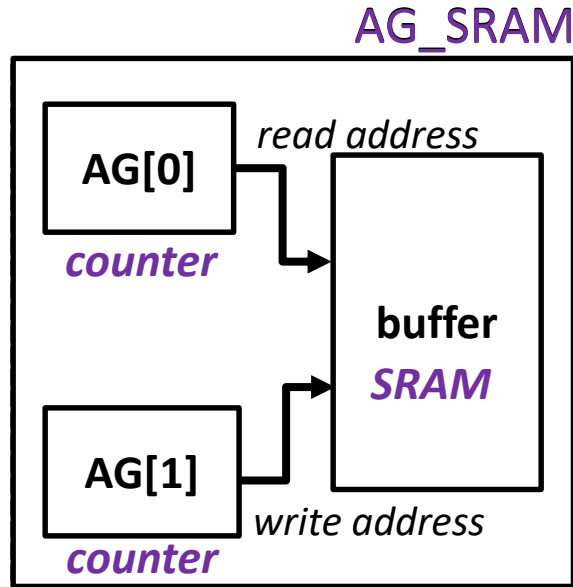
The fixed set of compound components is not sufficient to describe arbitrary accelerator designs

# Accelergy Overview

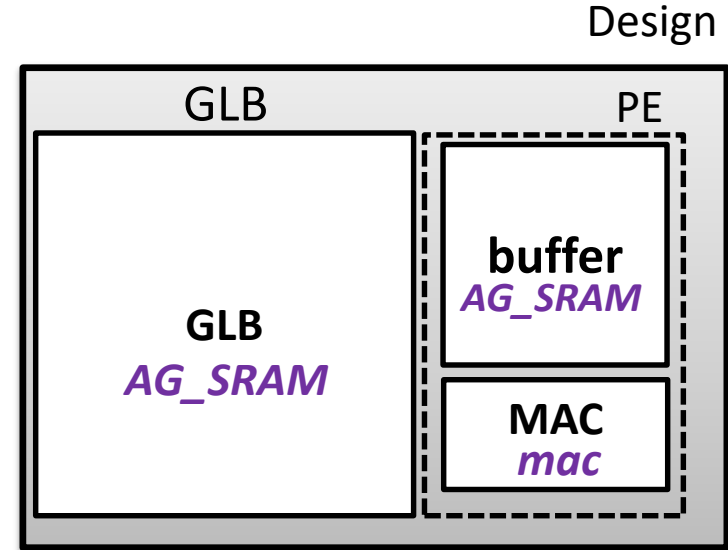
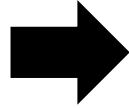
---

- An architecture-level energy estimator
- Flexibly characterizes various primitive components of different technologies
- **Succinctly models diverse and complicated designs**
- Improves estimation accuracy via fine-grained classification of operations
- Achieves 95% accuracy in evaluating a deep neural network (DNN) accelerator – Eyeriss [ISSCC 2016]

