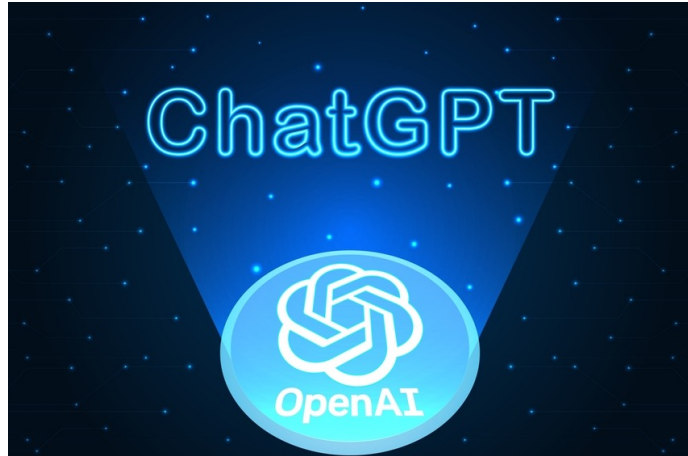# HighLight:
# Efficient and Flexible DNN Acceleration with Hierarchical Structured Sparsity

Yannan Nellie Wu[1], **Po-An Tsai**[2], Saurav Muralidharan[2], Angshuman Parashar[2], Vivienne Sze[1], Joel S. Emer[1,2]

[1]MIT, [2]NVIDIA

http://emze.csail.mit.edu/highlight

# Many Applications Involve DNNs



**Natural Language Processing**
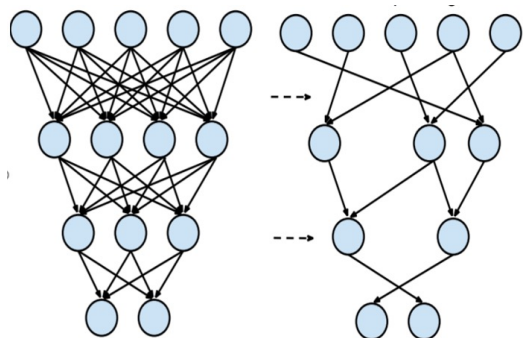


**Autonomous Navigation**



**Medical Imaging & Diagnostics**

data and computation intensive
subject to prediction accuracy & latency requirements

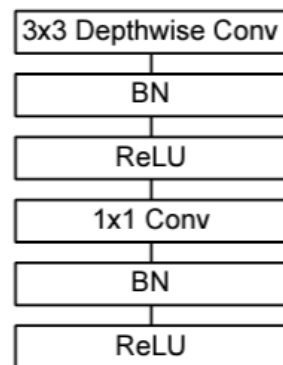**In great need of optimizations and accelerations**

# Different DNN Optimizations Introduce Different Sparsity

## Optimizations to Reduce Model Size



3x3 Depthwise Conv
BN
ReLU
1x1 Conv
BN
ReLU

Pruning
Techniques
*[Han, NeurIPS15]*

Depth-wise
Separable Layers
*[Howard, CVPR17]*

**Introduces
Sparse Weights**

**Introduces
Dense Weights**

# Different DNN Optimizations Introduce Different Sparsity

## Optimizations to Reduce Model Size



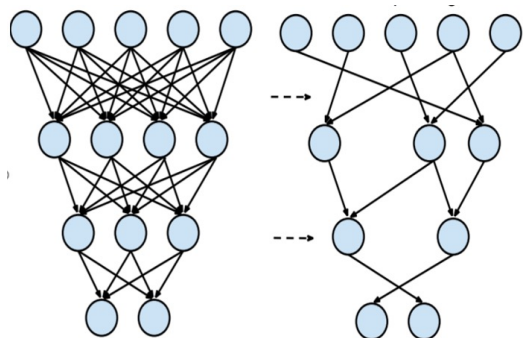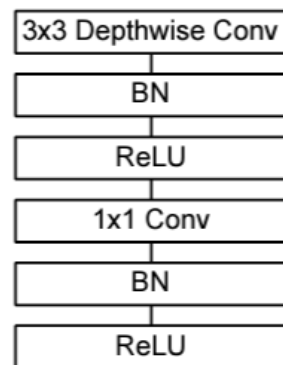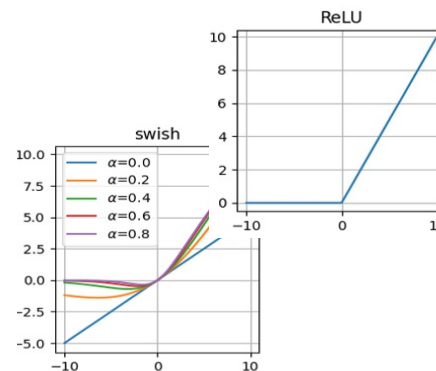### Pruning Techniques
*[Han, NeurIPS15]*

**Introduces Sparse Weights**



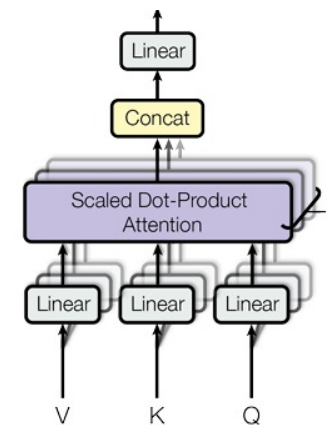### Depth-wise Separable Layers
*[Howard, CVPR17]*

**Introduces Dense Weights**

## Optimizations to Improve Accuracy



### Activation Functions
*[Apicella, NN21]*
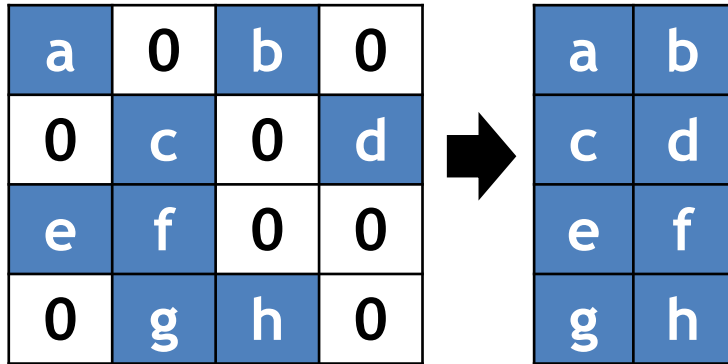
**Introduces Dense/Sparse Activations**



### Attention-based Modules
*[Vaswani, NeurIPS17]*

**Introduces Dense Act./Weights**

---

**Modern DNNs can weights and activations that are either dense or sparse with various sparsity degrees**

# High-Level Opportunities for Sparse DNNs

| a | 0 | b | 0 |
|---|---|---|---|
| 0 | c | 0 | d |
| e | f | 0 | 0 |
| 0 | g | h | 0 |

