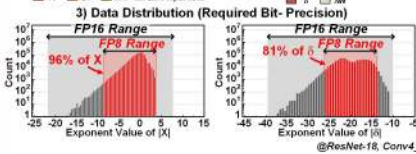


LNPU: An Energy-Efficient Deep-Neural-Network Training Processor with Fine-Grained Mixed Precision

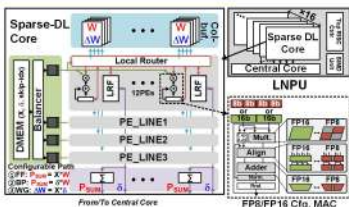
Jinsu Lee (jinsulee@kaist.ac.kr), Juhyoung Lee, Donghyeon Han, Jimmook Lee, Gwangtae Park and Hoi-Jun Yoo

DNN Training Analysis

- 1) ~43% of Operation is Zero Input MAC Operation
- 2) X and δ occupy ~78% of External Memory Access
- 3) FP8 Cover ~80% of overall data @ ResNet-18

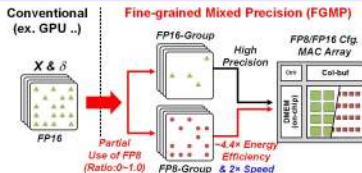


Overall Architecture

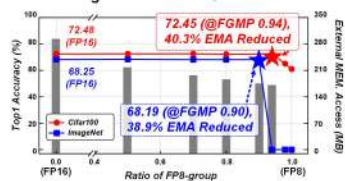


- 1) Fine-Grained Mixed Precision of float16 / float8
 - Improve Performance & Lower EMA
- 2) Sparse DL Core
 - Skip Zero Input MAC Op. caused by FGMP
- 3) Central Core
 - Global Mem. Manage for Irregular Core latency

Proposed Fine-Grained Mixed Precision

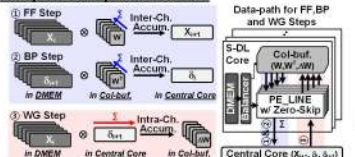


ResNet-18 Training Results with FGMP

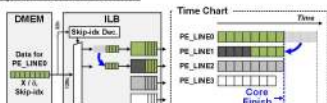


Sparse DL Core & Central Core

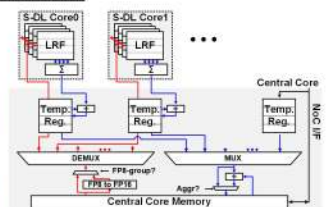
Datapath for Sparse DL Core



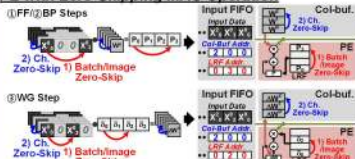
Input Load Balancer



Central Core



PE with Zero-Skipping MAC Operation



Implementation Results

Chip Photo & Summary

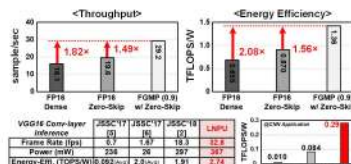


Specifications	
Technology	55nm 1P8M CMOS
Die Area	4mm x 4mm (16mm ²)
GRAM	372 KB
Supply Voltage	0.7V - 1.1V
Frequency	< 200MHz
Date Type	FP8, FP16
Power Consumption	43 mW @ 0.7V, 50MHz 367mW @ 1.1V, 200MHz
Power Efficiency (TFLUPS/W)	3918 / 1798

Autonomous Drone System w/ LUPU



Performance & Comparisons



*Effective TFLUPS/W with 80% input sparsity