# Accelergy: Succinctly Model Arbitrary Architecture

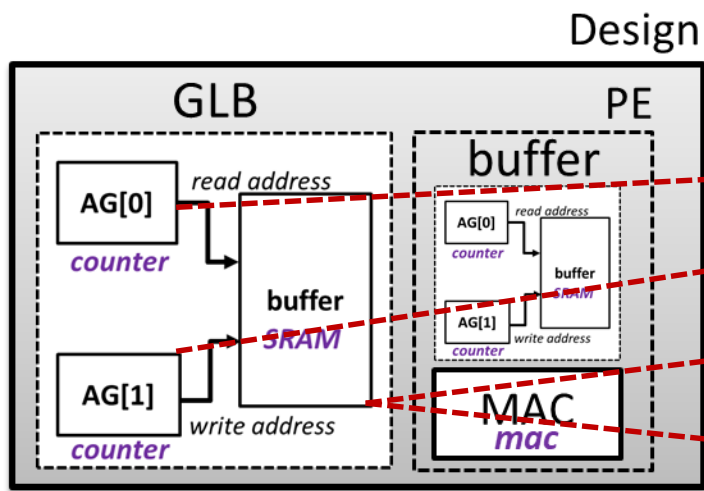


AG\_SRAM is an  
user-defined compound  
component



Architecture described with  
compound components and  
primitive components

# Accelergy: Succinctly Model Arbitrary Action Counts

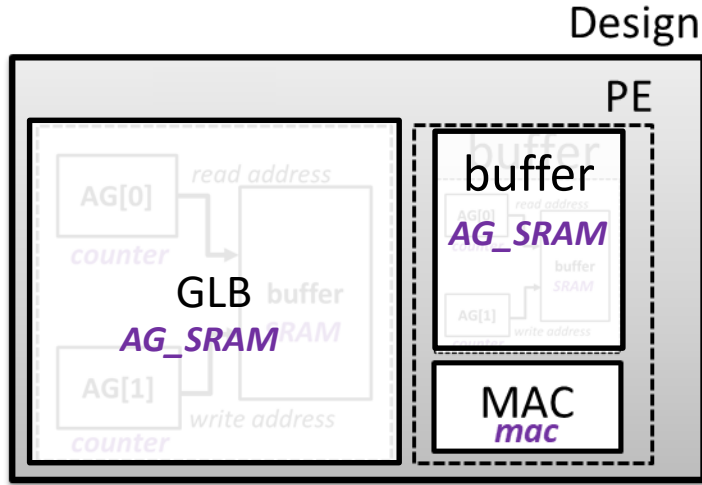


Action Counts

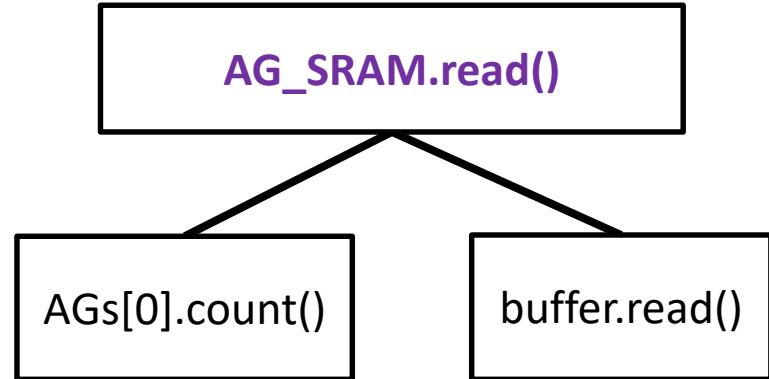
Name	Action	Counts
GLB.AG[0]	count()	50
GLB.AG[1]	count()	20
GLB.buffer	read()	50
GLB.buffer	write()	20
...		

**Tedious action counts in terms of  
primitive component actions**

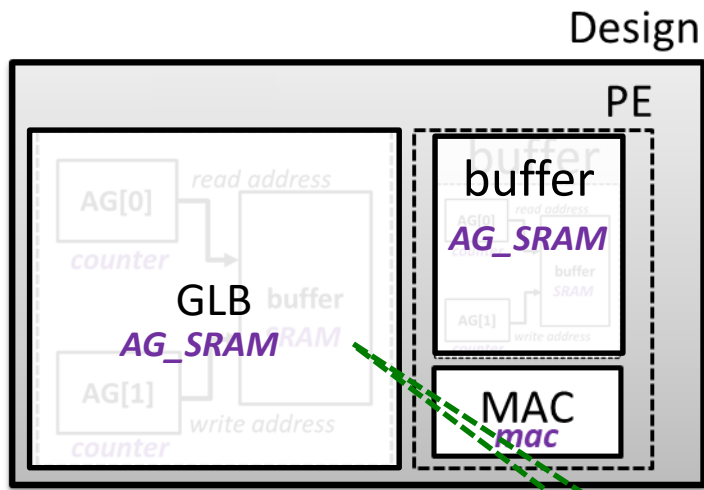
# Accelergy: Succinctly Model Arbitrary Action Counts