➡

| a | b |
|---|---|
| c | d |
| e | f |
| g | h |

$$x \times 0 = 0$$
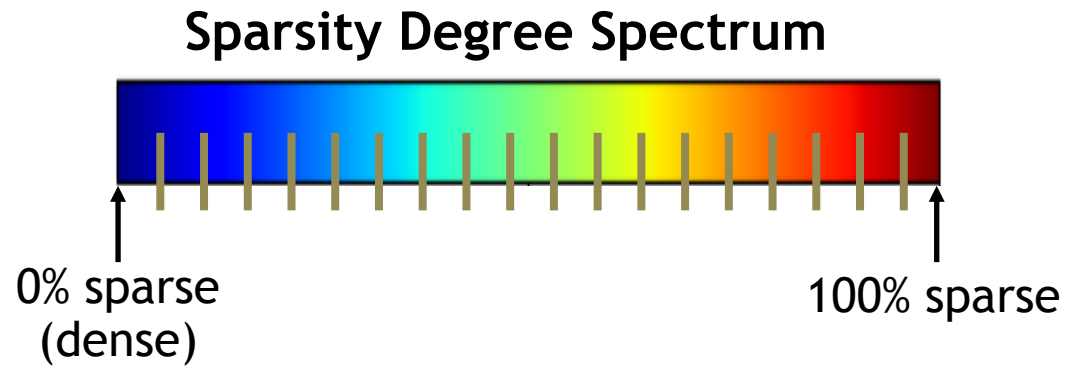$$x + 0 = x$$

**Zero Values
Can be Compressed Away**

**Ineffectual Operations
Can be Eliminated**

**Important to design sparse DNN accelerators
to exploit such opportunities**

# Requirements for an Ideal Sparse DNN Accelerator

**Sparsity Degree Spectrum**

0% sparse
(dense)

100% sparse

## Flexible

*exploit many sparsity degrees*

# Requirements for an Ideal Sparse DNN Accelerator

**Sparsity Degree Spectrum**

0% sparse
(dense)

100% sparse
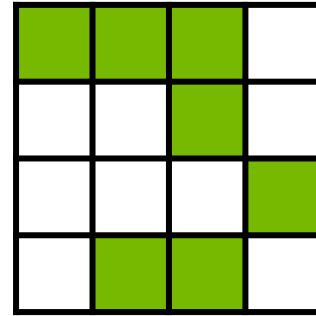
## Flexible

*exploit many sparsity degrees*

## Efficient

*low __sparsity tax__ for eliminating ineffectual operations*

# Existing Works Do Not Meet Such Requirements

## Unstructured Sparse Accelerators



*Unstructured sparse*

**Dual-Side Sparse Tensor Core (DSTC)**
*[Wang, ISCA21]*

Sparsity Degree Spectrum



Continuously Translated into Savings

− **High Sparsity Tax**
+ **Flexible**

## Structured Sparse Accelerators



*Per-row
2:4 structured sparse
(G:H pattern)*

**NVIDIA Sparse Tensor Core (STC)**
*[NVIDIA, TechReport20]*

Sparsity Degree Spectrum



0% sparse
(dense)

50% sparse (2:4)

+ **Low Sparsity Tax**
− **Inflexible**

# Naïve Way to Increase Flexibility Structured Sparse Designs

## Extend the Number of G:H Ratios Supported

Sparsity Degree Spectrum



50% sparse
**2 : 4**

0% sparse
(dense)

67% sparse
**2 : 6**

75% sparse
**2 : 8**

**Not Scalable**
Sparsity tax increases approximately in proportion to the number of sparsity degrees

# Our Proposal

# Efficient and Flexible DNN Acceleration with
# Hierarchical Structured Sparsity

# Hierarchical Structured Sparsity (HSS)

**Compose G:H sparsity patterns in a hierarchical fashion**

N-Rank HSS: G:H → G:H ... → G:H

*Rank N-1*   *Rank N-2*   *Rank 0*

*What does a 3:4→2:4 pattern look like?*



**Dense Vector**

# Hierarchical Structured Sparsity (HSS)

## Compose G:H sparsity patterns in a hierarchical fashion

*What does a* 3:4→2:4 *pattern look like?*

Rank1: 3 nonempty blocks out of the 4 blocks



**Vector with Rank1 Sparsity Applied**

# Hierarchical Structured Sparsity (HSS)

## Compose G:H sparsity patterns in a hierarchical fashion

*What does a 3:4→2:4 pattern look like?*

Rank1: 3 nonempty blocks out of the 4 blocks
Rank0: 2 nonzero values out of 4 values within the block



block0    block1    block2    block3

| 0 | | | 0 | | 0 | 0 | | 0 | 0 | 0 | 0 | | 0 | | 0 |

**Vector with Both Ranks' Sparsity Applied**

# Hierarchical Structured Sparsity (HSS)

**DNN Workloads Often Have Tensors with Multiple Dimensions**



**Per-Row** 3:4-2:4 Tensor

HSS can be applied to an arbitrary dimension in a multi-dimensional tensor

# HSS Introduces A Flexible Way to Express Sparsity Degrees

*4 sparsity degrees*

Rank 1

| 4:4 (0%) | 4:5 (20%) | 4:6 (33%) | 4:7 (43%) |
|---|---|---|---|

Sparsity Degree Spectrum

0%    20%    33%    43%

# HSS Introduces A Flexible Way to Express Sparsity Degrees

*4 sparsity degrees*

Rank 1

| 4:4 | 4:5 | 4:6 | 4:7 |
|------|------|------|------|
| (0%) | (20%) | (33%) | (43%) |

Rank 0

| 4:4 | 2:4 | 1:4 |
|------|------|------|
| (0%) | (50%) | (75%) |

*3 sparsity degrees*

## Sparsity Degree Spectrum



0%          50%

75%

# HSS Introduces A Flexible Way to Express Sparsity Degrees

*4 sparsity degrees*

Rank 1

| 4:4 (0%) | 4:5 (20%) | 4:6 (33%) | 4:7 (43%) |

Rank 0

| 4:4 (0%) | 2:4 (50%) | 1:4 (75%) |

*3 sparsity degrees*

**Multiplication of Fractions**

**4:5-2:4 (60%)**

**Sparsity Degree Spectrum**



**60%**

*17*

# HSS Introduces A Flexible Way to Express Sparsity Degrees

*4 sparsity degrees*

**Rank 1**

| 4:4 | 4:5 | 4:6 | 4:7 |
|-----|-----|-----|-----|
| (0%) | (20%) | (33%) | (43%) |

