



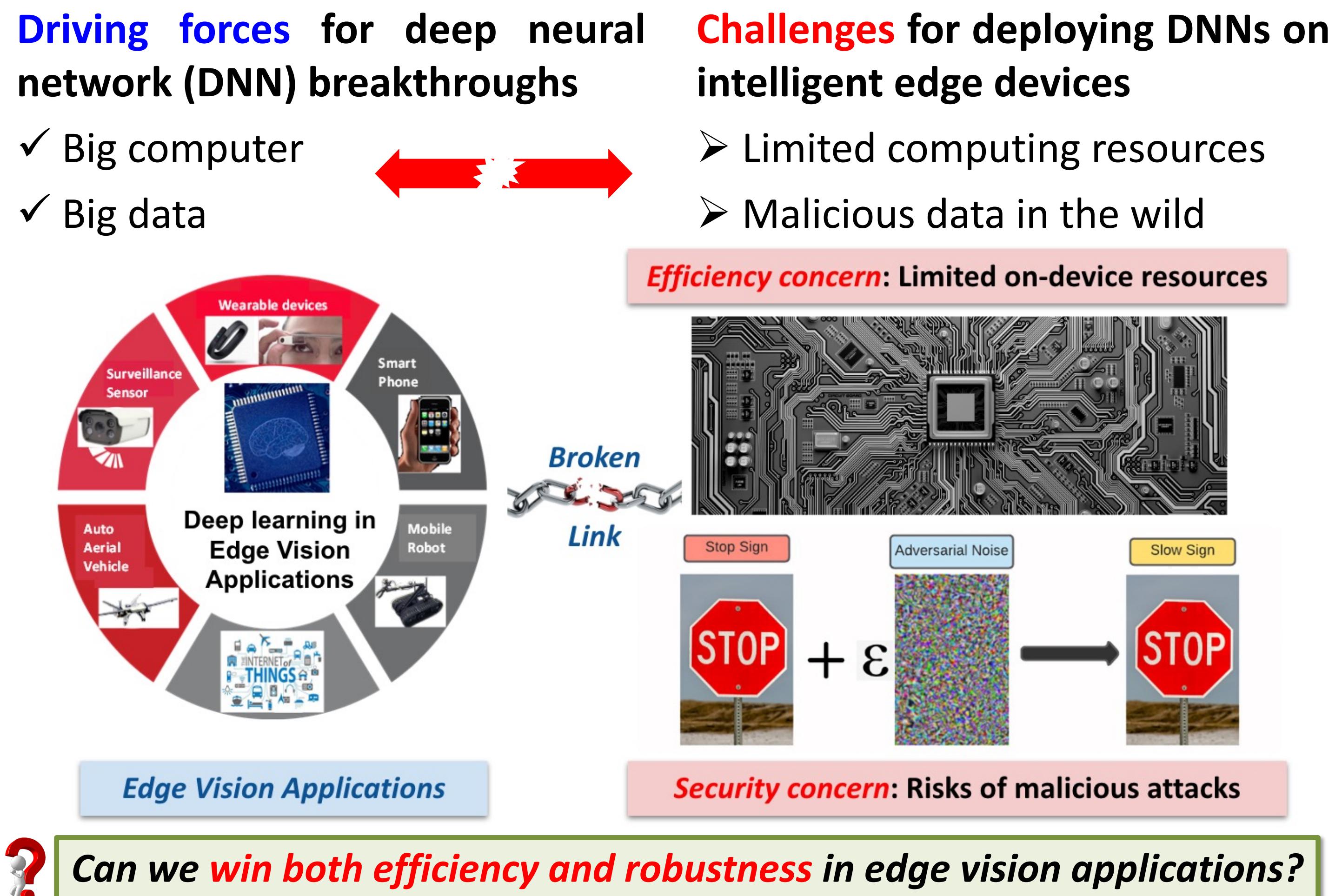
2-in-1 Accelerator: Enabling Random Precision Switch for Winning Both Adversarial Robustness and Efficiency