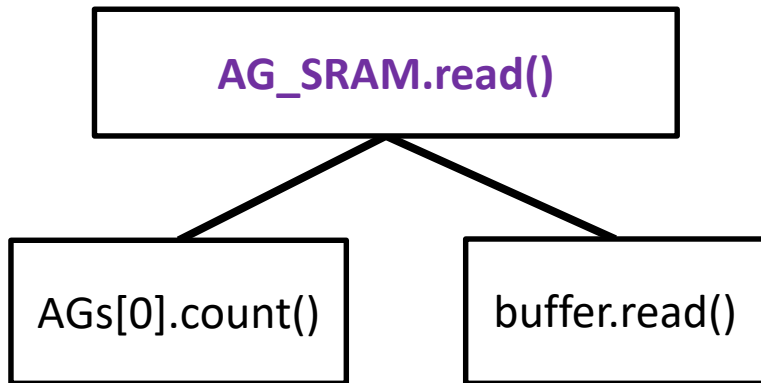
## User-defined compound actions



# Accelergy: Succinctly Model Arbitrary Action Counts



## User-defined compound actions



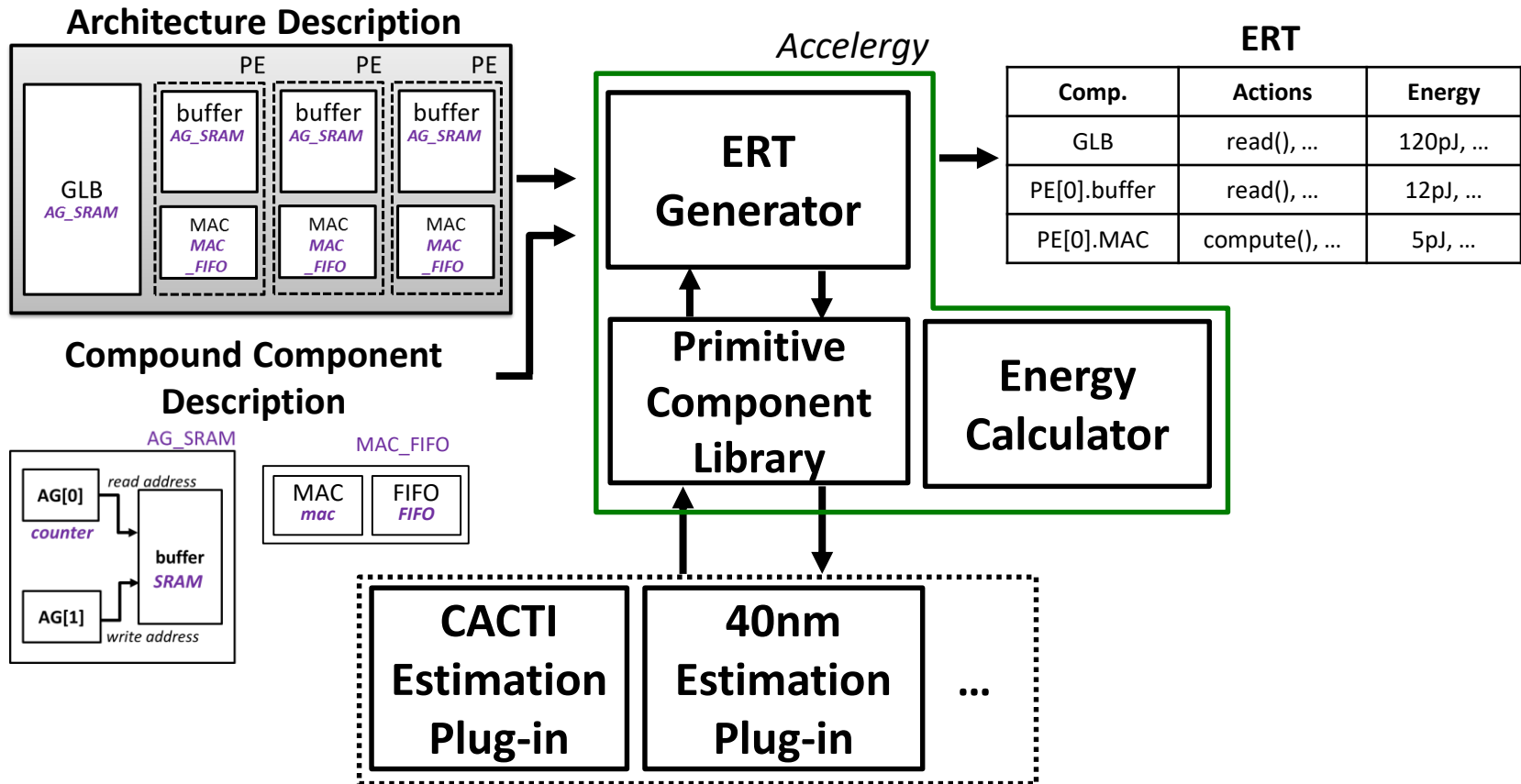
## Action Counts

Name	Action	Counts
GLB	read()	50
GLB	write()	20

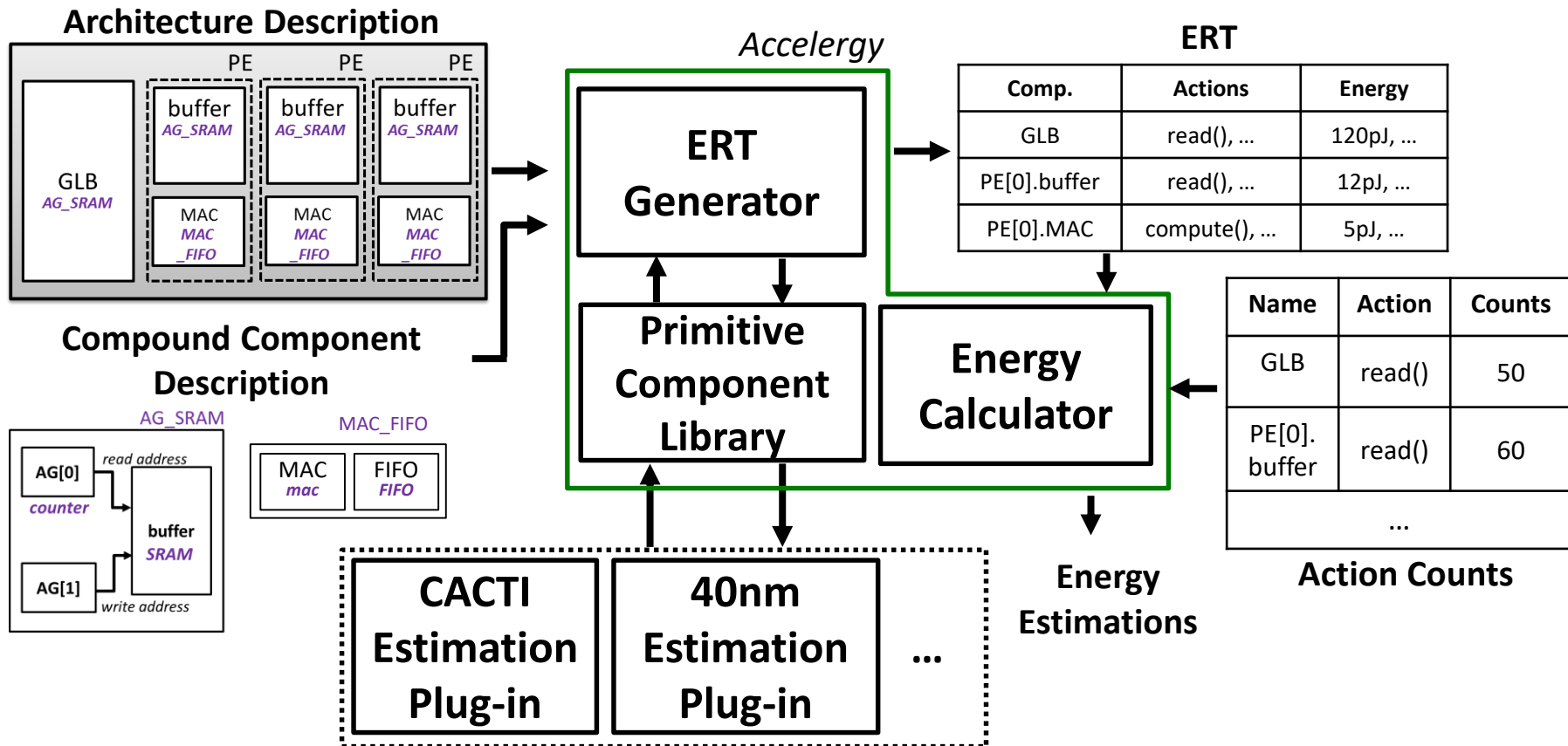
**Succinct action counts with compound component actions**



