# Algorithm-Hardware Co-Design for Data-Free Quantization and Inference using Contrastive Learning and Next Generation Arithmetic

**Presenter:** Akshat Ramachandran (**Advisor:** Tushar Krishna)

*Collaborators: Zishen Wan, Geonhwa Jeong, John Gustafson (ASU) and Souvik Kundu (Intel)*

Georgia Institute of Technology

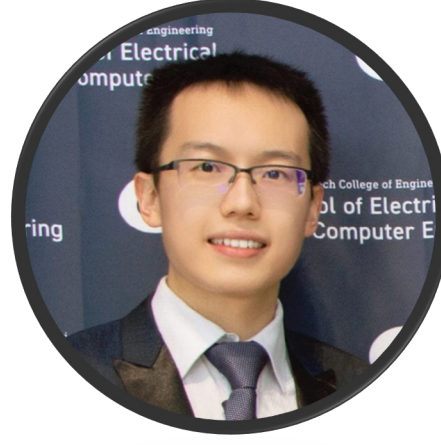**Contact:** akshat.r@gatech.edu (Akshat Ramachandran)

# Team

Students



Akshat Ramachandran (GT)

Zishen Wan (GT)

Geonhwa Jeong (GT)

PIs

Liaison

Tushar Krishna (GT)

John Gustafson (ASU)

Souvik Kundu (Intel)

**SRC**

# Executive Summary

- **Why ?**
  - Existing quantization data-formats lack adaptability and/or are inefficient for implementation in resource-constrained devices.
  - Privacy and security reasons limiting calibration data access for quantization.
  - Lack of a generalized quantization framework in prior art.

- **What ?**
  - Design of a composite data type called Logarithmic Posits (LP) that blends the adaptability of posits with hardware efficiency of Logarithmic Numbers.
  - Novel genetic-algorithm based data-free quantization framework leveraging contrastive learning for synthetic data-generation.
  - Unified mixed-precision Logarithmic Posit Accelerator (LPA) for high throughput DNN inference.
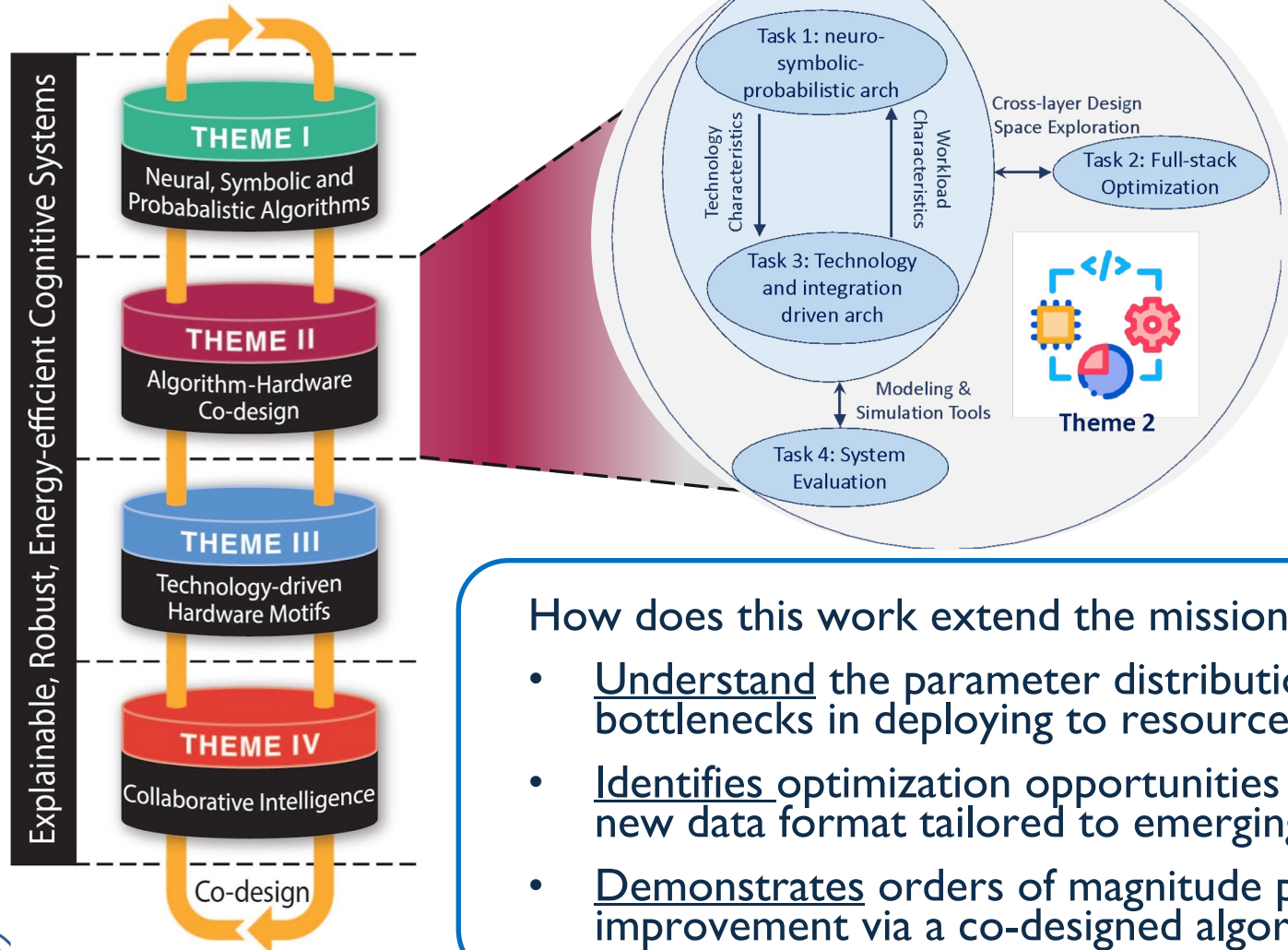
- **How ?**
  - Employs a co-design approach, integrating the development of the LP data type, the quantization framework, and the LPA architecture. This integrated approach allows for dynamic adaptation to DNN parameter distributions, resulting in highly efficient DNN inference.

**SRC**

# Relevance to the Center Goals



**Theme II Mission Statement**

The co-design of cognitive hardware driven by the evolution of cognitive workloads as well as the capabilities of future hardware technologies has the potential to unlock quantum improvements in processing efficiency
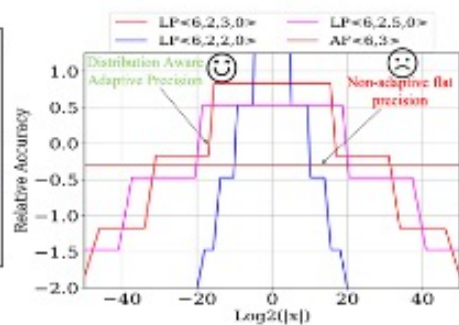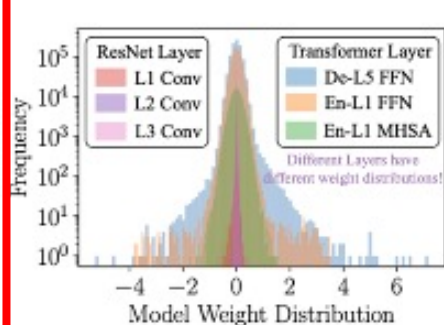
How does this work extend the mission of CoCoSys?

- <u>Understand</u> the parameter distribution variance of neuroinspired models and bottlenecks in deploying to resource-constrained devices.
- <u>Identifies</u> optimization opportunities over existing alternatives by demonstrating a new data format tailored to emerging neuroinspired models.
- <u>Demonstrates</u> orders of magnitude performance and energy efficiency improvement via a co-designed algorithm-hardware approach.
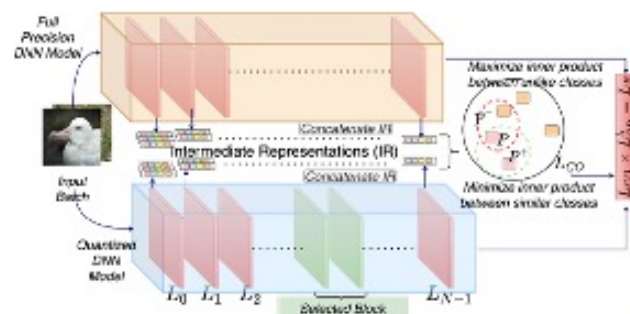
# Outline

Introduction

Automated Quantization

Data-Free Quantization

Next Generation Arithmetic: Logarithmic Posits

Hardware: Logarithmic Posit Accelerator
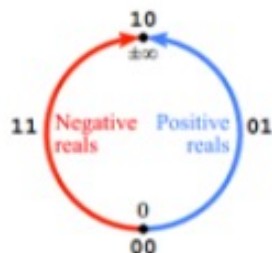
Hardware



SRC

# Outline

Introduction
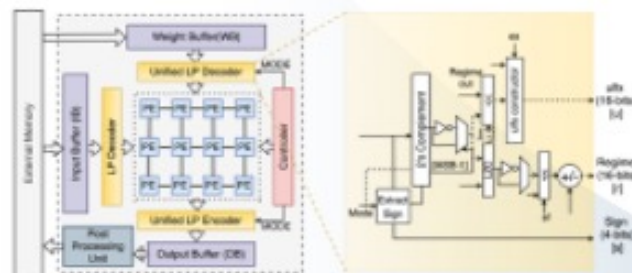
Automated Quantization

Data-Free Quantization

Next Generation Arithmetic: Logarithmic Posits

Hardware: Logarithmic Posit Accelerator

Hardware

SRC

4

# From Bulky DNNs to Sleek Edge Deployment!



- ❖ YoY increase in DNN sizes leads to escalating computational and storage demands!
- ❖ Limited compute, storage resources and energy budget of edge devices (e.g., phones) makes deployment challenging!

# From Bulky DNNs to Sleek Edge Deployment!



Source: Rahul, M et al. "A survey on deep neural network compression: challenges, overview, and solutions." arXiv preprint arXiv:2010.03954 (2020).

# From Bulky DNNs to Sleek Edge Deployment!

Model Compression techniques for efficient DNN deployment:

Focus of our work!



Pruning

Quantization

Source: Rahul, M et al. "A survey on deep neural network compression: challenges, overview, and solutions." arXiv preprint arXiv:2010.03954 (2020).

# Background: Types of Quantization Techniques

Quantization Aware Training (QAT):

    Higher Accuracy
    Improved Model Robustness

❌ Increased training complexity and time
❌ Large-scale data dependency

Post Training Quantization (PTQ):

    Speed and Simplicity
    Low data requirement (Can also be data-free)

❌ Potential accuracy drop
❌ Sensitivity to calibration dataset



Accuracy (%) — 86.9 / 85.5
Memory (GB) Per GPU — 29.8 / 8.6 / 21 / 3,180 — Training Time (min) — 4
#Data (K) — 393
QAT / PTQ

SRC

# Background: Types of Quantization Formats

**Uniform Quantization**

- Assigns same quantization step size cross the entire range of values.
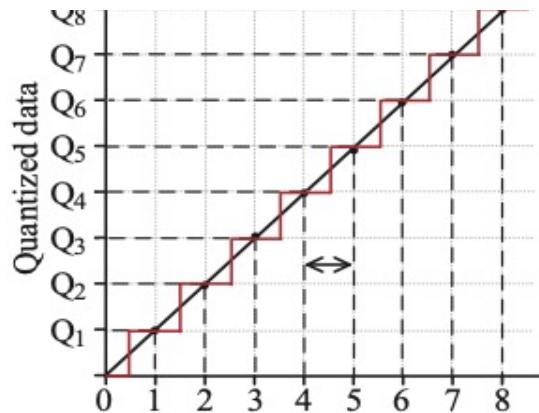- Results in simpler hardware but leads to higher quantization error for distributions with large variances.
- Example: Integers.

**Non-Uniform Quantization**

- Has variable distribution intervals, for selective quantization for significant data or extensive dynamic ranges.
- Potentially complex hardware and requires automated quantization algorithms.
- Example: IEEE-754 Floating-Point, Posits, Vector Quantization, Logarithmic Posits.

$$Q(X, \gamma, b) = clip(\left\lfloor \frac{X}{\gamma} \right\rceil, -2^{b-1} + 1, 2^{b-1} - 1)$$

$$Q(X, b, e, f) = sign(X) \times 2^e \times \mathbf{1}.\boldsymbol{f}$$





Source: Liu, Fangxin, et al. "Improving neural network efficiency via post-training quantization with adaptive floating-point." CVPR. 2021.

# Outline

**Algorithm**

**Introduction**



**Automated Quantization**



**Data-Free Quantization**



**Next Generation Arithmetic: Logarithmic Posits**



**Hardware: Logarithmic Posit Accelerator**



**Hardware**

SRC

# Problem: Non-smooth Loss Landscape

DNNs particularly depict a non-smooth quantization loss landscape.