**Rank 0**

| 4:4 | 2:4 | 1:4 |
|-----|-----|-----|
| (0%) | (50%) | (75%) |

*3 sparsity degrees*

**Multiplication of Fractions**

| 4:5-2:4 (60%) | 4:6-1:4 (83%) |
|---|---|

## Sparsity Degree Spectrum

60%          83%

# HSS Introduces A Flexible Way to Express Sparsity Degrees

**4 sparsity degrees**

Rank 1

| 4:4 (0%) | 4:5 (20%) | 4:6 (33%) | 4:7 (43%) |

Rank 0

| 4:4 (0%) | 2:4 (50%) | 1:4 (75%) |

**3 sparsity degrees**

**Multiplication of Fractions**

*12 sparsity degrees*

| 4:4-4:4 (0%) | 4:5-4:4 (20%) | 4:6-4:4 (33%) | 4:7-4:4 (43%) | 4:4-2:4 (50%) | 4:5-2:4 (60%) | 4:6-2:4 (67%) | 4:7-2:4 (71%) | 4:4-1:4 (75%) | 4:5-1:4 (80%) | 4:6-1:4 (83%) | 4:7-1:4 (86%) |

Sparsity Degree Spectrum

0%   20%   33%   43%   50%   **60%**   67%   71%   75% 80%   86%
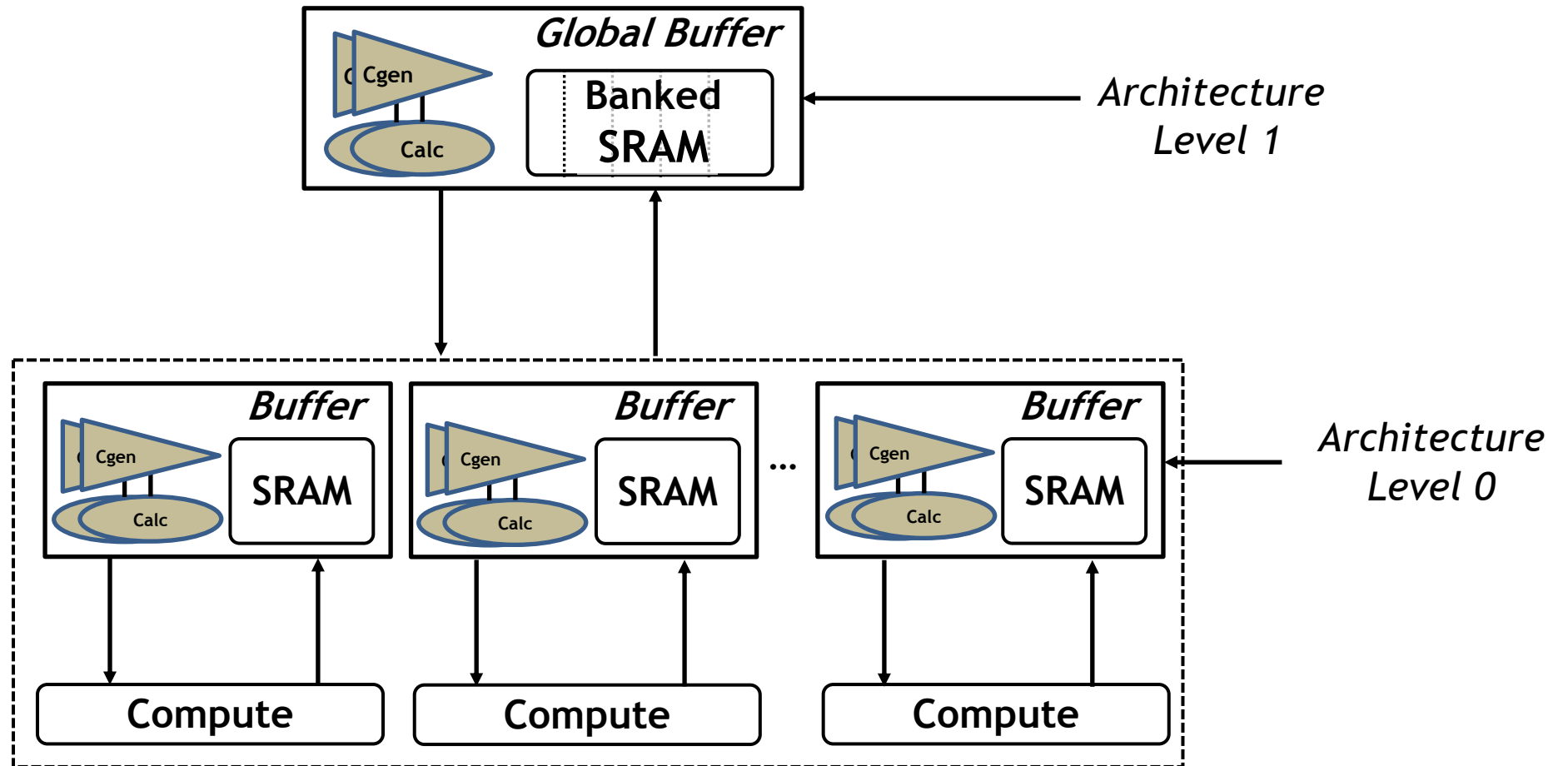                                                    83%

Fraction multiplication allows flexible representation of many sparsity degrees in a wide range