# Accelergy: Succinctly Model Complex Designs



# Accelergy: Succinctly Model Complex Designs



# Additional Challenge: Inaccurate Modeling of Energy/Action

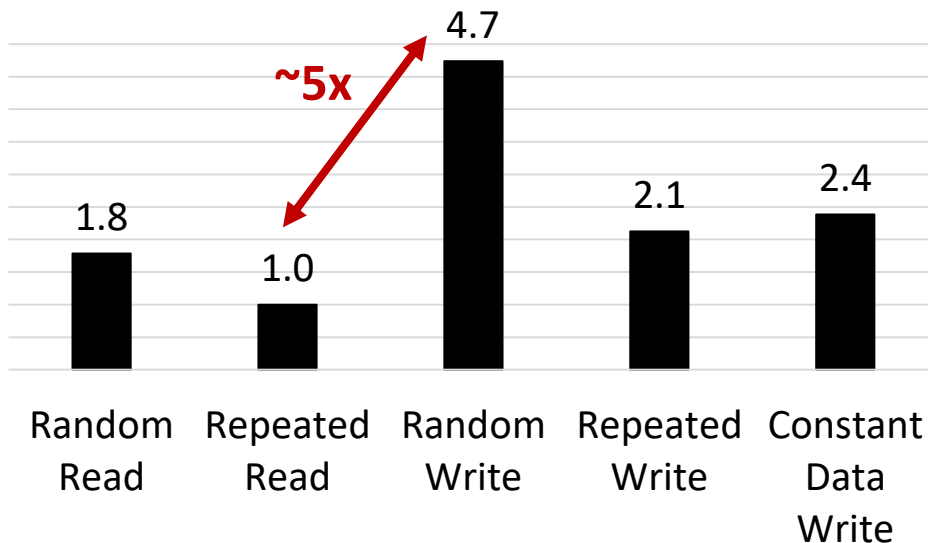
- Existing architecture-level energy estimators only model coarse action types

Component	Action	Energy
GLB	access()	100pJ
Buffer	access()	10pJ
ALU	compute()	5pJ

**Coarse-grained Actions**

**Coarse-grained estimations  
reduce estimation accuracies**

**Energy-Per-Actions of a Register File  
(normalized to idle)**



**Fine-grained Actions**

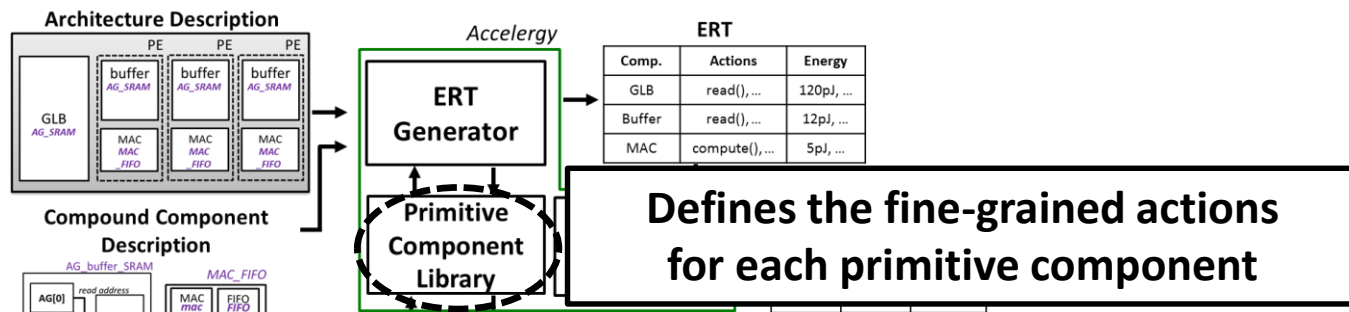
# Accelergy Overview

---

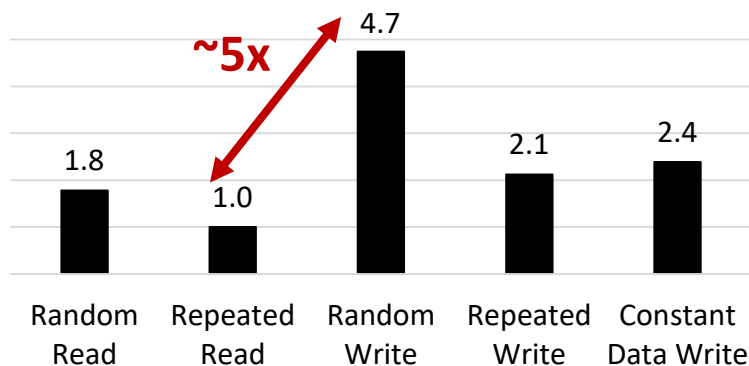
- An architecture-level energy estimator
- Flexibly characterizes various primitive components of different technologies
- Succinctly models diverse and complicated designs
- **Improves estimation accuracy via fine-grained actions**
- Achieves 95% accuracy in evaluating a deep neural network (DNN) accelerator – Eyeriss [ISSCC 2016]

# Accelergy: Fine-grained Action Classification

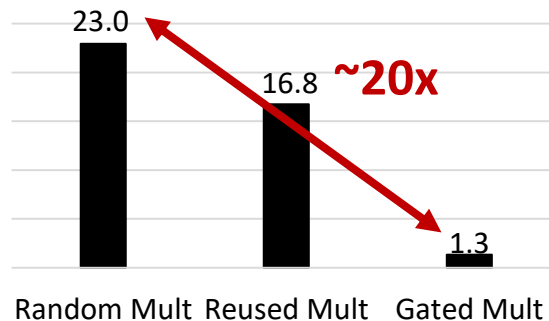
- Accurate estimation with a primitive component library