Jagged and non-smooth peaks arise from weight perturbations during the quantization process

# Problem: Gradient-Descent Cannot Traverse this Landscape

First and second-order gradient cannot be used satisfactorily to traverse this non-smooth loss landscape because of the presence of multiple local-minima.

GA-based search is able to traverse this loss landscape and optimally identify quantization parameters

# Our Solution

We therefore adopt a GA based layer-wise quantization framework.

Why layer wise ? Perturbing too many layers at a time can cause traversal of a high-dimensional search space and does not guarantee convergence.

Source: Ramachandran, A, et al. "Algorithm-Hardware Co-Design of Distribution-Aware Logarithmic-Posit Encodings for Efficient DNN Inference." arXiv preprint arXiv:2403.05465 (2024).

# Background: Genetic Algorithm

Consider a population of randomly rabbits: some individuals are potentially faster and smarter than others.

# Background: Genetic Algorithm

➢ Slower, dumber rabbits are likely to be caught and eaten by foxes.

➢ Fast, smart rabbits survive to do what rabbits to best: make more rabbits!!

# Background: Genetic Algorithm

➢ The rabbits that survive breed with each other to generate offspring, which starts to mix up their genetic traits

➢ – Fast rabbits might breed with fast rabbits

➢ – Fast rabbits with slow rabbits

➢ – Smart with not-so-smart, etc…

➢ Furthermore, nature occasionally throws in a "wild hare" because genes can mutate.

# Background: Genetic Algorithm



➢ At the end the rabbits that survive will be the fastest, strongest and smartest among the population because the fox would've eaten the slow, weak and not-so-smart ones.

➢ In this analogy, an individual rabbit represents a solution to the problem (i.e. a single point in the state space).

➢ The foxes represent the problem constraints – Solutions that do well are likely to survive.

➢ We create similar such notions for quantization.

**SRC**

# Proposed: Automated Quantization Framework

# Proposed: Automated Quantization Framework

# Proposed: Automated Quantization Framework

# Proposed: Automated Quantization Framework

# Proposed: Automated Quantization Framework



Source: Ramachandran, A, et al. "Algorithm-Hardware Co-Design of Distribution-Aware Logarithmic-Posit Encodings for Efficient DNN Inference." arXiv preprint arXiv:2403.05465 (2024).

# Automated Quantization Framework: Step 1



**Step 1: Candidate Initialization**

Fitness value of each candidate

A quantization solution comprises an encoded vector $\Delta$ of length $2N$ and each set of 2 values represent the 2 integer quantization parameters of a layer $l$.

$$Q(X, \gamma, b) = clip\left(\left\lfloor \frac{X}{\gamma} \right\rceil, -2^{b-1} + 1, 2^{b-1} - 1\right)$$

# Automated Quantization Framework: Fitness Function Calculation

➢ Experiments in self-supervised learning and our own experiments suggest that contrastive learning tend to smooth the loss landscape.

$$\mathcal{L}_{i,j}^C = -\log \frac{\sum_{p+} \exp(\lambda_{i,j}^p \cdot \lambda_{i,j}^{p+}/\tau)}{\sum_{p+} \exp(\lambda_{i,j}^p \cdot \lambda_{i,j}^{p+}/\tau) + \sum_{p-} \exp(\lambda_{i,j}^p \cdot \lambda_{i,j}^{p-}/\tau)}$$



With Contrastive Learning

**Minimize** angle with $o^+$
**Maximize** dissimilarity with $o^-$

Quantized Model's Output

Corresponding FP Batch

Source: Ramachandran, A, et al. "Algorithm-Hardware Co-Design of Distribution-Aware Logarithmic-Posit Encodings for Efficient DNN Inference." arXiv preprint arXiv:2403.05465 (2024).

# Automated Quantization Framework: Fitness Function Calculation



$$\mathcal{L}_{\mathrm{CR}} = \sum_{l \in N} \#\mathrm{PARAM}(\mathbf{H}_l^{\mathrm{FP}}) \times n_l$$

The contrastive component of fitness function aims to align the distribution of quantized model's intermediate representations closely with the FP model, while the compression ratio metric incentivizes lower bit-widths

# Automated Quantization Framework: Step 2



Full Precision DNN Model

Input Batch

Quantized DNN Model

Concatenate IR

Intermediate Representations (IR)

Concatenate IR

$L_0$ $L_1$ $L_2$ $\cdots$ $L_{N-1}$

Selected Block

Maximize inner product between unlike classes

$P^-$

$P$

$P^+$

$L_{CO}$

Minimize inner product between similar classes

$L_{CO} \times L_{CR}^\lambda = L_F$

**Step 2: Re-generation**

**Crossover**

**Mutation**

For each selected block perform regeneration and crossover.

$$b_{child} = \mathbf{random}(\mathbf{min}(b_{p1}, b_{p2}) - 1, \mathbf{max}(b_{p1}, b_{p2}) + 1)$$

$$\gamma_{child} = \mathbf{mean}(\gamma_{p1}, \gamma_{p2}) + \eta(-10^{-3}, 10^3)$$

SRC

# Automated Quantization Framework: Step 3



We create additional random parents and use the regenerated child in the previous stage as the other parent to generate five diverse children.

Source: Ramachandran, A, et al. "Algorithm-Hardware Co-Design of Distribution-Aware Logarithmic-Posit Encodings for Efficient DNN Inference." arXiv preprint arXiv:2403.05465 (2024).

# Automated Quantization Framework: Step 4



**Step 4: Evaluation and Population Update**

Evaluate all generated children on the calibration data → Obtain $L_F$ for each candidate → Add $C_N$ and best of $(C_{N1}, C_{N2}....)$ to population as $(\Delta_C, L_F)$

Source: Ramachandran, A, et al. "Algorithm-Hardware Co-Design of Distribution-Aware Logarithmic-Posit Encodings for Efficient DNN Inference." arXiv preprint arXiv:2403.05465 (2024).

# Activation Quantization

Activation quantization sensitivity closely aligns with that of the weight parameters producing them.

The output activation quantization parameters for layer i are determined as,

$$b_{act}[i] = \mathbf{min}(8, b[i] \times 2) \text{ and } \gamma_{act}[i] = \gamma_{act}[i-1] + \gamma[i].$$

# Challenges with Current Quantization Scheme



- Current PTQ method requires access to calibration data drawn from the training set.
- However, due to privacy and security concerns the data might not always be accessible.

# Outline

## Introduction



## Automated Quantization



## Data-Free Quantization

**Under Submission**



## Next Generation Arithmetic: Logarithmic Posits



## Hardware: Logarithmic Posit Accelerator



Hardware

SRC

# Prior Art: Data-Free Quantization for CNNs



- In CNNs, data-free quantization techniques synthesize data for calibrating the quantized model according to the batch normalization (BN) statistics of FP32.
- This is a reliable and efficient method for distilling data from the FP32 model for CNNs.

# Data-Free Quantization: Transformer-Based models (ViTs)



**Does not explore semantically meaningful interpatch relations!**

- ViTs or transformer-based models do not have a BN layer or any layer that holds statistics of the data it has been trained on.
- Previous methods rely on maximizing the entropy of the self-attention layer outputs.
- This ignores inter-patch relations and generates semantically non-informative and less realistic data.

Source: Li, Zhikai, et al. "Psaq-vit v2: Toward accurate and general data-free quantization for vision transformers." IEEE Transactions on Neural Networks and Learning Systems (2023).

# Data-Free Quantization



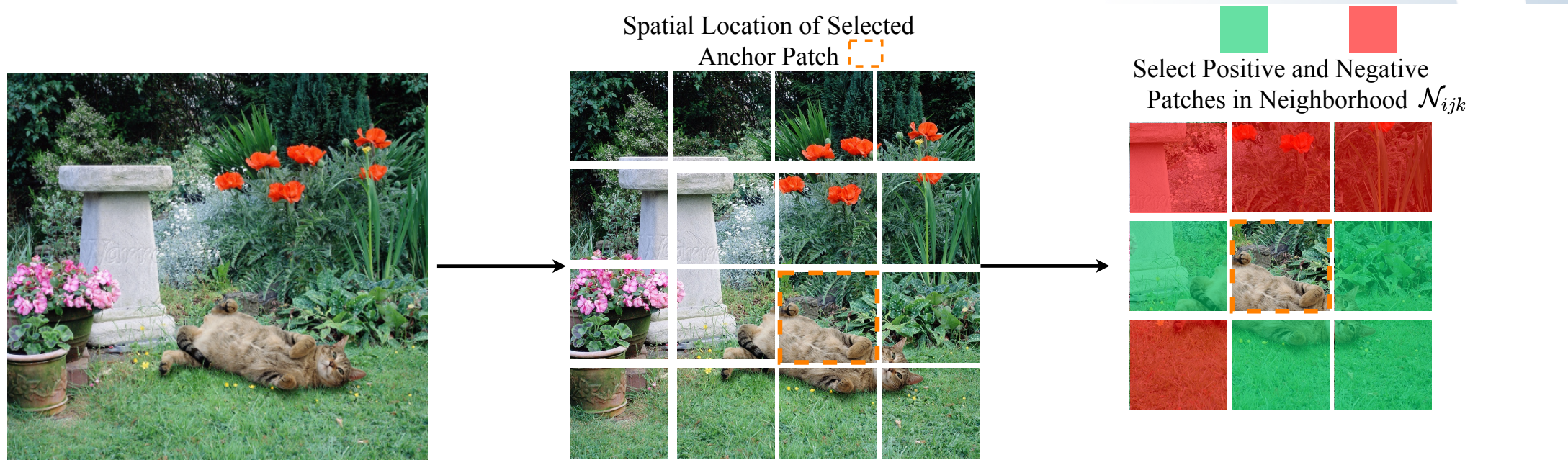**Smooth Neighborhood**

**Dramatic Loss Variations**

- Synthetic image
- Real-world image

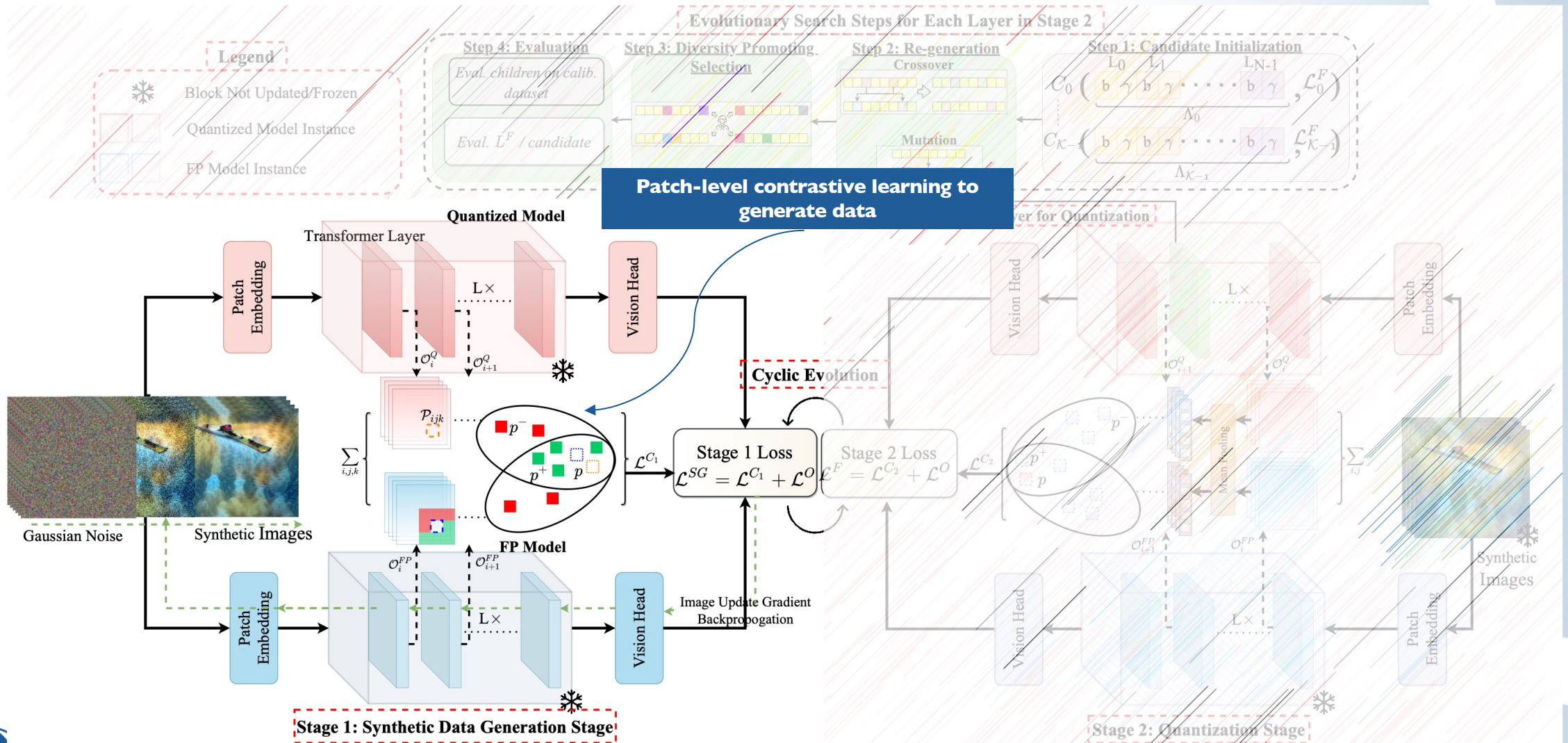- Using a global scheme without considering inter-patch relations does not result in realistic and robust images.
- The lower semantic content in generated synthetic images severely impacts quantized model generalizability.