# HSS Enables Modularized Acceleration

Modularity of HSS allows different architecture levels to accelerate for different HSS ranks
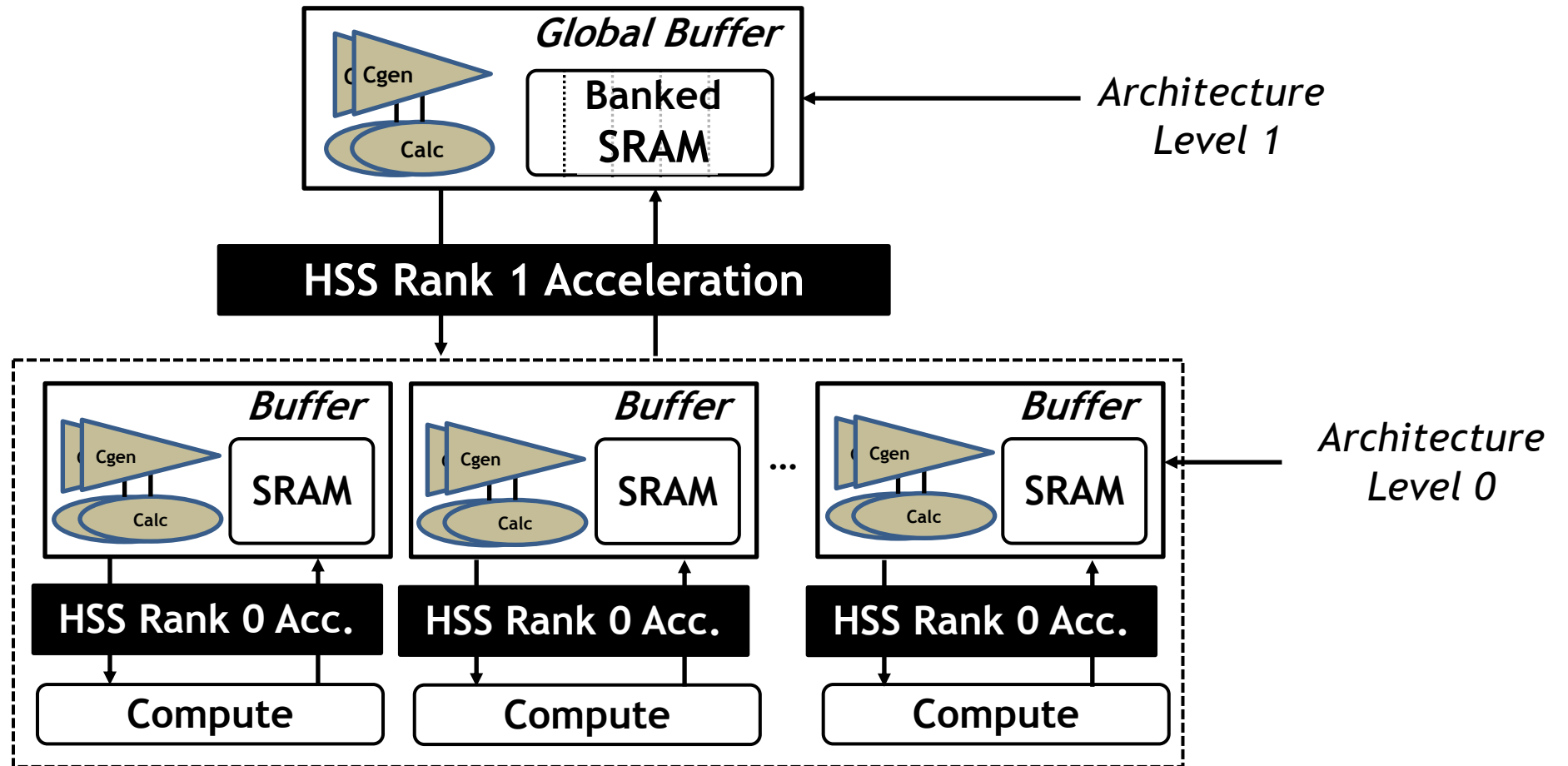
Example Accelerator Architecture Organization

# HSS Enables Modularized Acceleration

**Modularity of HSS allows different architecture levels to accelerate for different HSS ranks**
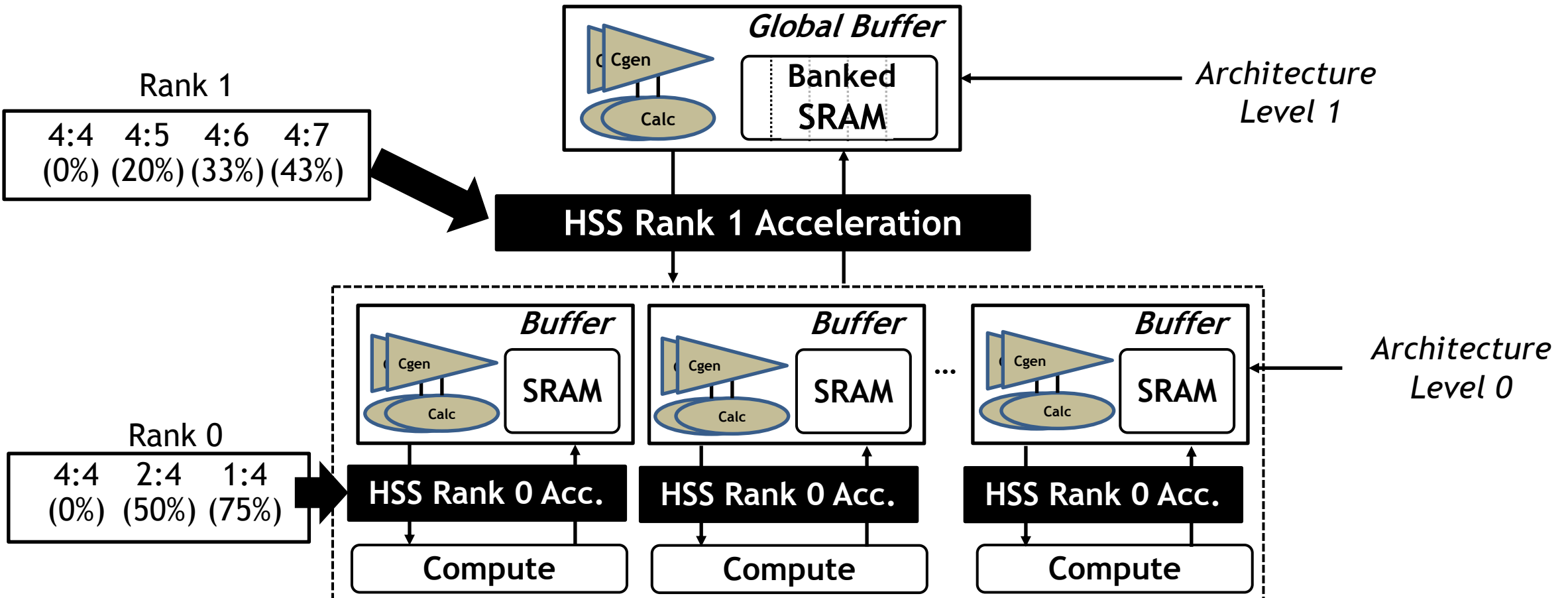
**Example Accelerator Architecture Organization**

# HSS Enables Modularized Acceleration

**Modularity of HSS allows different architecture levels to accelerate for different HSS ranks**

**Example Accelerator Architecture Organization**



Rank 1

| 4:4 | 4:5 | 4:6 | 4:7 |
|------|-------|-------|-------|
| (0%) | (20%) | (33%) | (43%) |

*Global Buffer*

Cgen

Calc

**Banked SRAM**

*Architecture Level 1*

**HSS Rank 1 Acceleration**

*Buffer* Cgen Calc SRAM

*Buffer* Cgen Calc SRAM

... 