Yonggan Fu, Yang Zhao, Qixuan Yu, Chaojian Li, and Yingyan (Celine) Lin

Georgia Institute of Technology

Accepted by MICRO 2021

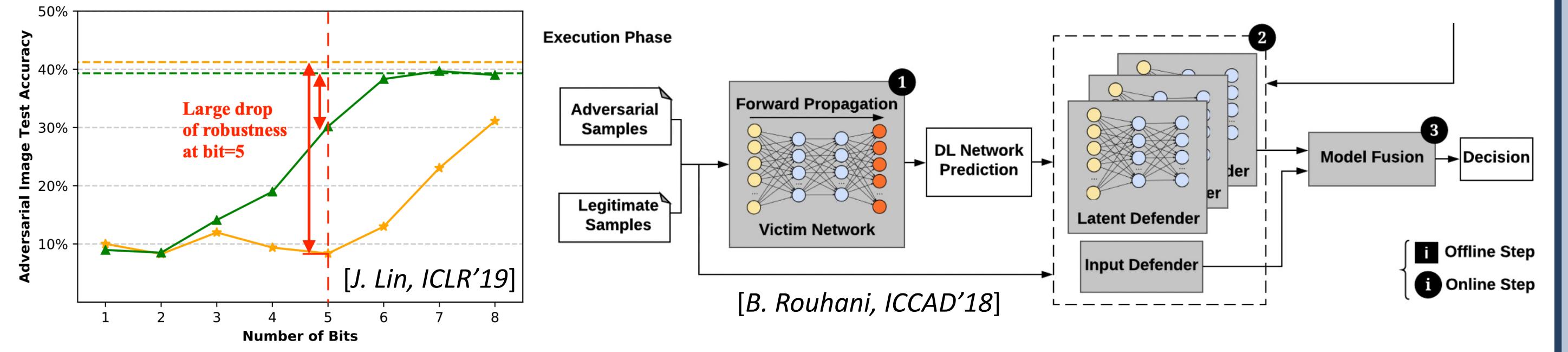
Background & Motivation



Limitations of Previous Works

Algorithm side: Robustness and efficiency are often contradictory

- Overparameterization is crucial for robustness [A. Madry, ICLR'18]
- Quantized DNNs are inferior in robustness [J. Lin, ICLR'19]



Hardware side: Robustness is achieved at the cost of efficiency

- Extra detection modules are needed: Deepfense [B. Rouhani, ICCAD'18], DNNGuard [Z. Wang, ASPLOS'20], Ptolemy [Y. Gan, MICRO'20]

2-in-1 Accelerator: Overview

2-in-1 Accelerator: Algorithm-Architecture Co-Design to Win both Efficiency and Adversarial Robustness

Algorithm: Random Precision Switch to enable random DNN quantization as an in-situ model switch

Architecture: A novel precision-scalable accelerator featuring a spatial-temporal MAC array with bit-level optimization

2-in-1 Accelerator: Algorithm

Motivating observation: Poor Adv. Transferability between Precisions

Visualize the transferability

- Robust accuracy measured under varied attack and inference precision pairs

Key observation

- Consistently poor adversarial transferability btw. precisions across different adversarial training methods

Key insight: Adversarial perturbations are shielded by quantization noises which cannot be effectively learned by gradient-based attacks

The Inspired Technique: Random Precision Switch (RPS)

Technique 1: RPS Inference

- Randomly select one inference precision
- ✓ Relatively stable natural accuracy
- ✓ Degraded attack effectiveness

Technique 2: RPS Training

- Equip the model with switchable BN
- ✓ Increase the difficulty of adv. transfer
- Randomly select one precision for training

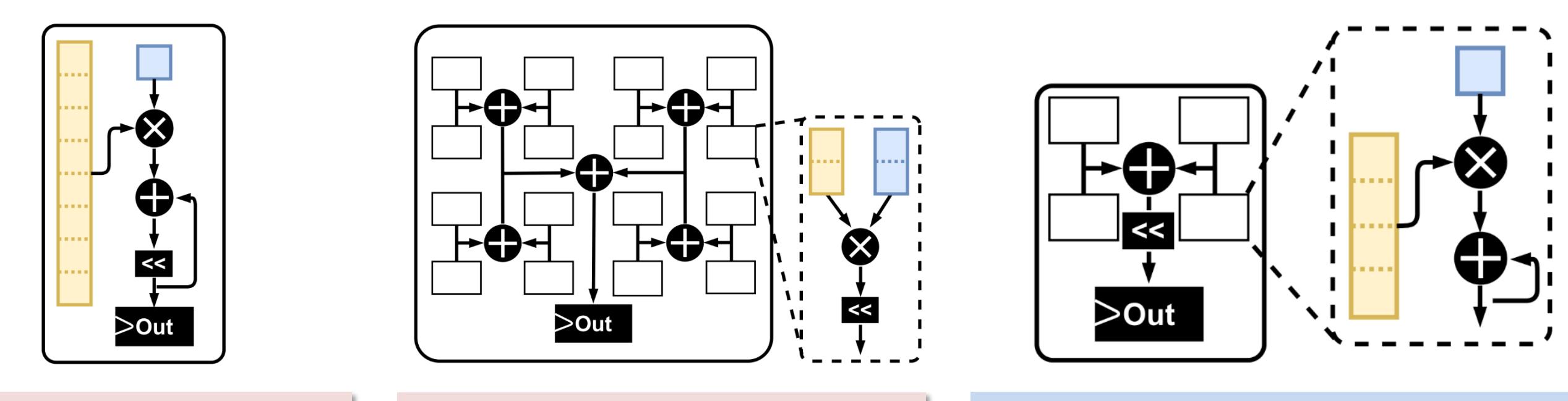
```

1: === RPS Training ===
2: Equip  $f_\theta$  with SBN
3: for epoch  $\in [1, T]$  do
4:   for  $(x, y) \in D_{train}$  do
5:     Randomly select a precision  $q$  from  $Set_Q$ 
6:     Obtain  $f_\theta^q$  by quantizing  $f_\theta$  to  $q$ -bit
7:      $\delta = 0$  or random initialized
8:     for  $t \in [1, 7]$  do
9:        $\delta = \text{clip}_\epsilon(\delta + \alpha \cdot \text{sign}(\nabla_\theta f_\theta^q(x + \delta), y))$ 
10:    end for
11:    $\theta = \theta - \nabla_\theta f_\theta^q(x + \delta), y)$ 
12: end for
13: end for
14: === RPS Inference ===
15: for  $x_{adv} \in D_{adv}$  do
16:   Randomly select a precision  $q$  from  $Set_Q$ 
17:   Obtain  $f_\theta^q$  by quantizing  $f_\theta$  to  $q$ -bit
18:   Evaluate  $\hat{y} = f_\theta^q(x_{adv})$ 
19: end for