Source: Li, Zhikai, et al. "Psaq-vit v2: Toward accurate and general data-free quantization for vision transformers." IEEE Transactions on Neural Networks and Learning Systems (2023).

# Data-Free Quantization: Proposed Approach



Spatial Location of Selected Anchor Patch

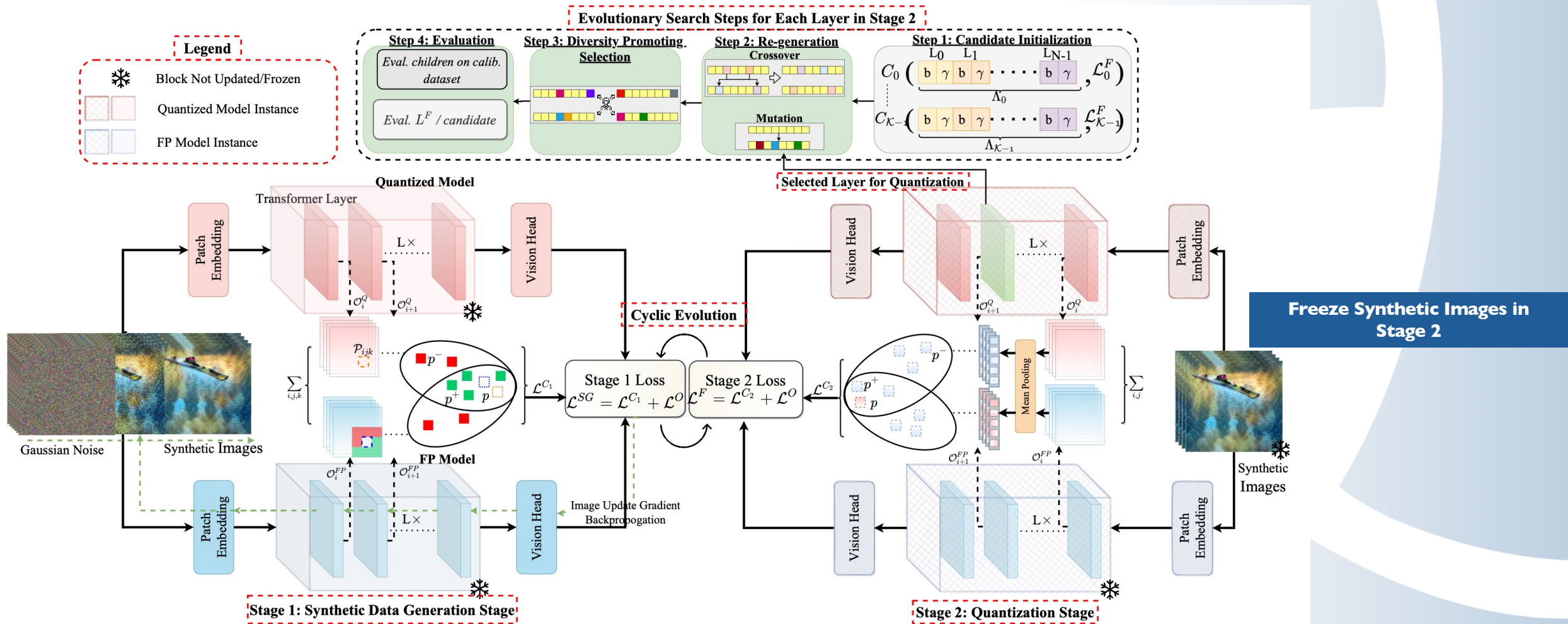Select Positive and Negative Patches in Neighborhood $\mathcal{N}_{ijk}$

- We propose a data-generation scheme that leverages the architectural characteristics of ViTs i.e, patch-level attention and the inherent properties of real-images to generate semantically rich and meaningful data.

- We develop a contrastive learning scheme that treats semantically similar patches in a neighborhood as positive and others as negative.
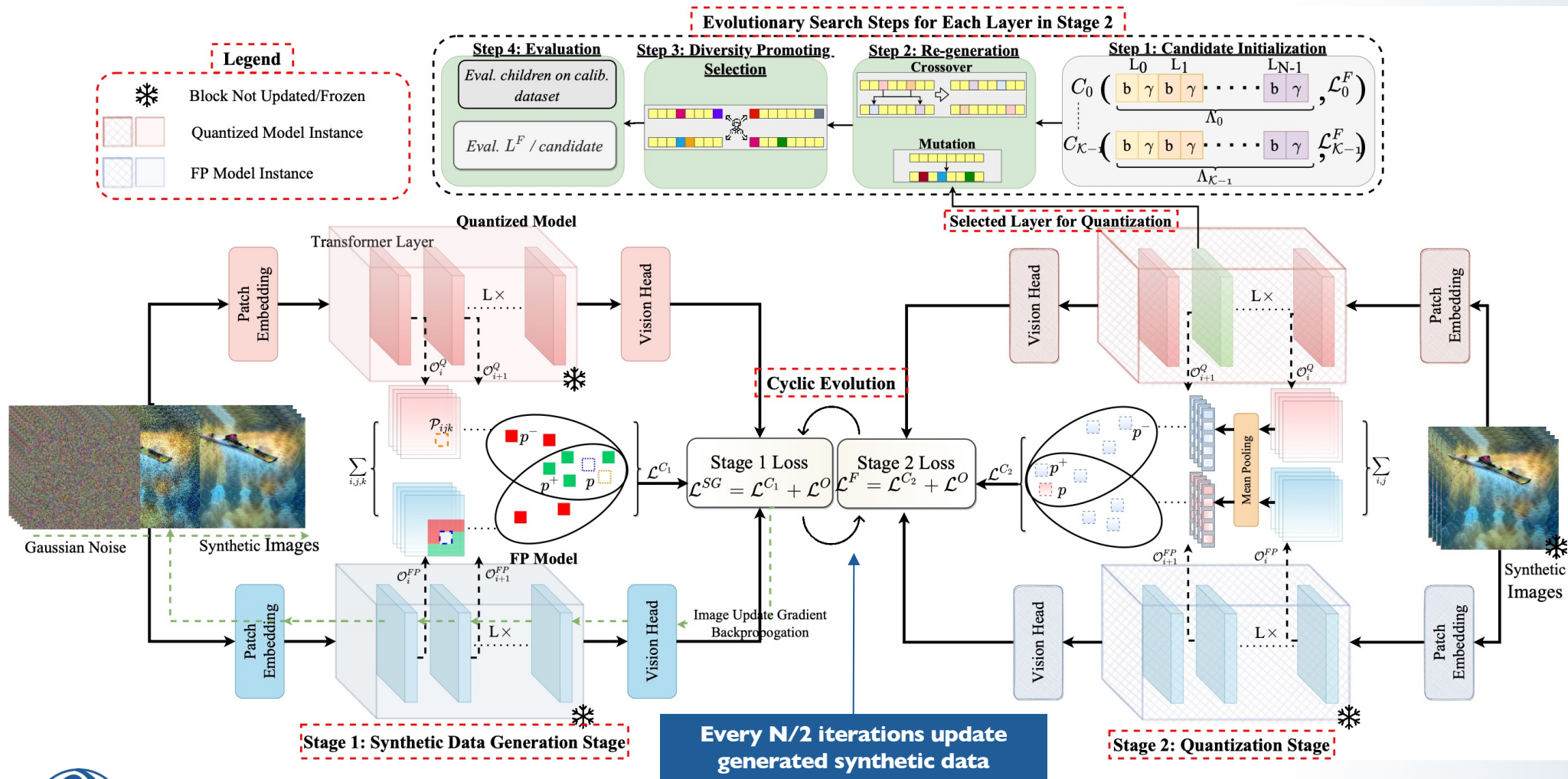
# Data-Free Quantization: Complete Pipeline
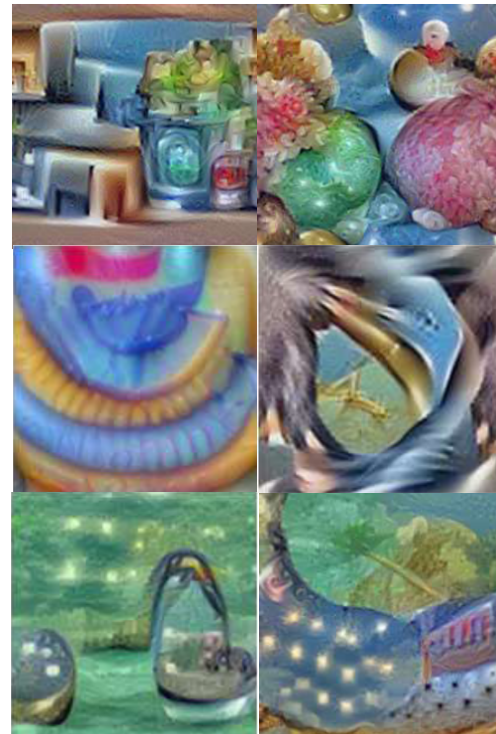
# Data-Free Quantization: Complete Pipeline
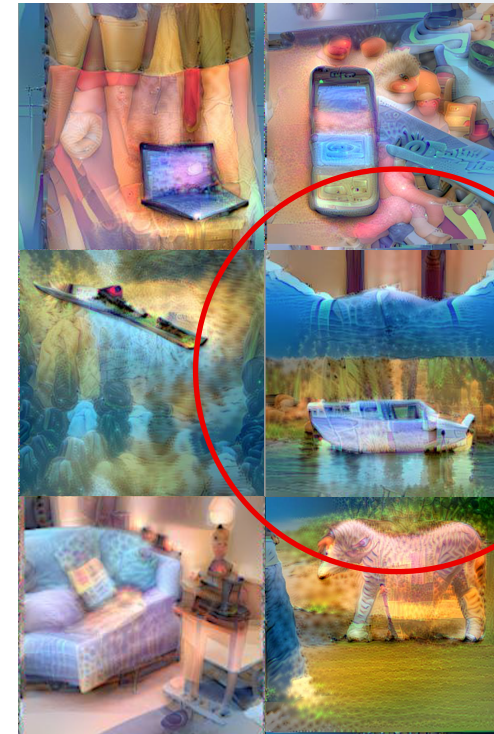
# Data-Free Quantization: Complete Pipeline

# Data-Free Quantization: Generated Samples



**PSAQ-ViT v1**

**PSAQ-ViT v2**

**Ours**

Notice the semantic meaningfulness!