*Buffer* Cgen Calc SRAM

*Architecture Level 0*

Rank 0

| 4:4 | 2:4 | 1:4 |
|------|-------|-------|
| (0%) | (50%) | (75%) |

**HSS Rank 0 Acc.**

**HSS Rank 0 Acc.**

**HSS Rank 0 Acc.**
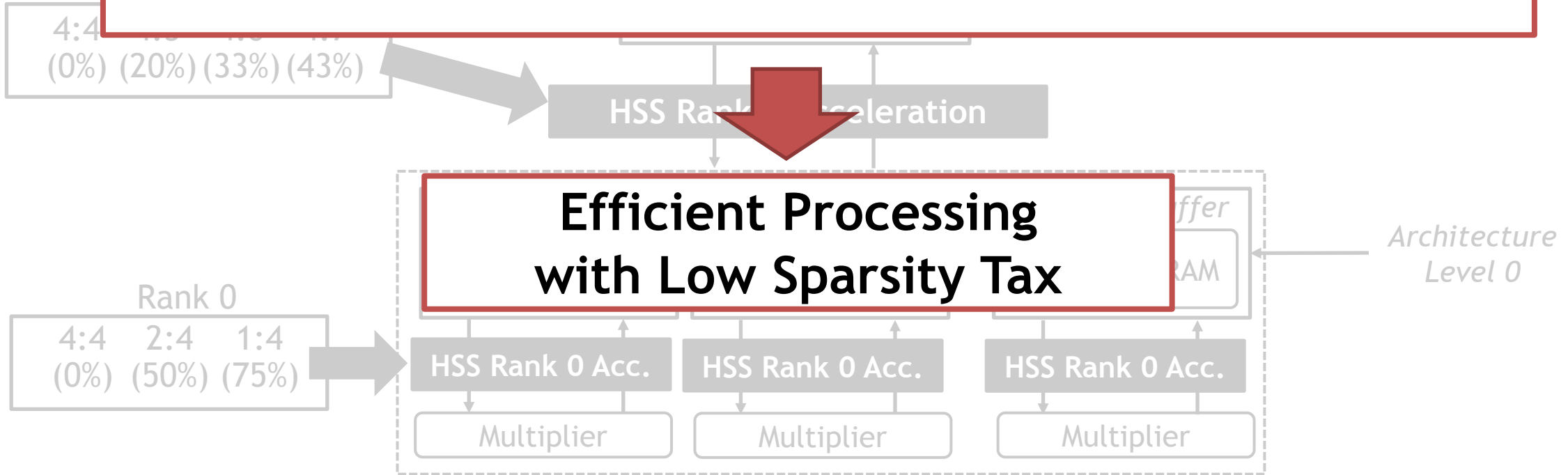
Compute

Compute

Compute

**Each level only needs to accelerate for a few sparsity degrees**
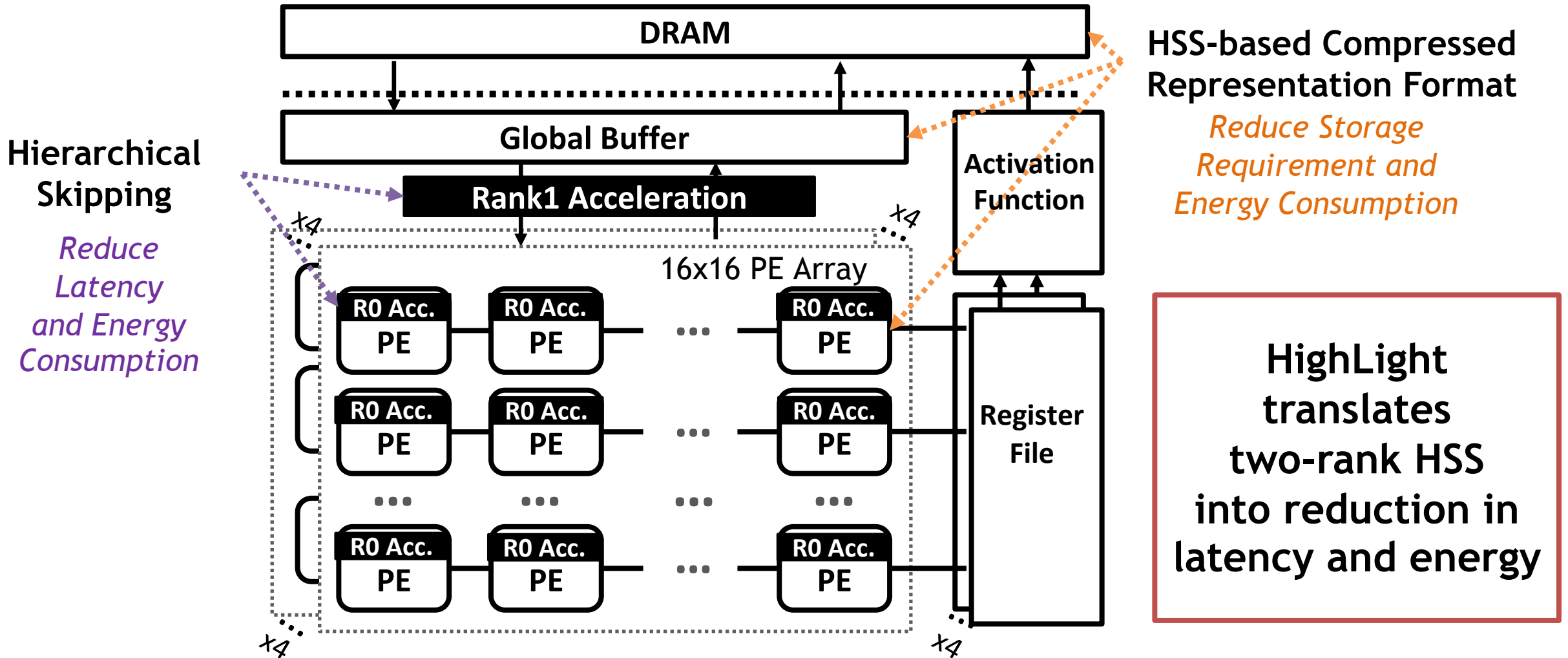
# HSS Enables Modularized Acceleration

Modularity of HSS allows different architecture levels to accelerate for different HSS ranks

**Simple Acceleration at Each Architecture Level Leads to Low Hardware Overhead**

4:4  (0%) (20%) (33%) (43%)

HSS Rank Acceleration

**Efficient Processing with Low Sparsity Tax**

Buffer

RAM

Architecture Level 0

Rank 0

4:4    2:4    1:4
(0%)  (50%)  (75%)

HSS Rank 0 Acc.    HSS Rank 0 Acc.    HSS Rank 0 Acc.