Fine-grained memory action types

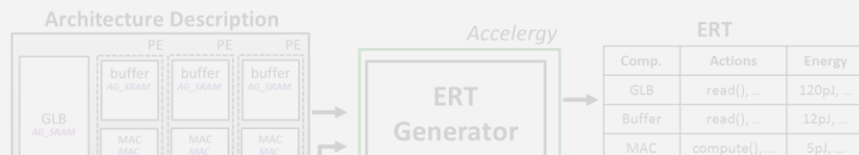


Fine-grained multiplier action types



# Accelergy: Fine-grained Action Classification

- Accurate estimation with a primitive component library\*



**Detailed methodology for generating fine-grained action types in paper**



# Accelergy Overview

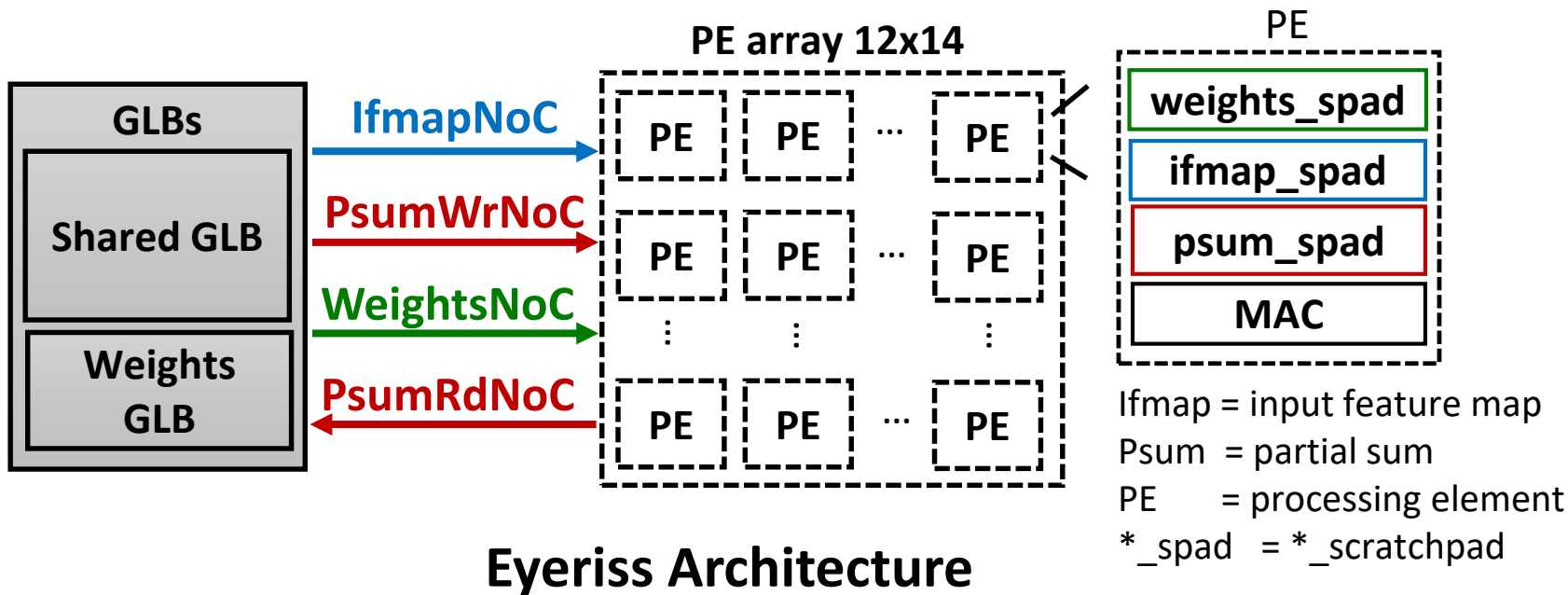
---

- An architecture-level energy estimator
- Flexibly characterizes various primitive components of different technologies
- Succinctly models diverse and complicated designs
- Improves estimation accuracy via fine-grained actions
- **Achieves 95% accuracy in evaluating a deep neural network (DNN) accelerator – Eyeriss [ISSCC 2016]**

# Energy Evaluations on Eyeriss

- Experimental Setup:

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





# Energy Evaluations on Eyeriss

- **Experimental Setup:**

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

## Zero-gating optimization

If there is a 0 ifmap data

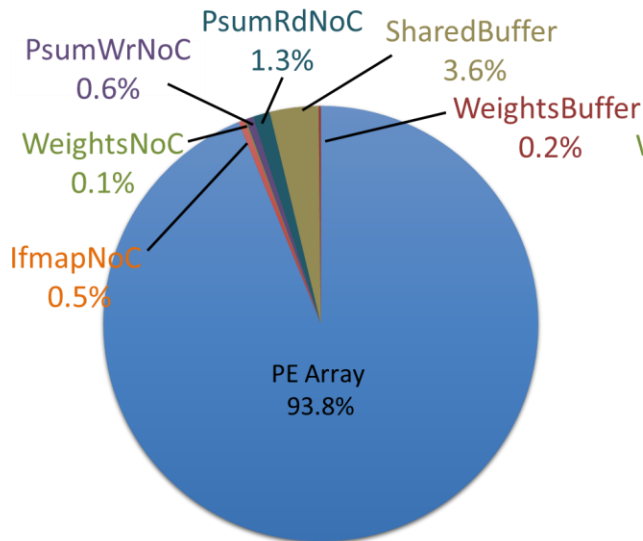
- Gate on reading the weights data => gated-read
- Gate on computing the MAC => gated-MAC