# Data-Free Quantization Performance

| Model | Method | Data | #Images | W/A | Top-1 | W/A | Top-1 |
|---|---|---|---|---|---|---|---|
| | Baseline | - | - | 32/32 | 84.53 | 32/32 | 84.53 |
| | PSAQ-ViT v1 | S | 32 | 8/8 | 37.36 | 4/8 | 25.34 |
| | PTQ4ViT | R | 32 | 8/8 | **84.25** | 4/8 | 67.16 |
| ViT-B | FQ-ViT | R | 1000 | 8/8 | 83.31 | 4/8 | **78.73** |
| | RepQ-ViT | R | 32 | 8/8 | 81.45 | 4/8 | 76.29 |
| | **CLAMP-ViT (Ours)** | S | 32 | 8/8 | 84.1̲5̲ | 4/8 | **78.73** |
| | Baseline | - | - | 32/32 | 72.21 | 32/32 | 72.21 |
| | PSAQ-ViT v1 | S | 32 | 8/8 | 71.56 | 4/8 | 65.57 |
| | PSAQ-ViT v2 | S | 32 | 8/8 | 72.17 | 4/8 | 68.61 |
| DeiT-T | FQ-ViT | R | 1000 | 8/8 | 71.61 | 4/8 | 66.91 |
| | RepQ-ViT | R | 32 | 8/8 | 72.05 | 4/8 | 68.75 |
| | **CLAMP-ViT (Ours)** | S | 32 | 8/8 | **72.17** | 4/8 | **69.93** |
| | Baseline | - | - | 32/32 | 79.85 | 32/32 | 79.85 |
| | PSAQ-ViT v1 | S | 32 | 8/8 | 76.92 | 4/8 | 73.23 |
| | PSAQ-ViT v2 | S | 32 | 8/8 | **79.56** | 4/8 | 76.36 |
| DeiT-S | PTQ4ViT | R | 32 | 8/8 | 79.47 | 4/8 | - |
| | FQ-ViT | R | 1000 | 8/8 | 79.17 | 4/8 | 76.93 |
| | RepQ-ViT | R | 32 | 8/8 | 79.55 | 4/8 | 76.75 |
| | **CLAMP-ViT (Ours)** | S | 32 | 8/8 | 79.5̲5̲ | 4/8 | **77.03** |

- **~2.2% Improvement over DFQ and ~1% Improvement over PTQ methods!**

**SRC**

# Data-Free Quantization Performance

| Model | Method | Data | #Images | W/A | Top-1 | W/A | Top-1 |
|-------|--------|------|---------|-----|-------|-----|-------|
| | Baseline | - | - | 32/32 | 84.53 | 32/32 | 84.53 |
| | PSAQ-ViT v1 | S | 32 | 8/8 | 37.36 | 4/8 | 25.34 |
| | PTQ4ViT | R | 32 | 8/8 | **84.25** | 4/8 | 67.16 |
| ViT-B | FQ-ViT | R | 1000 | 8/8 | 83.31 | 4/8 | **78.73** |
| | RepQ-ViT | R | 32 | 8/8 | 81.45 | 4/8 | 76.29 |
| | **CLAMP-ViT (Ours)** | S | 32 | 8/8 | 84.19 | 4/8 | **78.73** |
| | Baseline | - | - | 32/32 | 72.21 | 32/32 | 72.21 |
| | PSAQ-ViT v1 | S | 32 | 8/8 | 71.56 | 4/8 | 65.57 |
| | PSAQ-ViT v2 | S | 32 | 8/8 | 72.17 | 4/8 | 68.61 |
| DeiT-T | FQ-ViT | R | 1000 | 8/8 | 71.61 | 4/8 | 66.91 |
| | RepQ-ViT | R | 32 | 8/8 | 72.05 | 4/8 | 68.75 |
| | **CLAMP-ViT (Ours)** | S | 32 | 8/8 | **72.17** | 4/8 | **69.93** |
| | Baseline | - | - | 32/32 | 79.85 | 32/32 | 79.85 |
| | PSAQ-ViT v1 | S | 32 | 8/8 | 76.92 | 4/8 | 73.23 |
| | PSAQ-ViT v2 | S | 32 | 8/8 | **79.56** | 4/8 | 76.36 |
| DeiT-S | PTQ4ViT | R | 32 | 8/8 | 79.47 | 4/8 | - |
| | FQ-ViT | R | 1000 | 8/8 | 79.17 | 4/8 | 76.93 |
| | RepQ-ViT | R | 32 | 8/8 | 79.55 | 4/8 | 76.75 |
| | **CLAMP-ViT (Ours)** | S | 32 | 8/8 | 79.55 | 4/8 | **77.03** |

- **For W4/A8 our method shows significant performance boost over all the existing alternatives despite be**
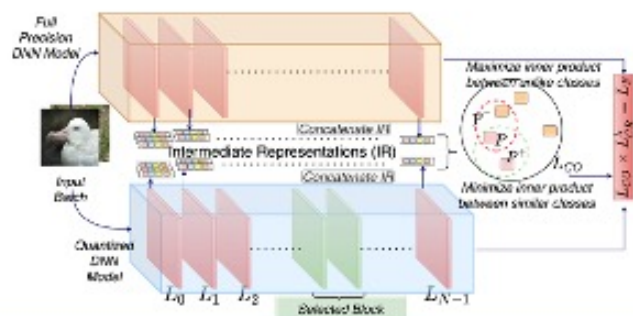
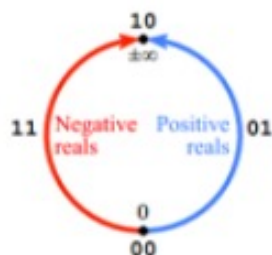# Outline

## Introduction



## Automated Quantization
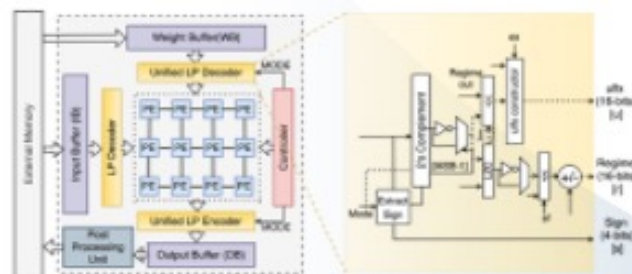


## Data-Free Quantization



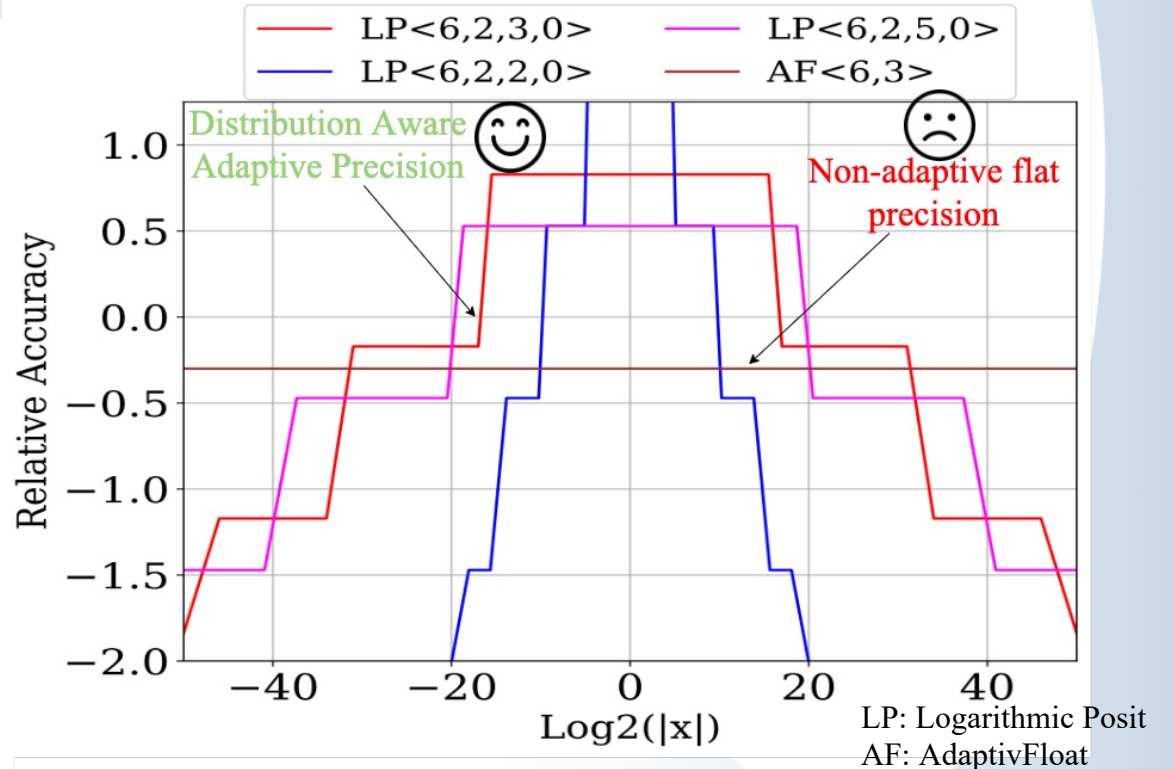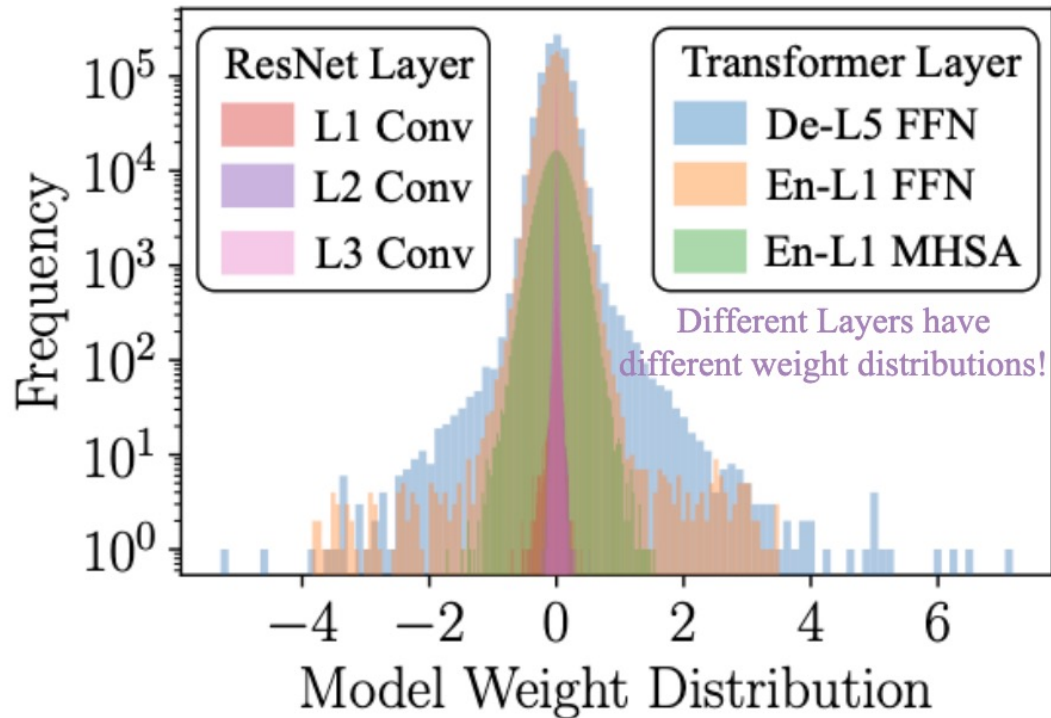## Next Generation Arithmetic: Logarithmic Posits



**To be presented at DAC'24**

## Hardware: Logarithmic Posit Accelerator



Hardware

SRC
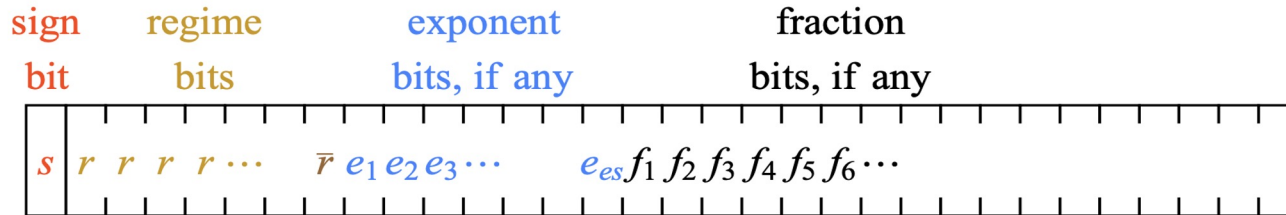
# Challenges with Traditional Quantization



- ➤ Uniform Quantization techniques lack the dynamic range and distributional variance required of DNNs.
- ➤ Floating-point based non-uniform quantization techniques fail to adapt to the different DNN distributions and have flat accuracy.
- ➤ Vector Quantization is adaptive but introduces additional codebook overhead.

# Next Generation Arithmetic: Posits

sign | regime | exponent | fraction
bit | bits | bits, if any | bits, if any

$$s \; r \; r \; r \; r \cdots \; \overline{r} \; e_1 \, e_2 \, e_3 \cdots \; e_{es} f_1 \, f_2 \, f_3 \, f_4 \, f_5 \, f_6 \cdots$$

Generic n-bits Posit format

| float | Sign | Exponent (Size ES) | | Fraction (Size F) |
|---|---|---|---|---|
| | s | $e_1 e_2 e_3 \ldots e_{es}$ | | $f_1 f_2 f_3 \ldots f_F$ |
| posit | Sign | Regime (Run Length K) | Exponent (if any) (Size ES) | Fraction (if any) |
| | s | $r_1 r_2 r_3 \ldots \overline{r_k}$ | $e_1 e_2 e_3 \ldots$ | $f_1 f_2 f_3 \ldots$ |

$$\texttt{posit val} = (-1)^{sign} * (2^{2^{es}})^k * 2^{expo} * 1.frac$$

➢ Possesses the unique field called regime that dynamically varies between [2, n-1] bits. It dynamically varies the exponent and fraction fields to give tapered accuracy and varied dynamic ranges.

➢ The regime is a run-length encoding of m 0s (1s) terminated by a 1(0), respectively, or by the final bit. The regime value $k$ is determined as $k = -m$ if the first bit of the regime is 0, or $k = m-1$ otherwise.