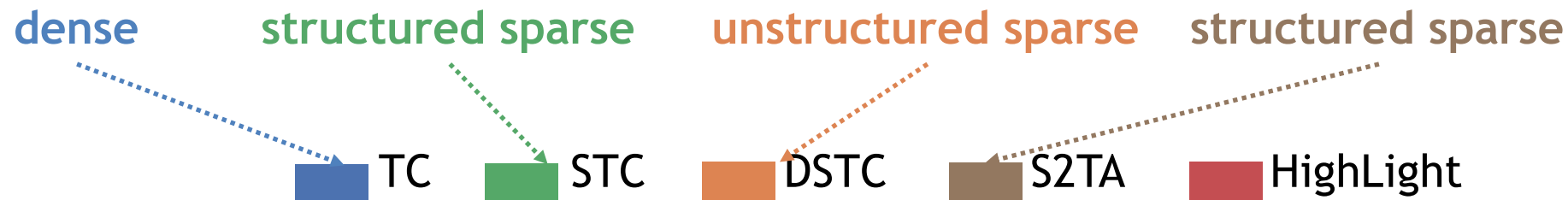Multiplier    Multiplier    Multiplier

Each level only needs to accelerate for a few sparsity degrees

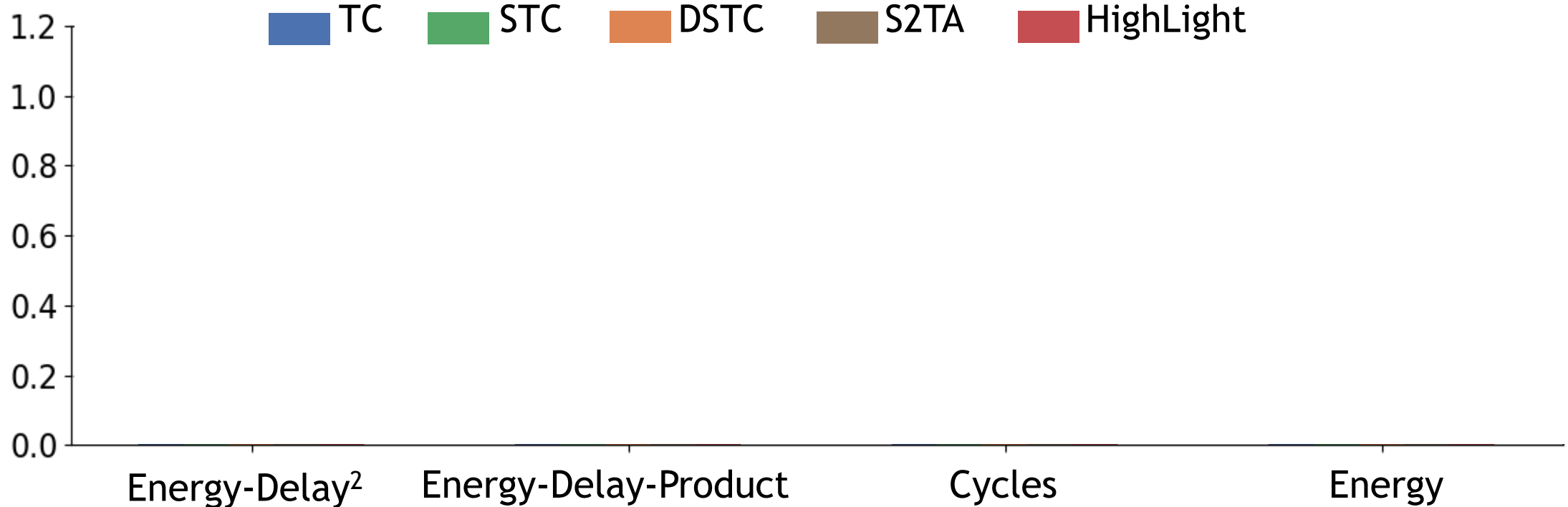# HighLight: Flexible and Efficient Sparse DNN Accelerator

# Experimental Results

# We Compare HighLight with Representative Designs

dense     structured sparse     unstructured sparse     structured sparse

TC     STC     DSTC     S2TA     HighLight

# Geomean Across Various Hardware Performance Metrics

**We evaluate the designs with synthetic workloads with different sparsity degrees ranging from 0%-75%**

**Geomean**
*(lower is better)*

■ TC   ■ STC   ■ DSTC   ■ S2TA   ■ HighLight

| | 1.2 |
|---|---|
| | 1.0 |
| | 0.8 |
| | 0.6 |
| | 0.4 |
| | 0.2 |
| | 0.0 |

Energy-Delay$^2$     Energy-Delay-Product     Cycles     Energy

# Geomean Across Various Hardware Performance Metrics

We evaluate the designs with synthetic workloads with different sparsity degrees ranging from 0%-75%
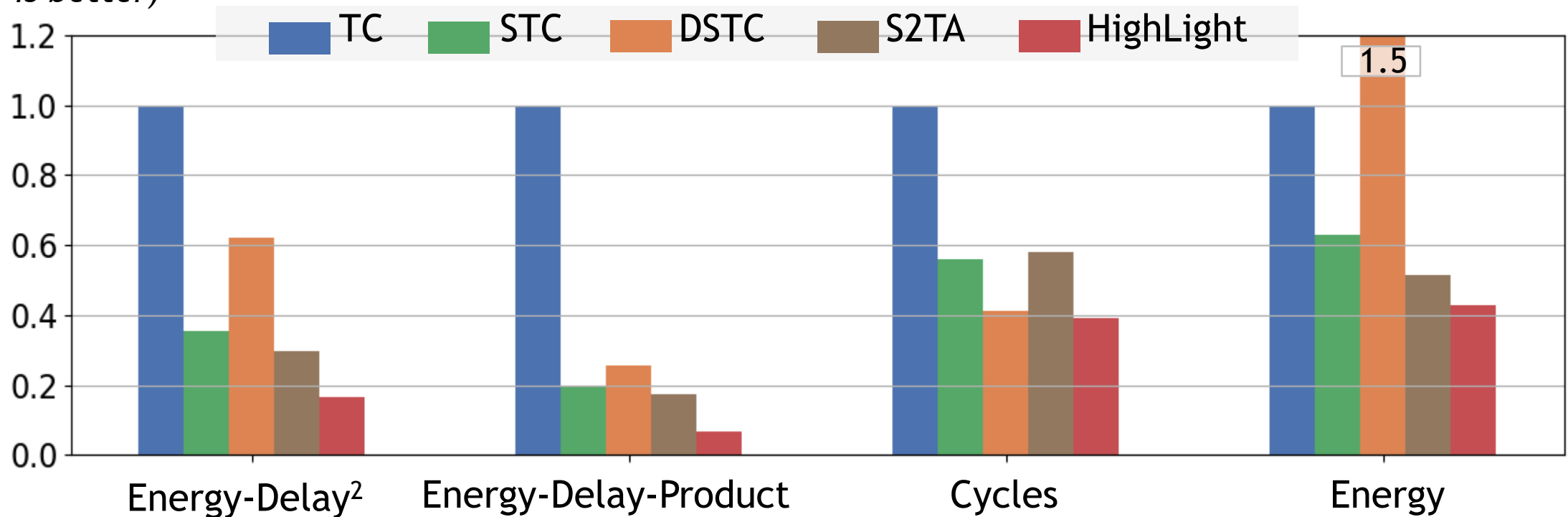
**Geomean**
*(lower is better)*