Eyeriss Architecture

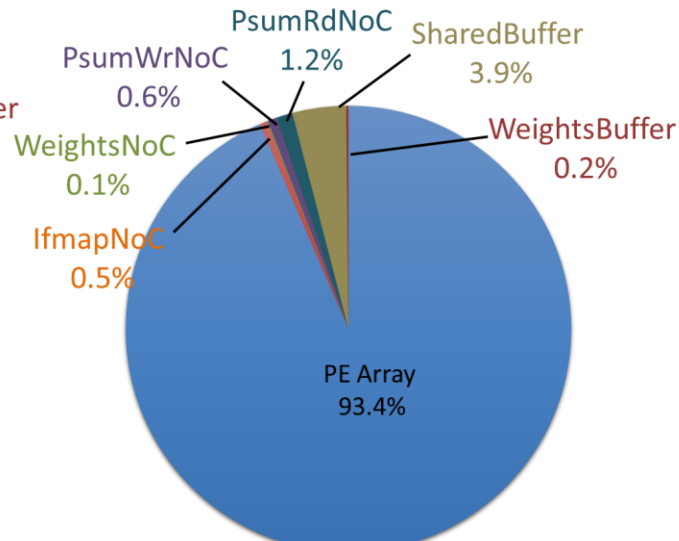
PE = processing element  
\*\_spad = \*\_scratchpad

# Total Energy and Coarse Energy Breakdown

- 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.



Ground Truth Energy Breakdown

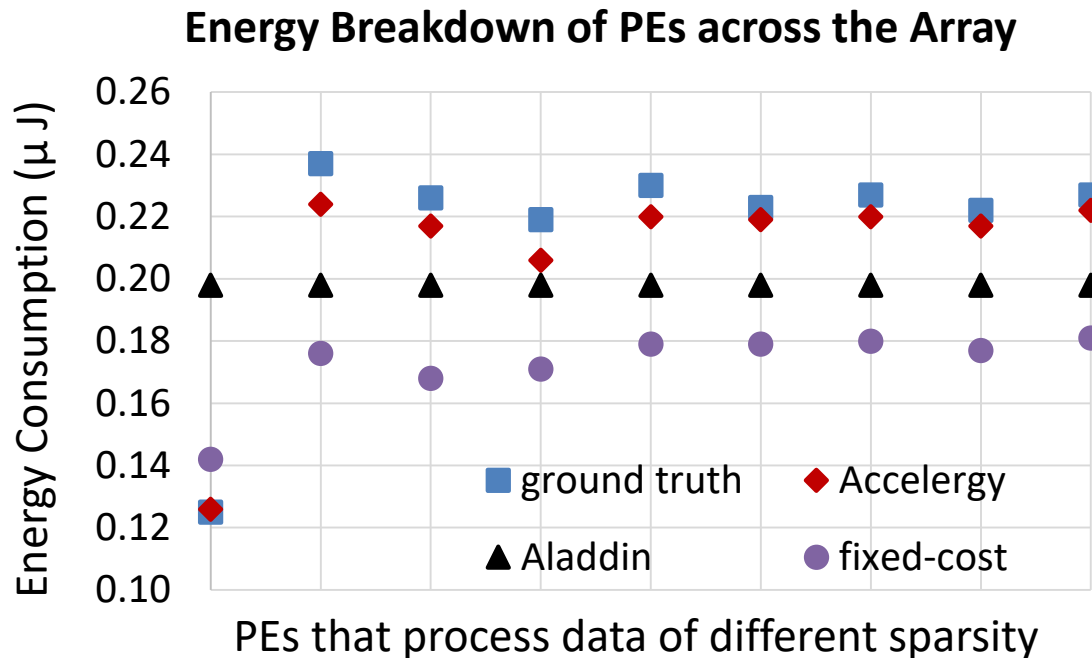


Accelergy Energy Breakdown

*\*Total energy might not add up to exact 100.0% due to rounding*

# PE Array Energy Breakdown

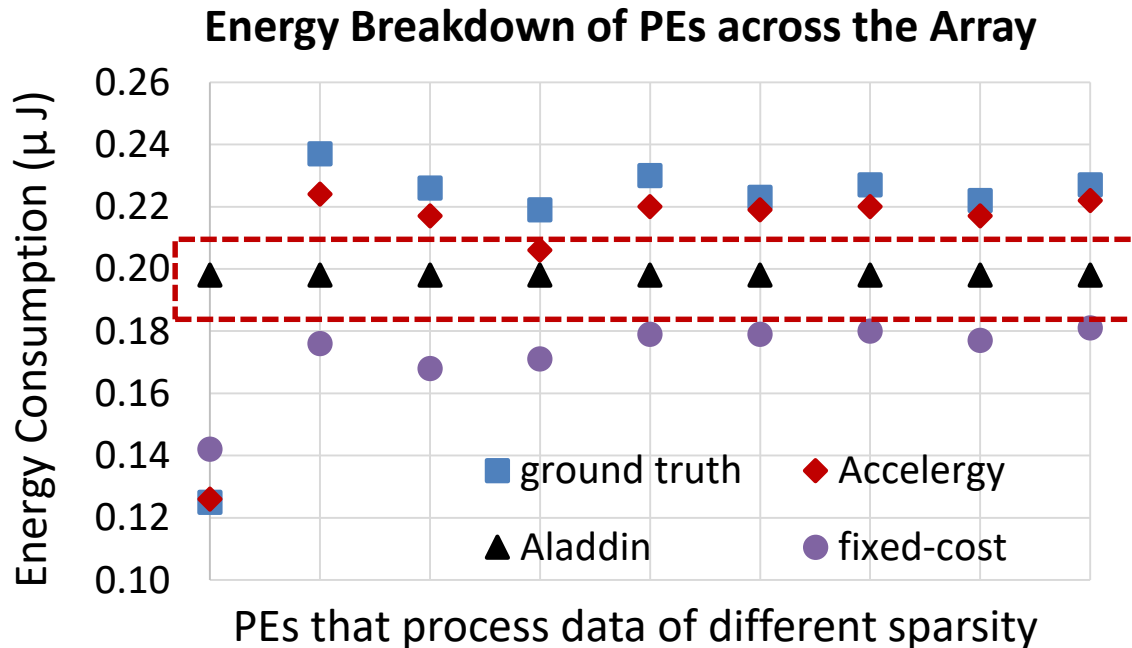
- Comparisons with existing work: Aladdin and fixed-cost



