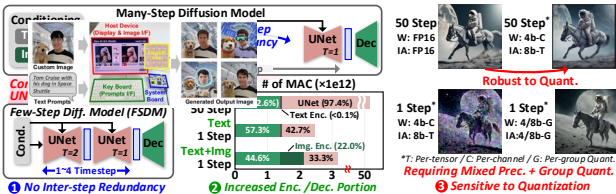


# EdgeDiff: Multi-modal Few-step Diffusion Model Accelerator with Mixed-Precision and Reordered Group-Quantization for On-device Generative AI

Sangjin Kim, Jungjun Oh, Jeonggyu So, Yuseon Choi, Sangyeob Kim, Dongseok Im, Gwangtae Park, and Hoi-Jun Yoo

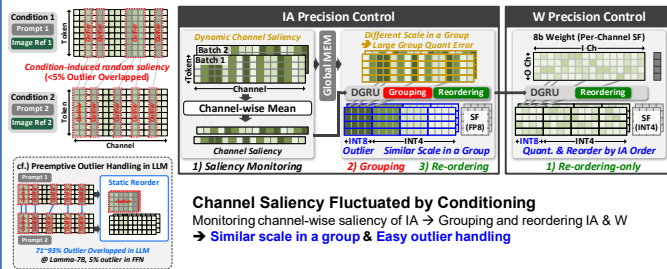
## Motivation

Trending Algorithm Opt. for DM: Few-Step Diffusion Model with Knowledge Distillation



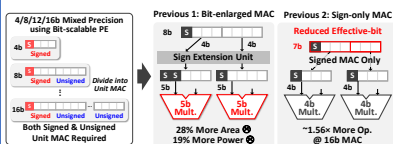
EdgeDiff: Accelerate All Encoder/Unet/Decoder of FSDM with Mixed Precision + Group Quant.

## Proposed Condition-aware Reordered Mixed Precision

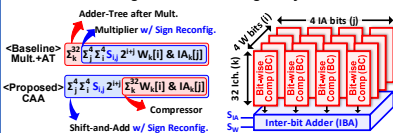


## Proposed HW Opt. for Mixed Precision and Group Quantization

### Micro-arch. for Mixed Precision MAC



Unit MAC should support both signed & unsigned MAC  
 Previous: Bit-enlarged MAC units or sign-only MAC units

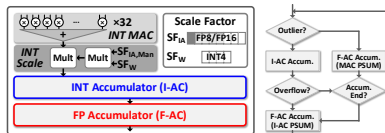


Reducing Overhead by Changing Accum. Order  
 Baseline: Inter-bit (mult. w/ sign reconfig. ⊗) → Inter-ch. (add)  
 CAA: Inter-ch. (compressor) → Inter-bit (add w/ sign reconfig. ⊗)

### Micro-arch. for Group Quantization

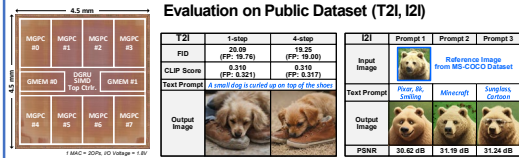


FP scaling factor (SF) → Power-intensive FP operation



Tiering into INT-Accum. & FP-Accum.  
 (Tier 1: I-AC) Checking Range & In-range Data Accumulation  
 (Tier 2: F-AC) Conditional Operation for Power Reduction

## Implementation Results (Chip, System)



### Demonstration System with Custom Image