Legend: TC, STC, DSTC, S2TA, HighLight

Categories: Energy-Delay$^2$, Energy-Delay-Product, Cycles, Energy

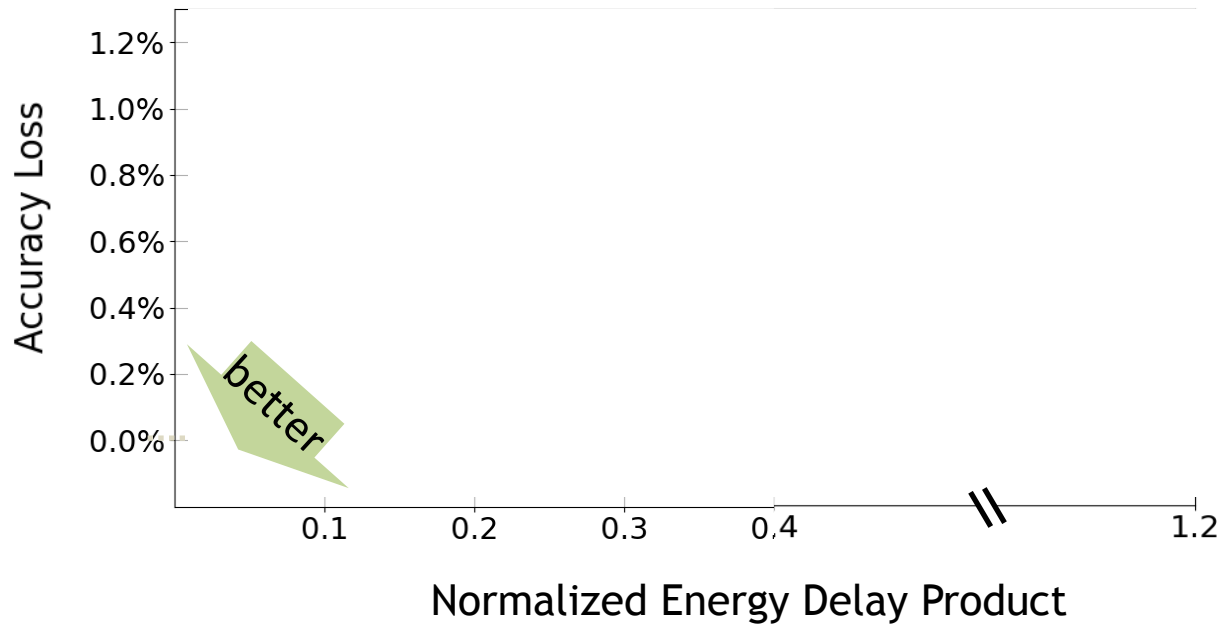**HighLight is efficient across evaluated metrics**

# Accuracy-Energy Delay Product Pareto Frontier

We evaluate the designs with representative DNNs pruned to different sparsity degrees, each with its respective sparsity structure (if any)
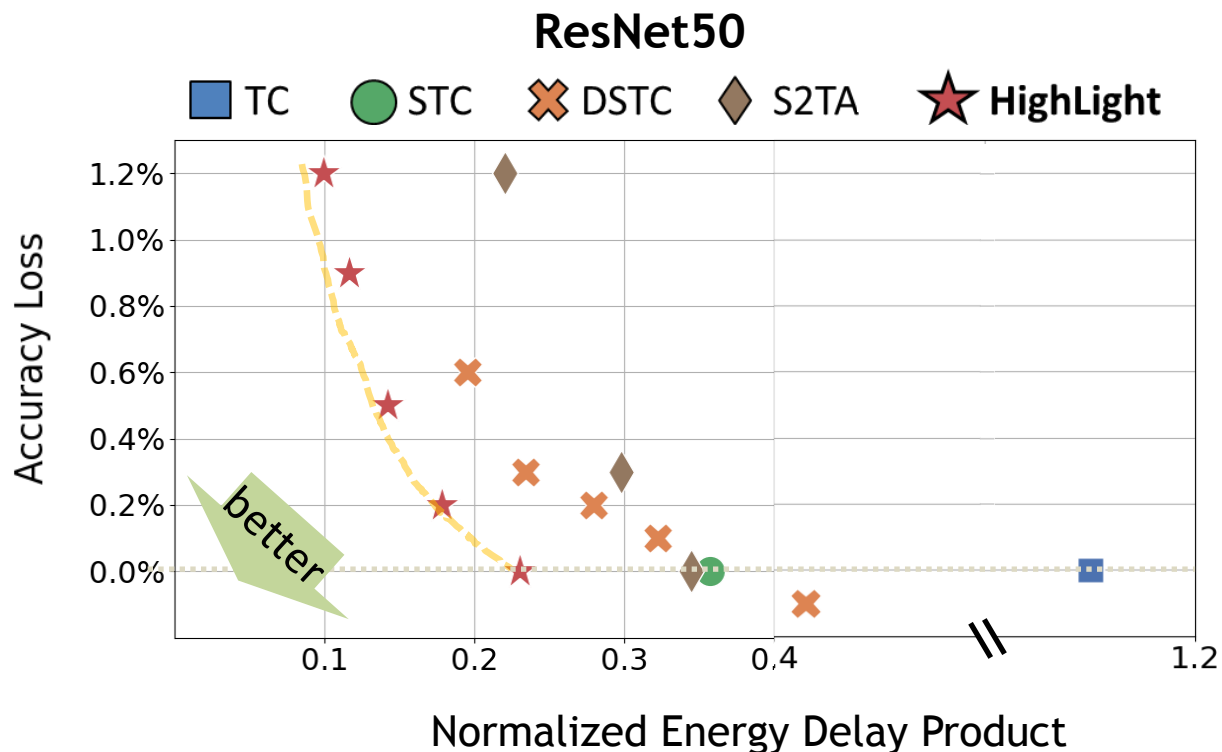
# Accuracy-Energy Delay Product Pareto Frontier