# PE Array Energy Breakdown

- Comparisons with existing work: Aladdin and fixed-cost

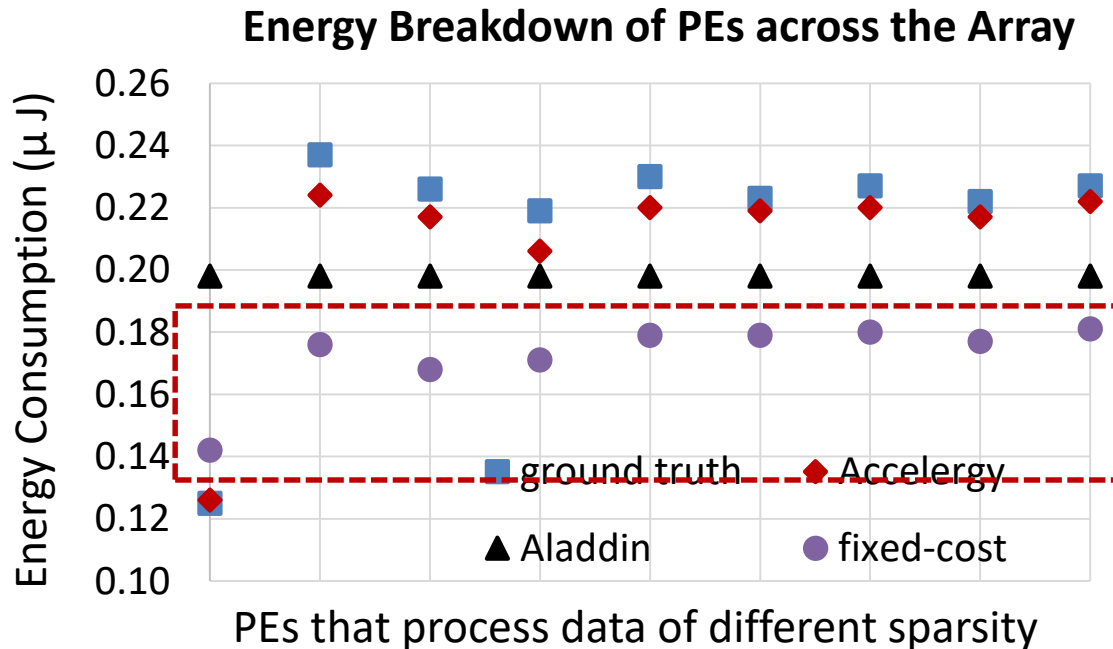
*Not aware of the fine-grained actions related to zero-gating optimization*



# PE Array Energy Breakdown

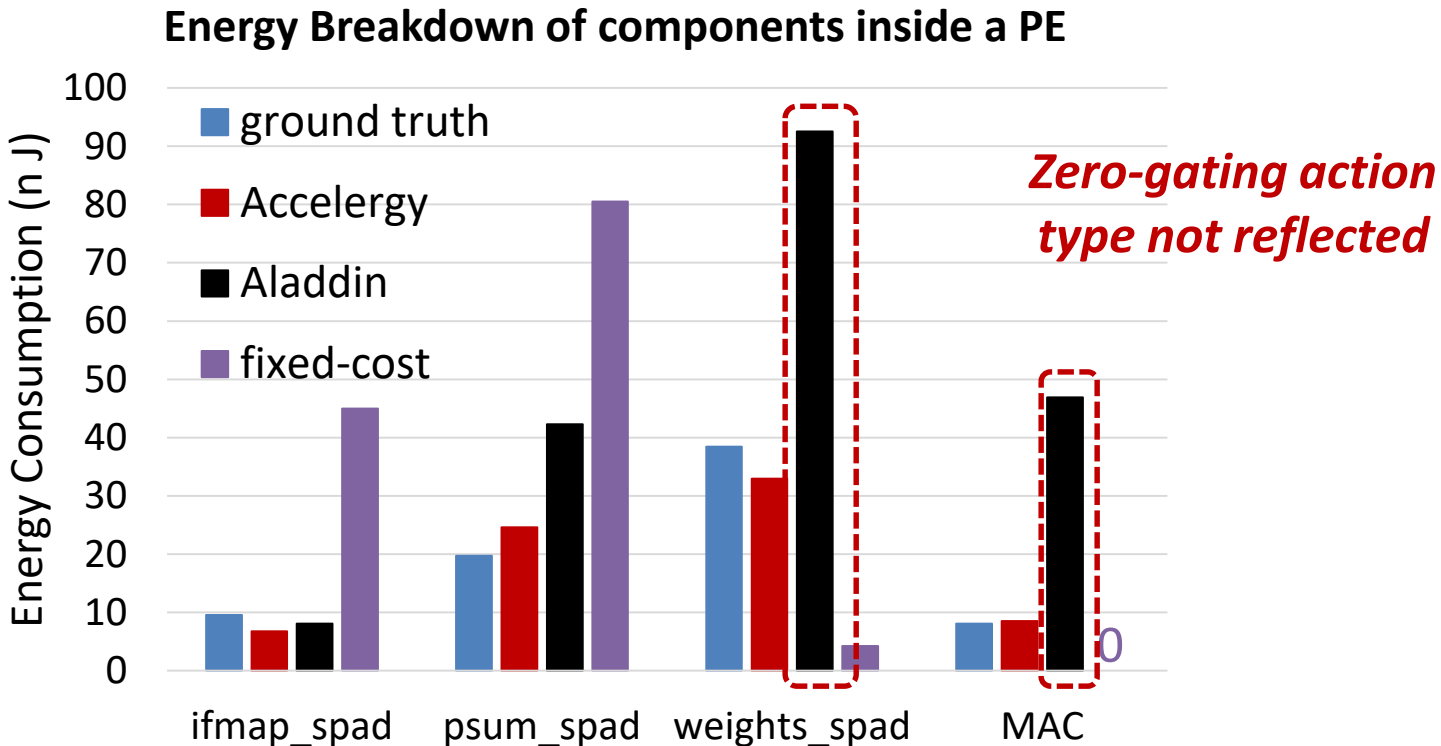
- Comparisons with existing work: Aladdin and fixed-cost