# Logarithmic Posits

$$x\langle n, es, rs, sf \rangle = (-1)^{sign} \times 2^{2^{es} \times k - sf} \times 2^{\mathbf{ulfx}}$$

LP value decoding

➤ A composite datatype that blends the adaptability of Posits with the hardware-efficiency of Logarithmic Number System (LNS).

➤ Parameterizations introduced for adaptability:
➤ **Bits (n):** Identify optimal precision for a DNN layer.
➤ **Exponent Size (es):** Controls dynamic range.
➤ **Regime Size (rs):** Controls distribution shape.
➤ **Scale Factor (sf):** Adjusts distribution position.

➤ Express standard fraction and exponent in the logarithmic domain as a unified fixed-point exponent of the power of two as $2^{ulfx}$, where ulfx=e+f.

**SRC**

# Floating-Point v. Posit v. Logarithmic Posit Decoding

*Binary Number:* 01001110

| **Floating-Point (E4M3), Bias = 7** | **Posit (ES = 2)** | **Logarithmic Posit (ES = 1, RS = 7, SF = 0)** |

### Floating-Point (E4M3), Bias = 7

*0_1001_110*

**Sign:** +
**E-b:** *1001 - 0111 = 9 - 7 = 2*
**M:** *1.110 = 1.75*

**Value** = +1 * 1.75 * $2^2$

### Posit (ES = 2)

*0_10_01_110*

**Sign:** +
**Regime:** *10 = 0*
**E:** *01 = 1*
**M:** 1.110 = 1.75

**Value** = $+1 * (2^{(2^2)})^0 * 2^1 * 1.75$

### Logarithmic Posit (ES = 1, RS = 7, SF = 0)

*0_10_0_1110*

**Sign:** +
**Regime:** *10 = 0*
**ulfx:** *0.1110 = 0.875*

**Value** = $+1 * (2^{(2^2)})^0 * 2^0 * 2^{0.875}$

**SRC**

# Floating-Point v. Posit v. Logarithmic Posit Decoding

*Binary Number:* 01001110

### Floating-Point (E4M3), Bias = 7

$0\_1001\_110$

**Sign:** +
**E-b:** *1001 - 0111 = 9 - 7 = 2*
**M:** *1.110 = 1.75*

**Value** $= +1 * 1.75 * 2^2$

### Posit (ES = 2)

$0\_10\_01\_110$

**Sign:** +
**Regime:** *10 = 0*
**E:** *01 = 1*
**M:** 1.110 = 1.75

**Value** $= +1 * (2^{(2^2)})^0 * 2^1 * 1.75$

### Logarithmic Posit (ES = 1, RS = 7, SF = 0)

$0\_10\_0\_1110$

**Sign:** +
**Regime:** *10 = 0*
**ulfx:** *0.1110 = 0.875*

**Value** $= +1 * (2^{(2^2)})^0 * 2^0 * 2^{0.875}$

# Floating-Point v. Posit v. Logarithmic Posit Decoding

*Binary Number:* 01001110

**Floating-Point (E4M3), Bias = 7**

$0\_1001\_110$

**Sign:** +

**E-b:** *1001 - 0111 = 9 - 7 = 2*

**M:** *1.110 = 1.75*

**Value** $= +1 * 1.75 * 2^2$

**Posit (ES = 2)**

$0\_10\_01\_110$

**Sign:** +
**Regime:** *10 = 0*
**E:** *01 = 1*
**M:** 1.110 = 1.75

**Value** $= +1 * (2^{(2^2)})^0 * 2^1 * 1.75$

**Logarithmic Posit (ES = 1, RS = 7, SF = 0)**

$0\_10\_0\_1110$

**Sign:** +
**Regime:** *10 = 0*
**ulfx:** *0.1110 = 0.875*

**Value** $= +1 * (2^{(2^2)})^0 * 2^0 * 2^{0.875}$

**SRC**

# Advantages of Logarithmic Posits
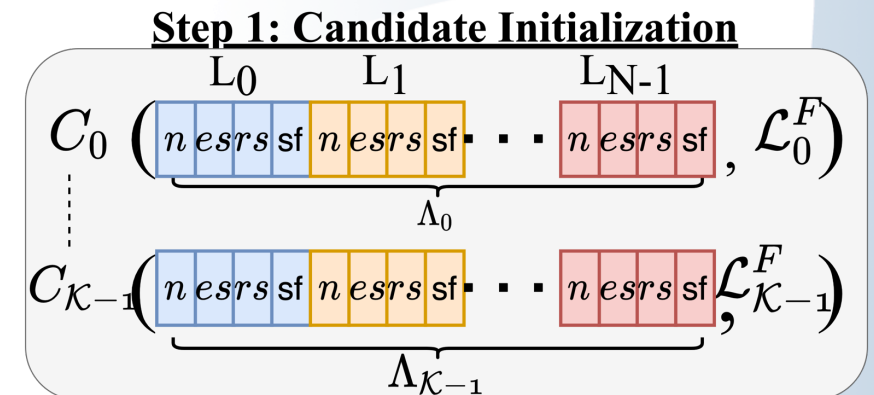


LP: Logarithmic Posit
AF: AdaptivFloat

➢ By parameterizing the individual bitfields of LP, we are able to adapt the LP representation distribution to the required DNN data distribution.

➢ While also enjoying the hardware efficiency of LNS.

**SRC**

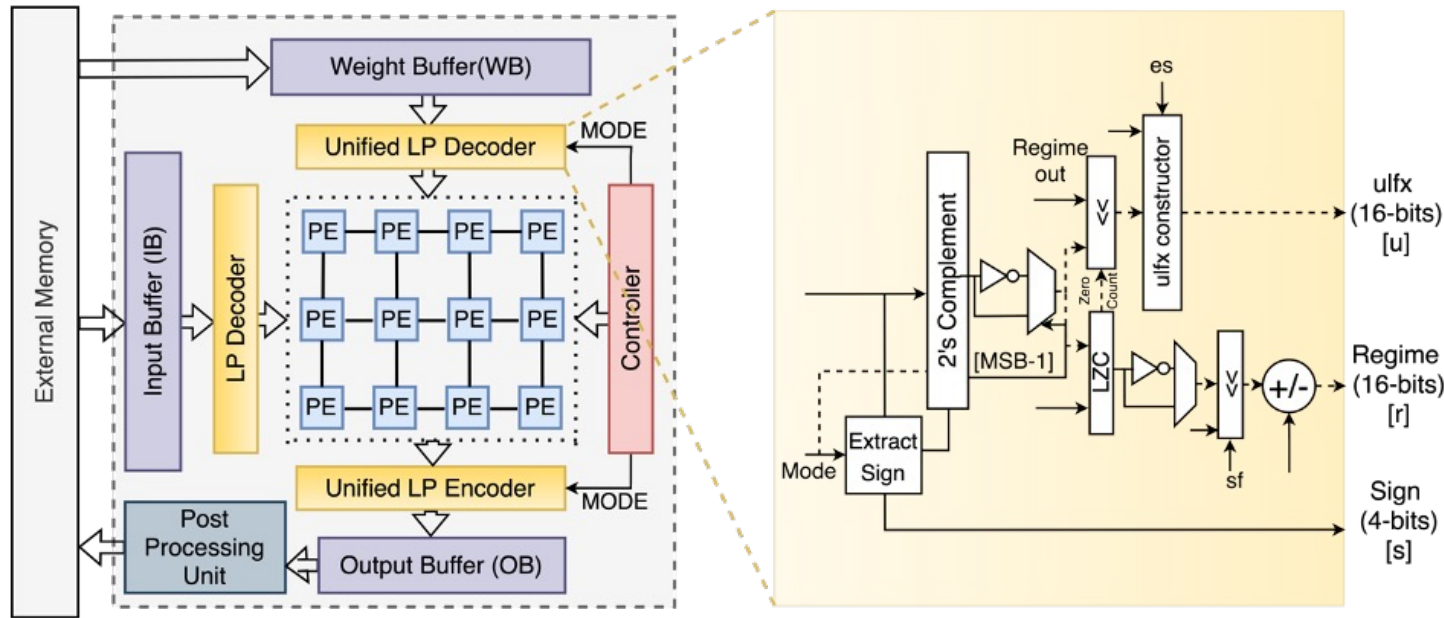# Changes in Algorithm to support LP

Uniform Quantization Candidates
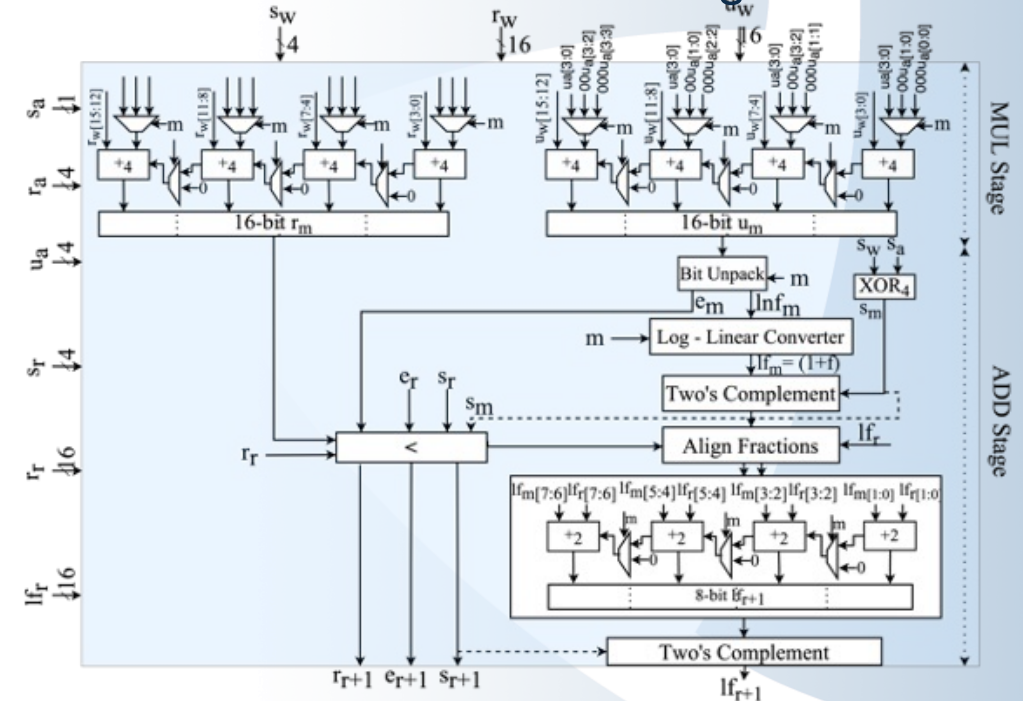
Logarithmic Posit Quantization Candidates

**Step 1: Candidate Initialization**



$$C_0 \left( \underbrace{\boxed{b \mid \gamma \mid b \mid \gamma} \cdots \boxed{b \mid \gamma}}_{\Lambda_0}, \mathcal{L}_0^F \right)$$

$$C_{\mathcal{K}-1} \left( \underbrace{\boxed{b \mid \gamma \mid b \mid \gamma} \cdots \boxed{b \mid \gamma}}_{\Lambda_{\mathcal{K}-1}}, \mathcal{L}_{\mathcal{K}-1}^F \right)$$

**Step 1: Candidate Initialization**

$$C_0 \left( \underbrace{\boxed{n \; es \; rs \; sf} \boxed{n \; es \; rs \; sf} \cdots \boxed{n \; es \; rs \; sf}}_{\Lambda_0}, \mathcal{L}_0^F \right)$$

$$C_{\mathcal{K}-1} \left( \underbrace{\boxed{n \; es \; rs \; sf} \boxed{n \; es \; rs \; sf} \cdots \boxed{n \; es \; rs \; sf}}_{\Lambda_{\mathcal{K}-1}}, \mathcal{L}_{\mathcal{K}-1}^F \right)$$

Our algorithm is general and can be easily applied to any quantization format by simply changing the candidate vector!