We evaluate the designs with representative DNNs pruned to different sparsity degrees, each with its respective sparsity structure (if any)
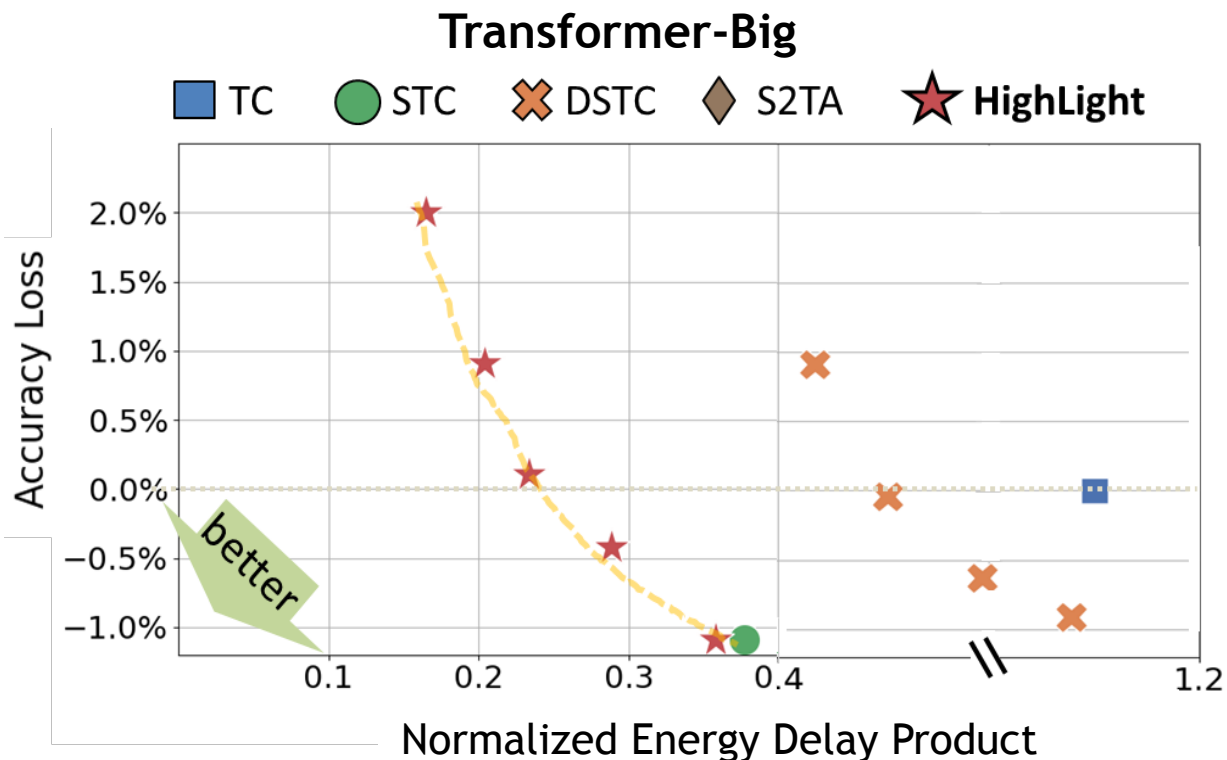
**ResNet50**

# Accuracy-Energy Delay Product Pareto Frontier

**We evaluate the designs with representative DNNs pruned to different sparsity degrees, each with its respective sparsity structure (if any)**
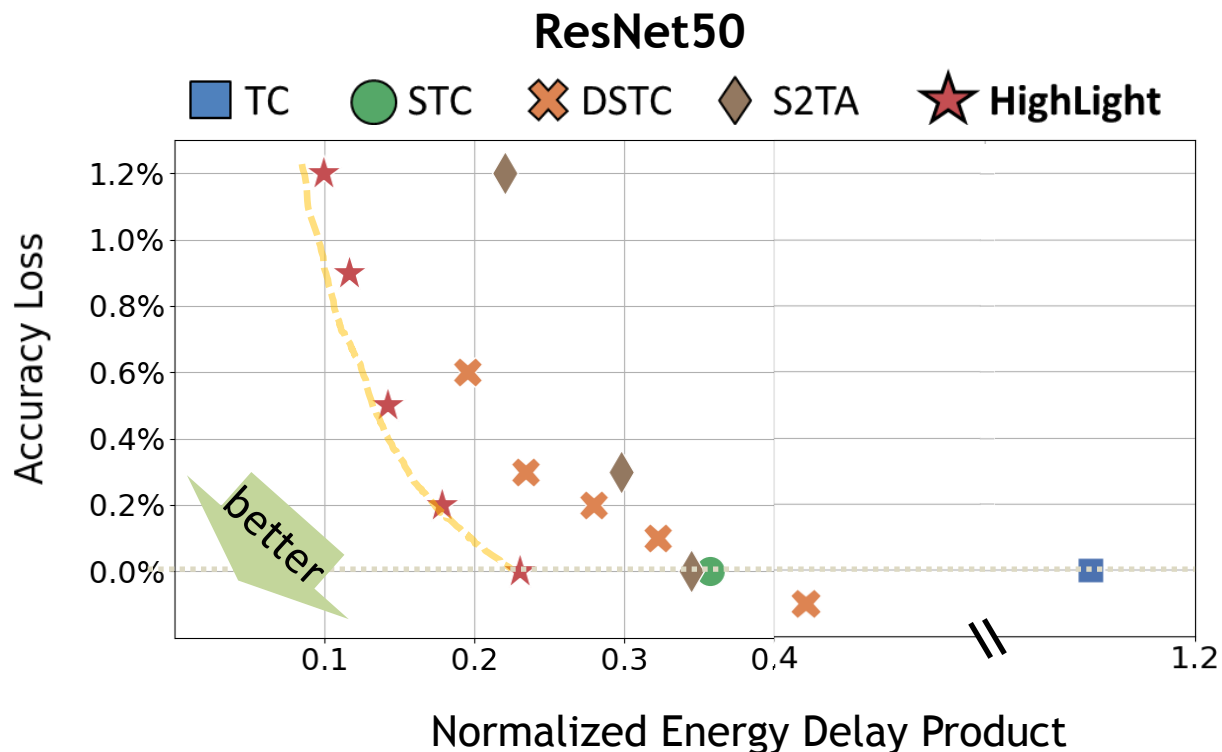


ResNet50

# Accuracy-Energy Delay Product Pareto Frontier

We evaluate the designs with representative DNNs pruned to different sparsity degrees, each with its respective sparsity structure (if any)



**HighLight sits on the accuracy-energy delay product pareto frontier**

# More Details in Paper!

- **How to systematically represent the diverse sparsity patterns in DNNs?**

    - Short answer: sparsity specification via fibertree abstraction.

- **What does HighLight's energy and area sparsity tax breakdowns look like?**

    - Short answer: low sparsity tax as HighLight independently accelerates simple sparsity patterns at different architecture levels.

- ...

# Summary

http://emze.csail.mit.edu/highlight

## Hierarchical Structured Sparsity (HSS)

- Composed of multiple levels of simple sparsity patterns
- Allows flexible expression of diverse sparsity degrees

## HighLight Accelerator

- Supports two-rank HSS for a few degrees at each level
- Implements low-overhead support for each rank at different architecture levels
- Ensures both efficiency and flexibility