*Inaccurate energy  
characterization of  
components*



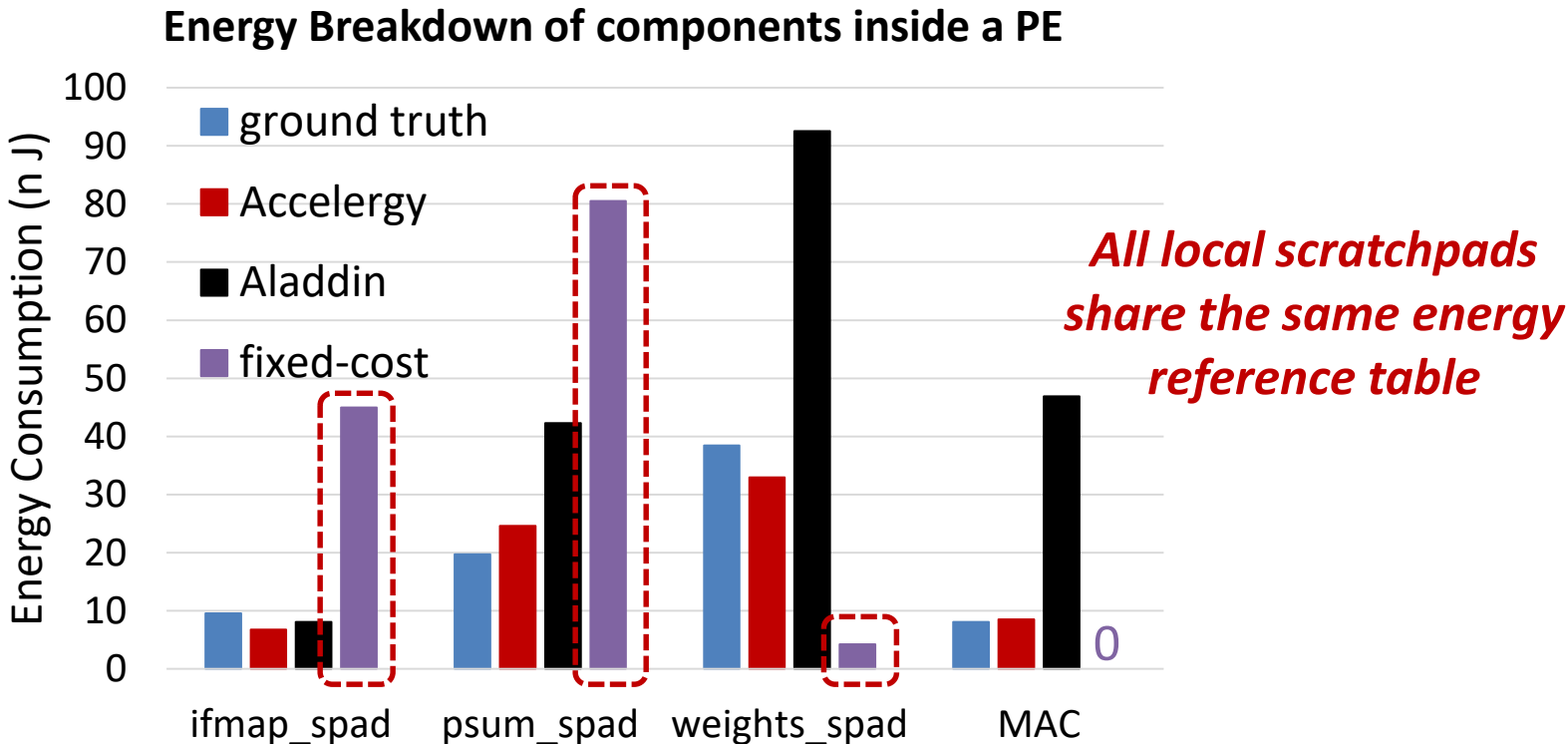
# PE Energy Breakdown

- Comparisons with existing work: Aladdin and fixed-cost



# PE Energy Breakdown

- Comparisons with existing work: Aladdin and fixed-cost



# Conclusion

---

- **Accelergy is an architecture-level energy estimator that**
  - **Accelerates accelerator design space exploration**
  - **Provides flexibility to**
    - Describe a diverse range of accelerator designs
    - Support estimation of different technologies, e.g., CMOS, RRAM, optical
  - **Achieves high accuracy energy estimations**
    - 95% accurate for the Eyeriss accelerator
- **Open-source code available at: <http://accelergy.mit.edu>**