# Outline

## Introduction



## Automated Quantization



## Data-Free Quantization



## Next Generation Arithmetic: Logarithmic Posits



## Hardware: Logarithmic Posit Accelerator

**To be presented at DAC'24**

SRC

# Hardware: Logarithmic Posit Accelerator

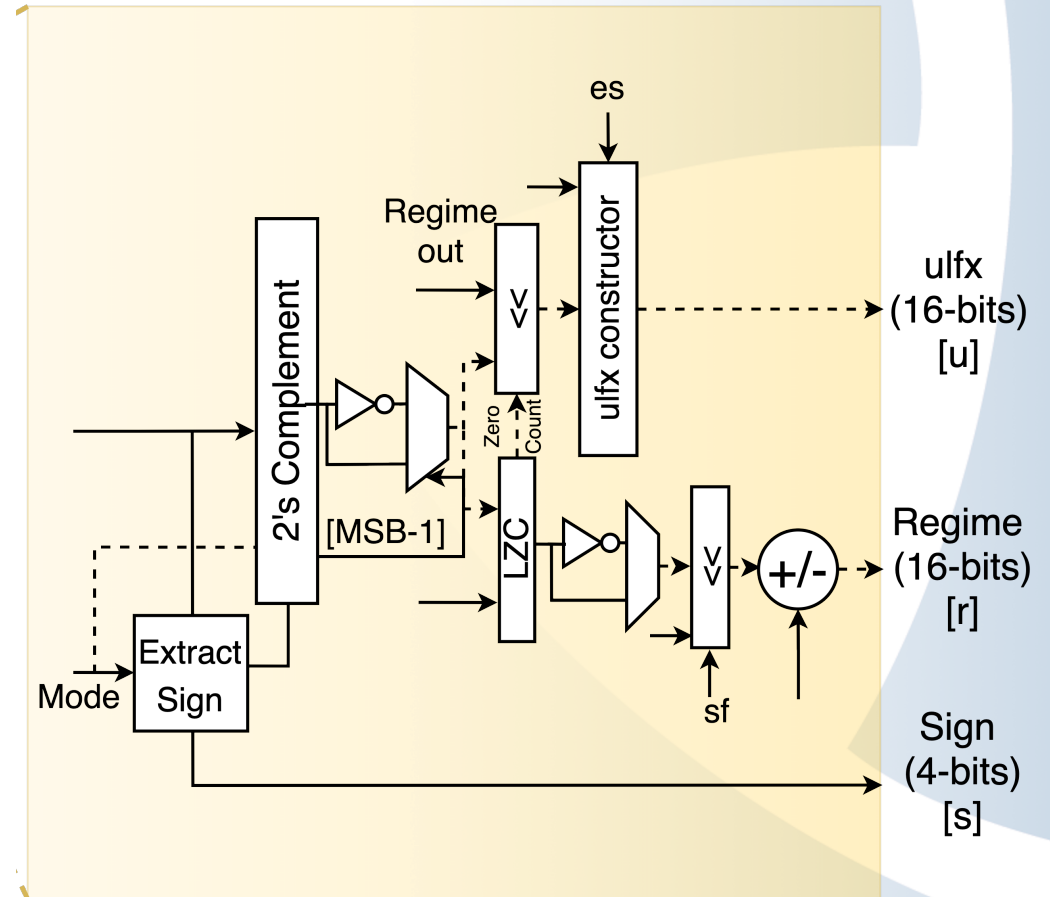Modifications to a WS Systolic Array

Mixed-Precision PE Design

# Modifications to Systolic Array
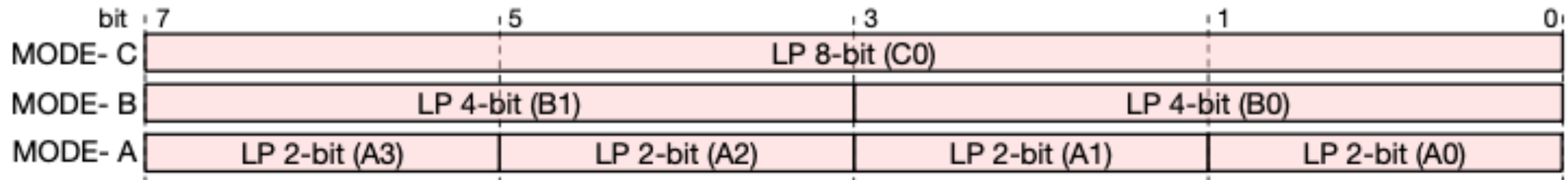
# Modifications to Systolic Array

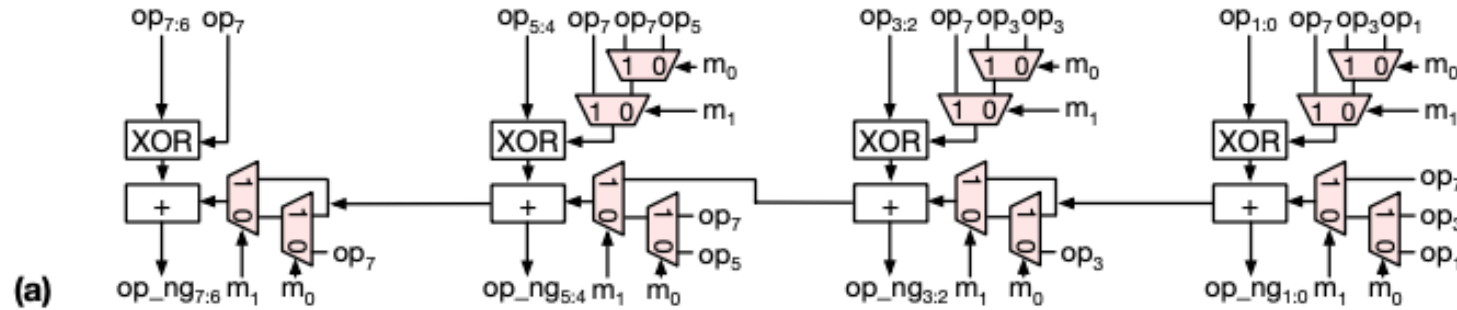Decoder/Encoder Overhead is 1.03% of compute area for an 8x8 systolic array.

Exhibits weak scaling!!

# Concept of MODE

The MODE field is used to identify the precision for computation!

# Implementing Mixed-Precision Components

# Mixed-Precision PE Architecture

Mixed-precision LP multiply-accumulate unit made entirely of 4-bit integer adder building blocks.
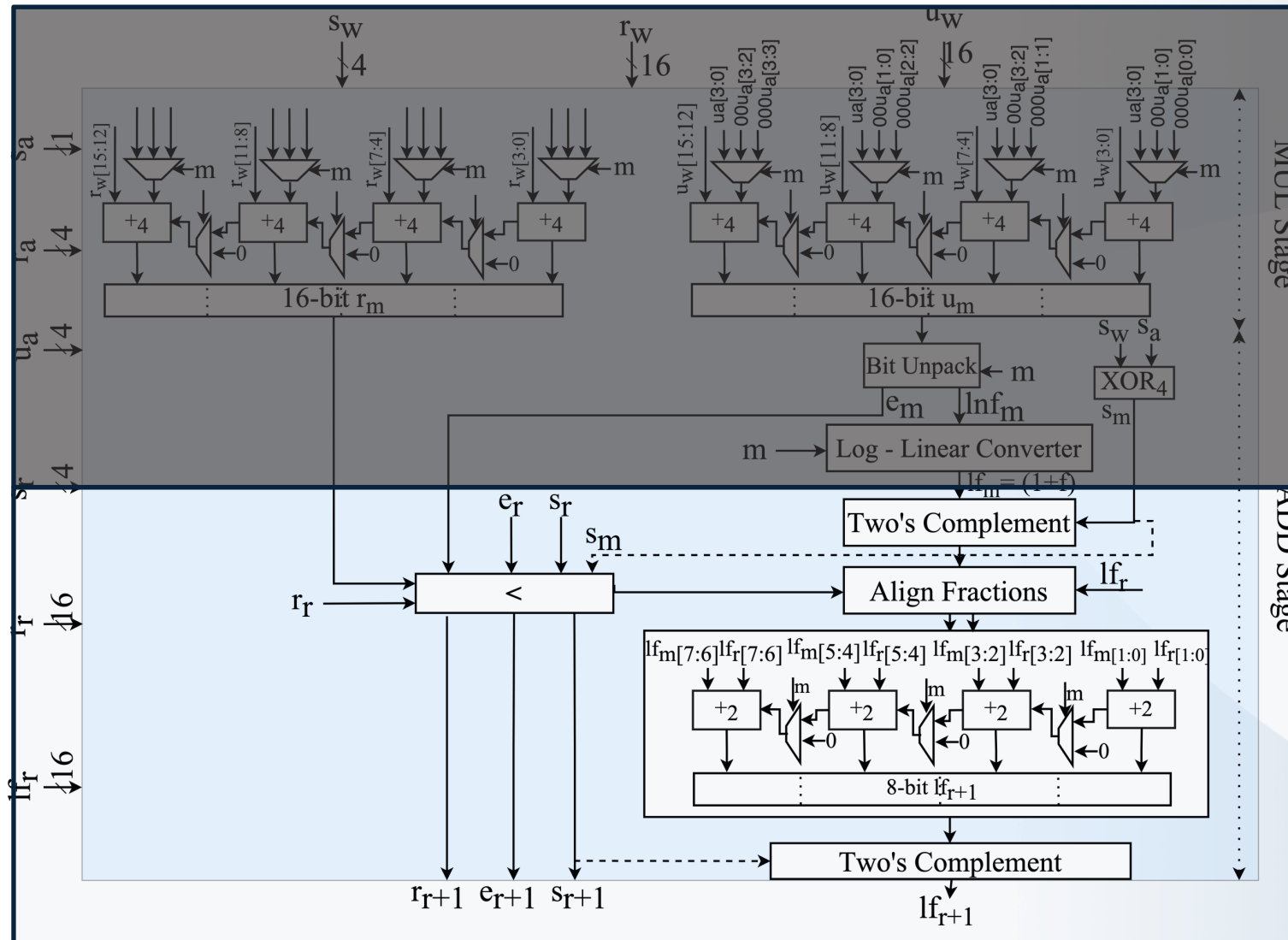
# Mixed-Precision PE Architecture



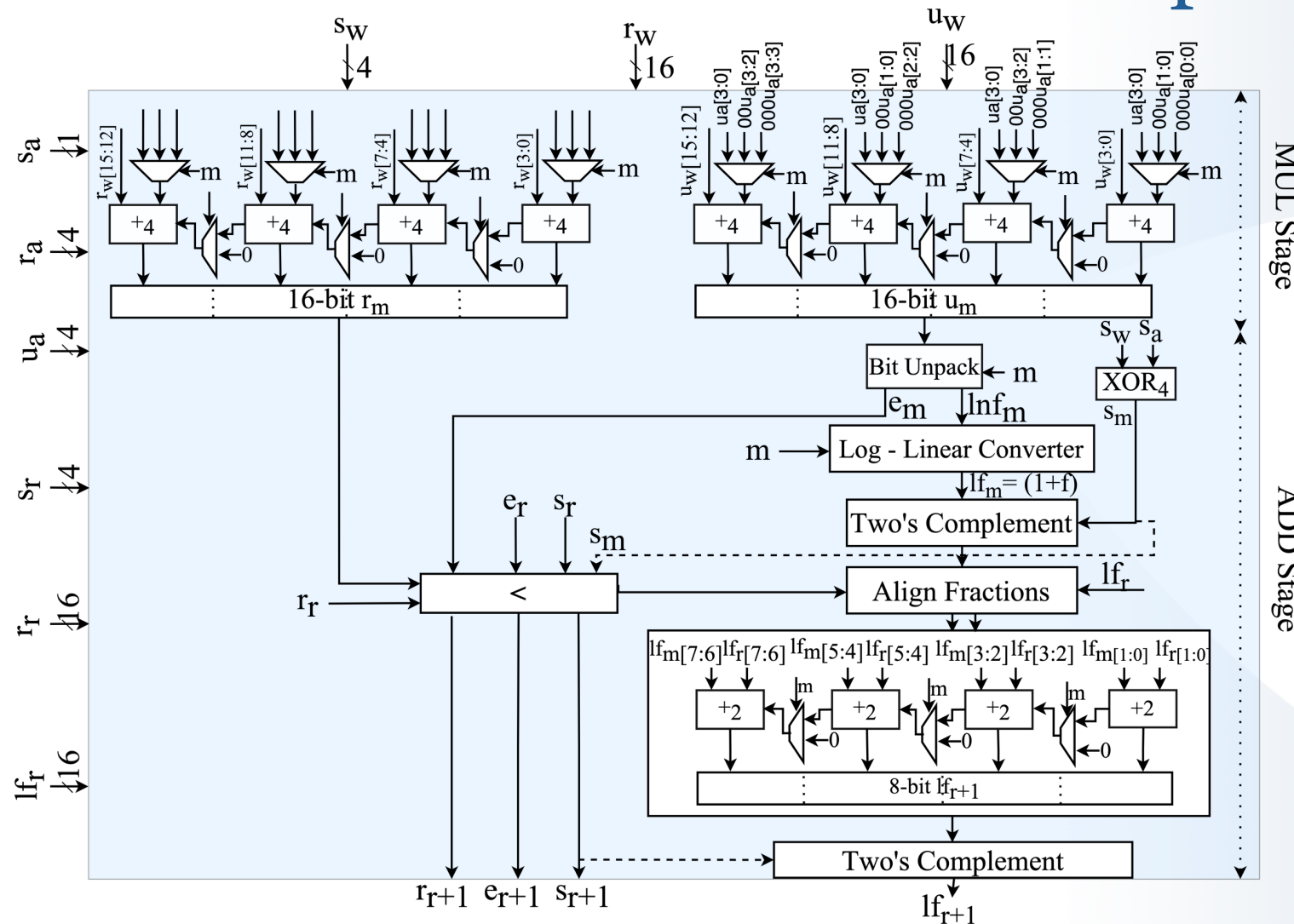Log-Linear Converter Implemented as Combinational Logic!
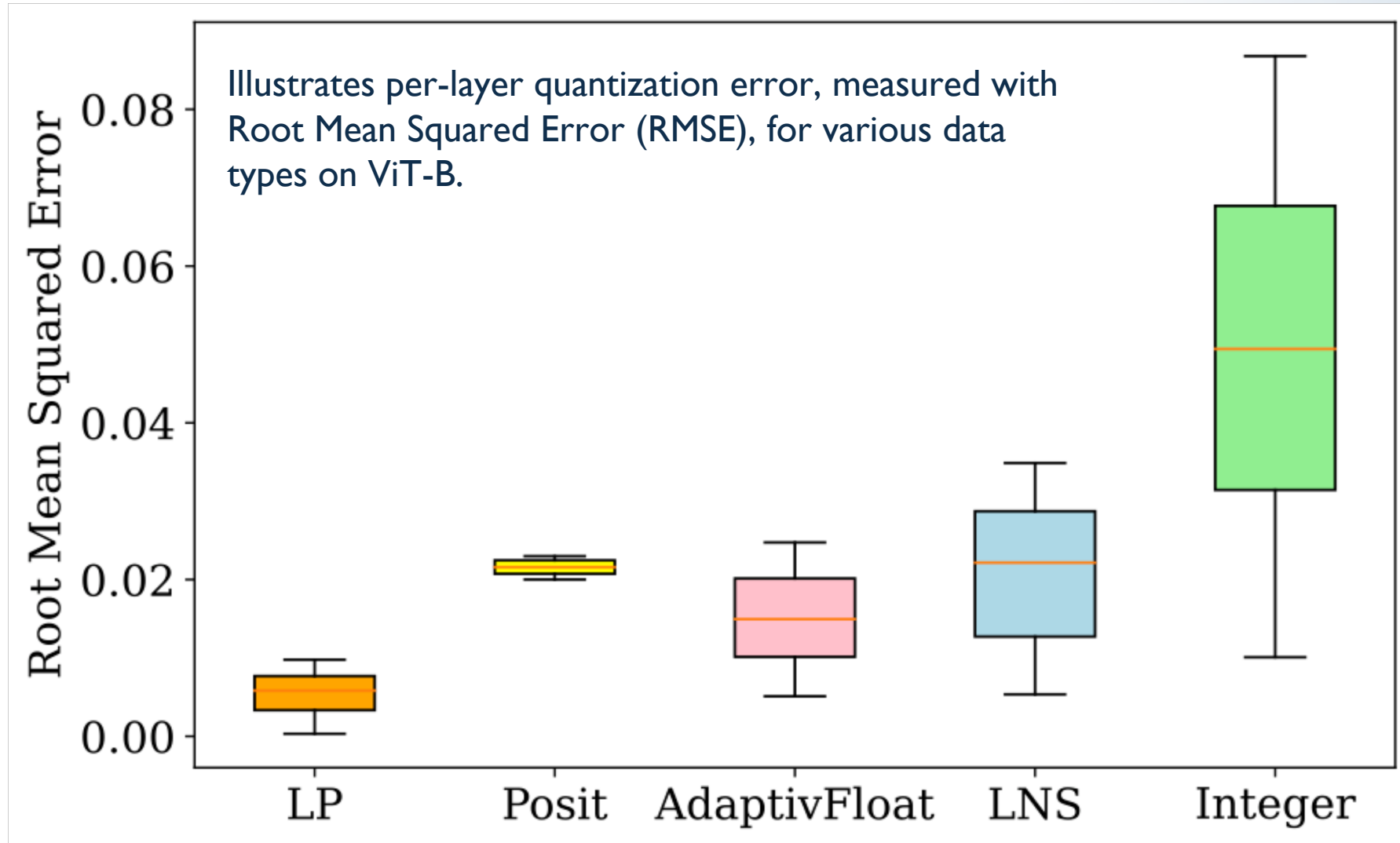
# Mixed-Precision PE Architecture



Linear Domain Addition made up of 2-bit adders.

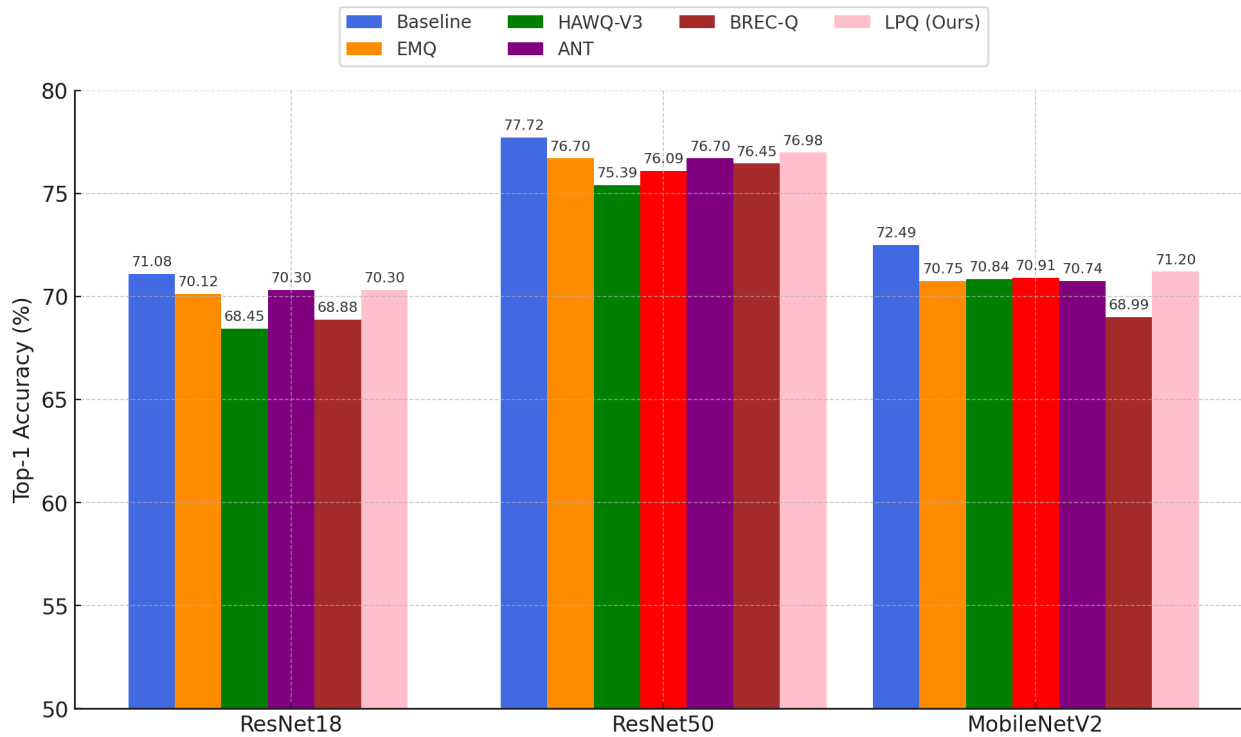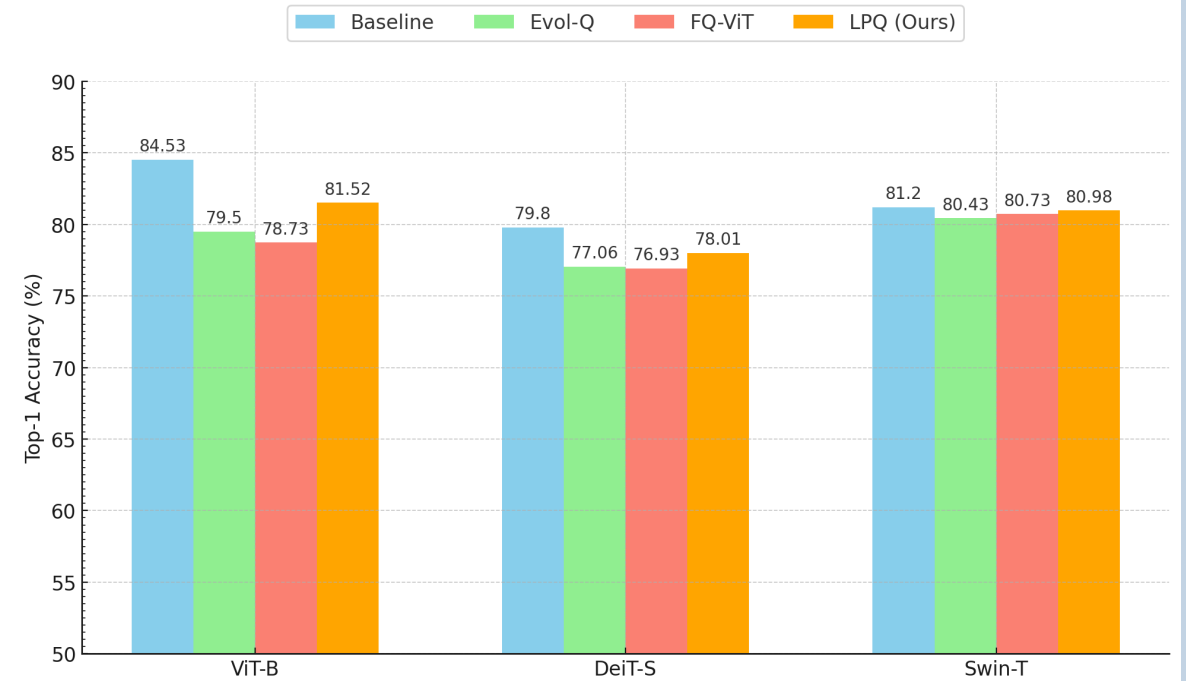# Mixed-Precision PE Architecture: Complete Flow

# Number Format Comparison



Illustrates per-layer quantization error, measured with Root Mean Squared Error (RMSE), for various data types on ViT-B.
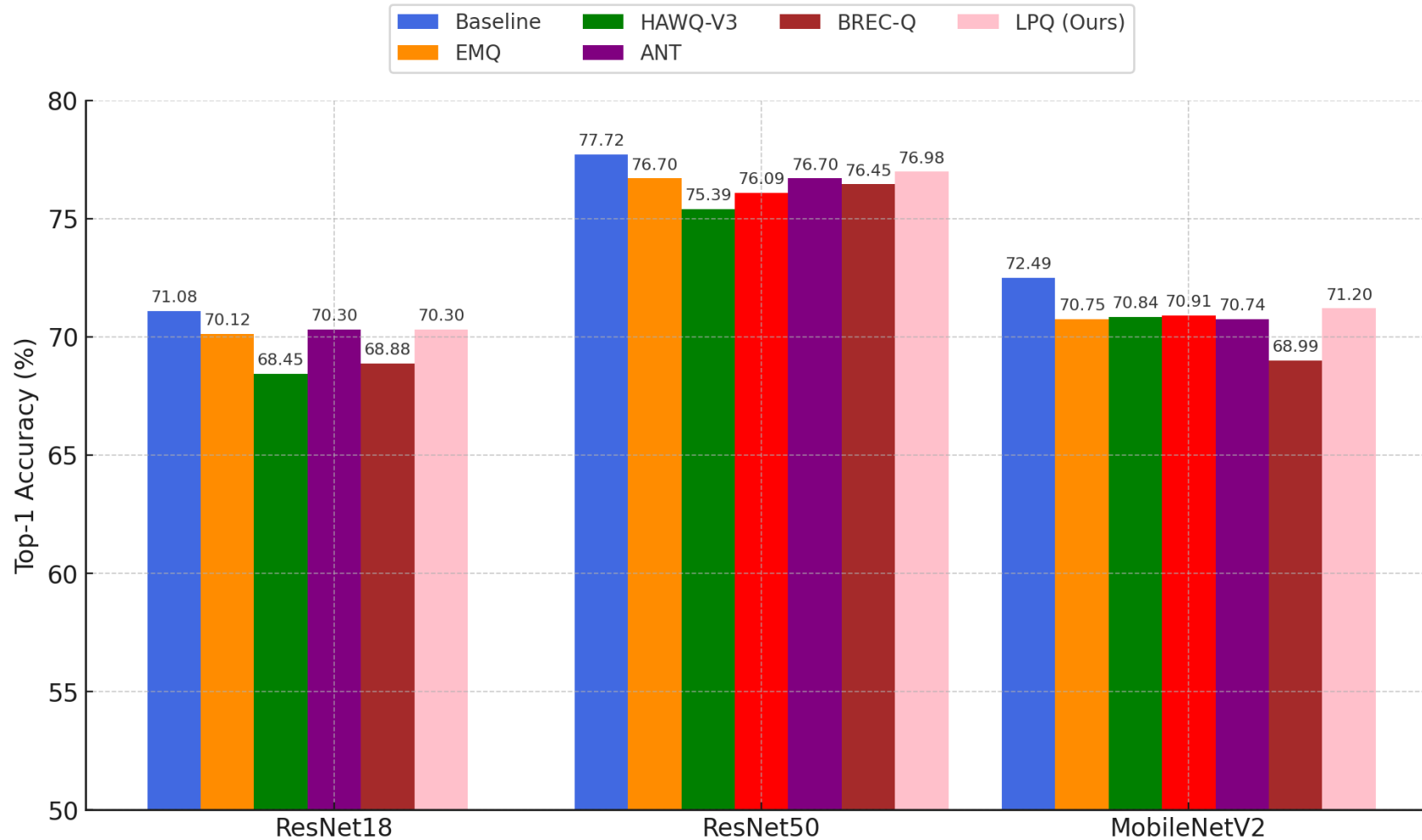
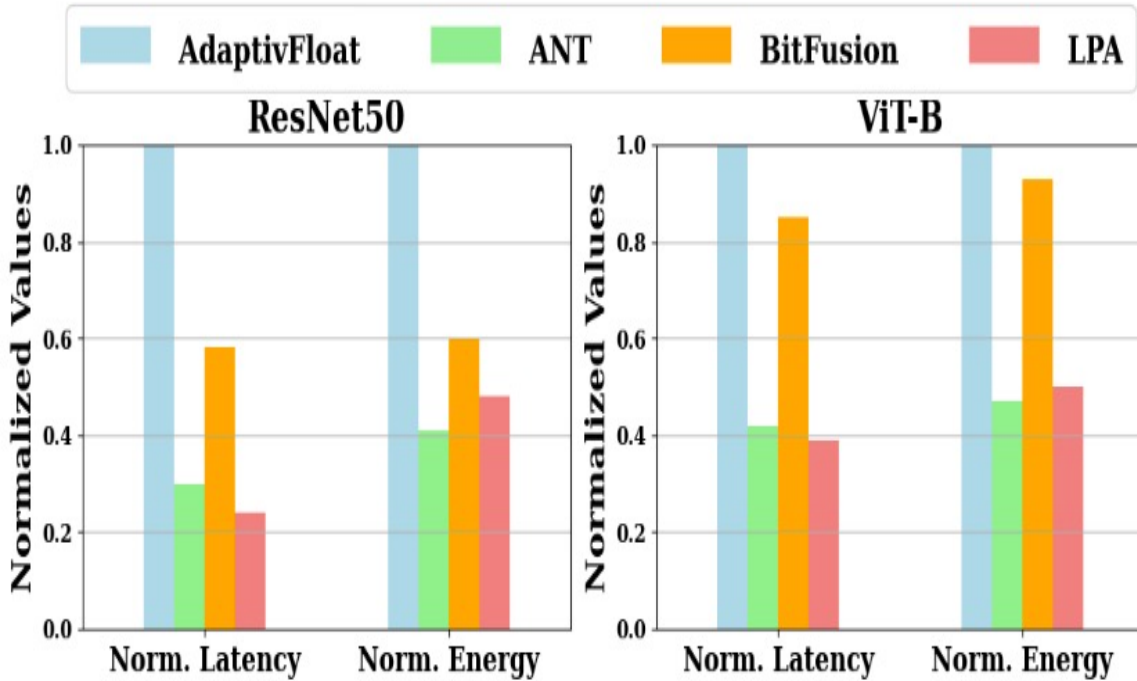# Mixed-Precision Quantization Results: CNN & ViT

# Mixed-Precision Quantization Results: CNN & ViT



- **Highlights:**

- On average <1% accuracy degradation compared to FP baseline.

- Mixed-Precision LPQ achieves average weight/activation bitwidth of 4.2/5.5.

- 90% reduction in model size.
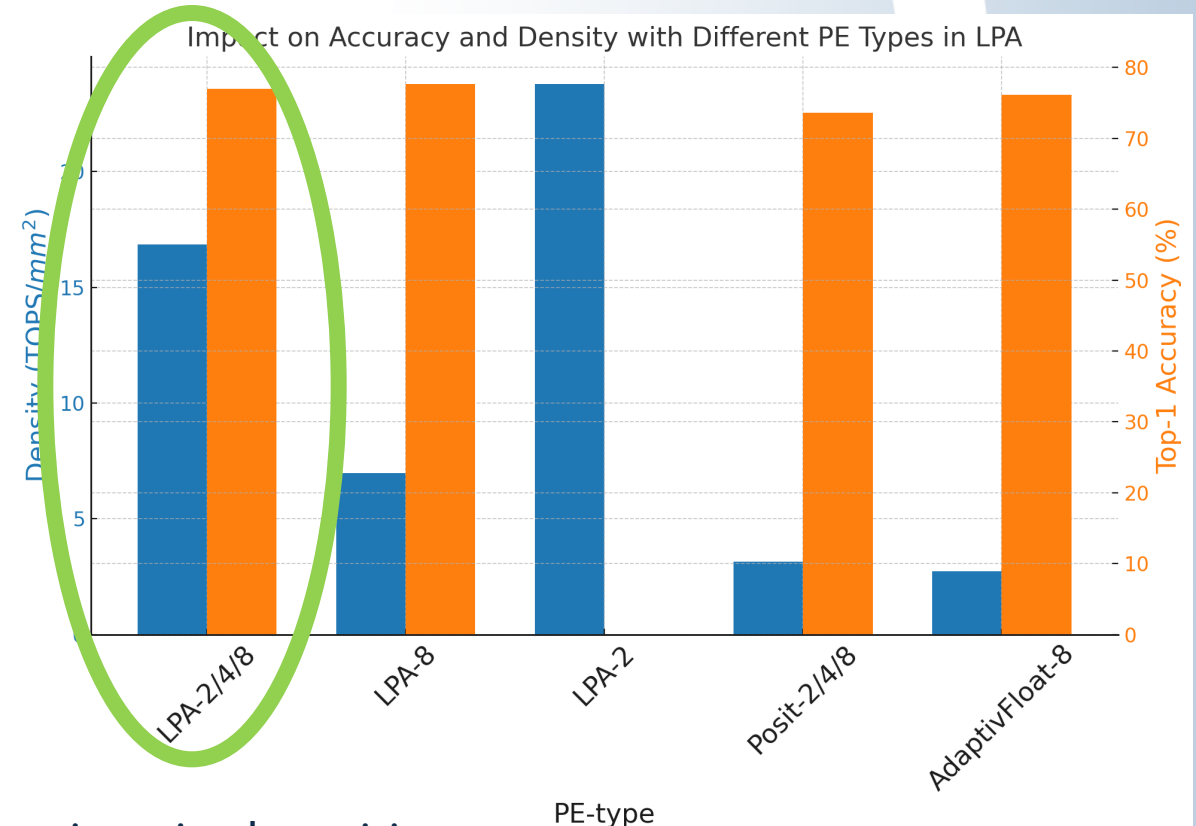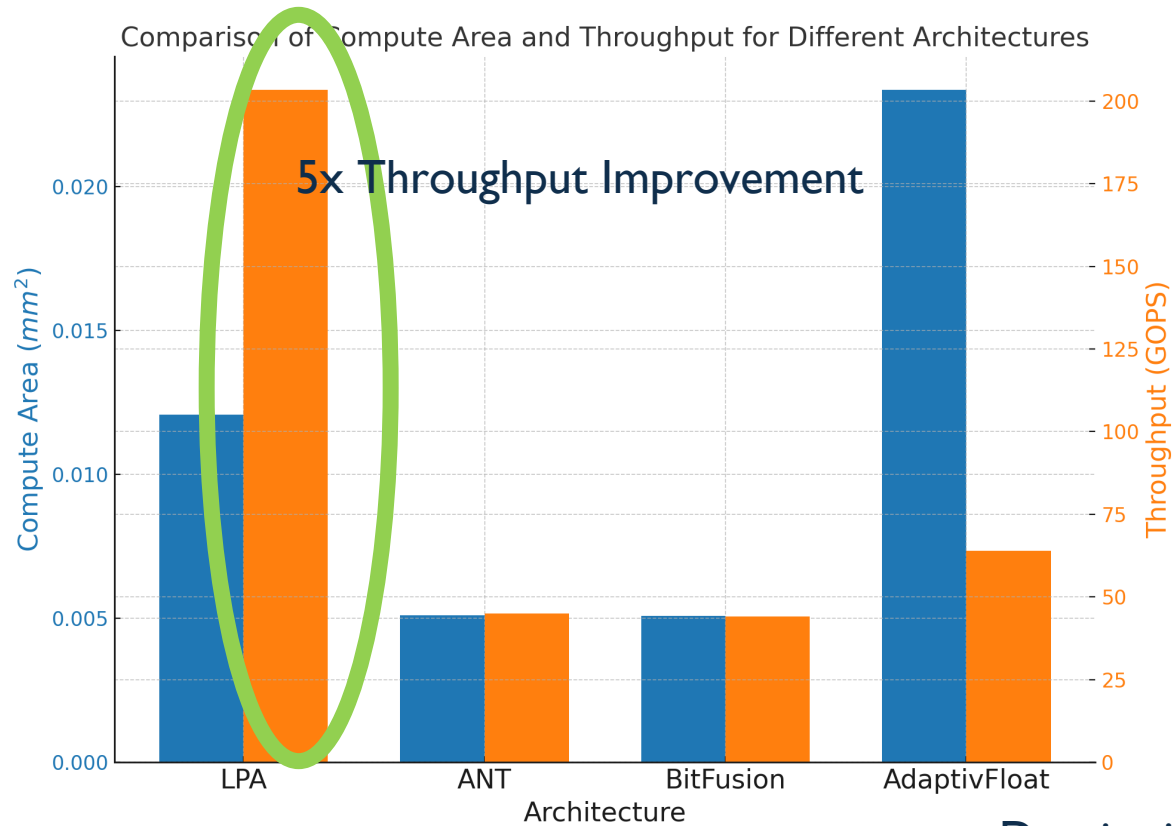
# Logarithmic Posit Accelerator Performance



LPA exhibits the lowest latency across models, with a modest increase in energy consumption over ANT attributed to overheads due to native mixed- precision support and conversion logic.

| Architecture | Component (Area) | Compute Area $(\mu m^2)$ | Throughput (GOPS) | Compute Density (TOPS/$mm^2$) | Total Area $(mm^2)$ |
|---|---|---|---|---|---|
| LPA | Decoder(5.2 $\mu m^2$) Encoder (9.4 $\mu m^2$) 2/4/8-bit PE (187.43 $\mu m^2$) | 12078.72 | **203.4** | **16.84** | 4.212 |
| ANT | Decoder(4.9 $\mu m^2$) 4/8-bit Int PE (79.57 $\mu m^2$) | 5102.28 | 44.95 | 8.81 | **4.205** |
| BitFusion | 2/4/8-bit PE | **5093.75** | 44.01 | 8.64 | **4.205** |
| AdaptivFloat | 8-bit PE | 23357.14 | 63.99 | 2.74 | 4.223 |

Despite ANT and BitFusion exhibiting lower area when compared with LPA for the same number of PEs, LPA results in proportionately higher performance per unit area (TOPS/$mm2$) for mixed-precision DNN inference.

# Logarithmic Posit Accelerator Performance



Comparison of Compute Area and Throughput for Different Architectures

5x Throughput Improvement

Impact on Accuracy and Density with Different PE Types in LPA

Despite incorporating mixed-precision support, LPA-2/4/8 achieves accuracy tending to the ideal scenario for both metrics, demonstrating a balanced trade-off.

SRC

# Summary

## CONTEXT

Current Approach and Challenges

- Approach: Algorithm-Hardware Co-Design is a promising area towards efficient DNN inference.
- Challenges: Inefficient data formats and lack of generalizable automated techniques.

## APPROACH

The proposed approach:

- Develop an adaptive, hardware-friendly data format.
- Identify existing limitations of automated quantization algorithms.
- Optimize existing hardware to support next-generation data formats.

## IMPACT

How do current results advance SOTA?

- <1% Accuracy drop across DNN model families post-quantization with > 15% higher compression ratio than competing methods.
- 2x improvement in PPA and energy-efficiency.

## Next Steps

- Extend: Verify adaptability and generalizability to LLMs and VLMs.
- Profile: Deepen understanding of data distribution of modern DNN models to improve existing data-format parameterization.
- Improve runtime of existing algorithm for faster quantization of large-scale models.

SRC