```

2-in-1 Accelerator: Architecture

2-in-1 Accelerator's Architecture: A Spatial-Temporal MAC Array

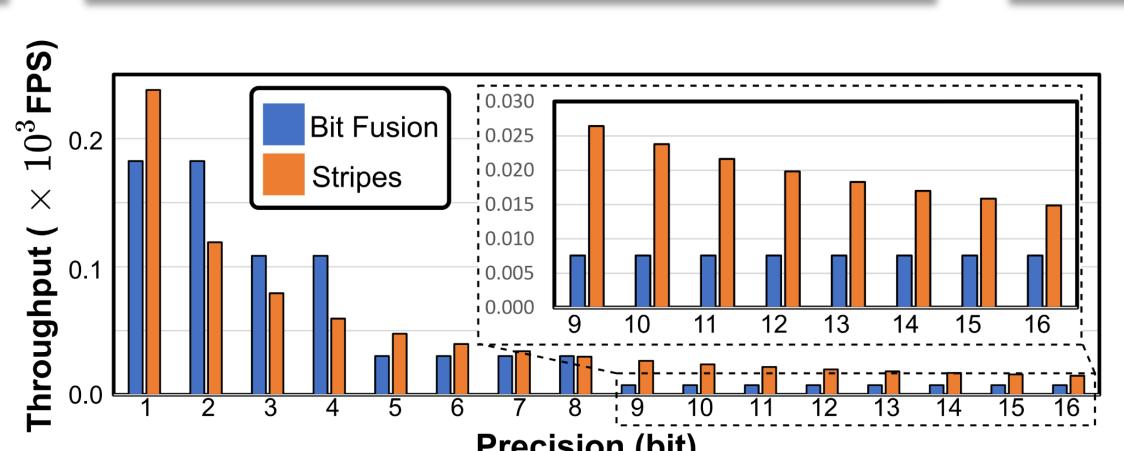


Inferior area-efficiency: Shift-add can take 90% of the total area in a 16-bit MAC unit

Inferior flexibility: Only support a limited set of precisions, e.g., 2/4/8-bit in Bit Fusion [H. Sharma, ISCA'18]

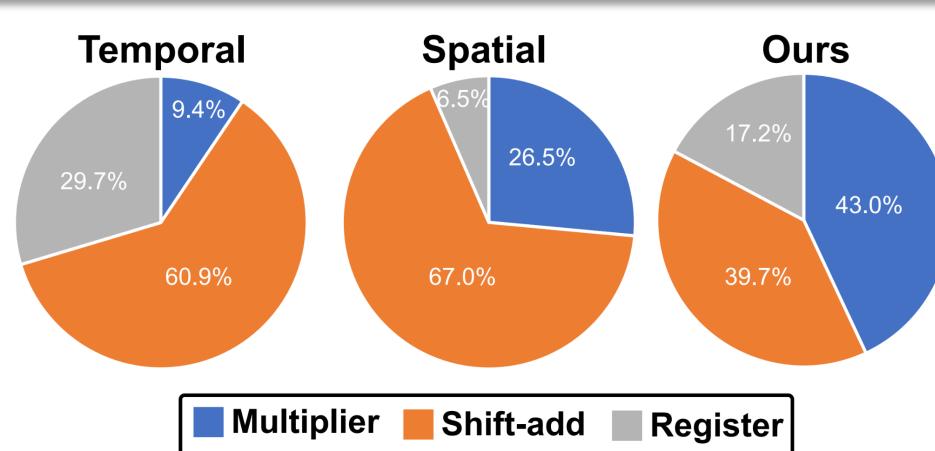
The dilemma between performance vs. flexibility

SOTA temporal design Strips [P. Judd, MICRO'16] vs. spatial design Bit Fusion



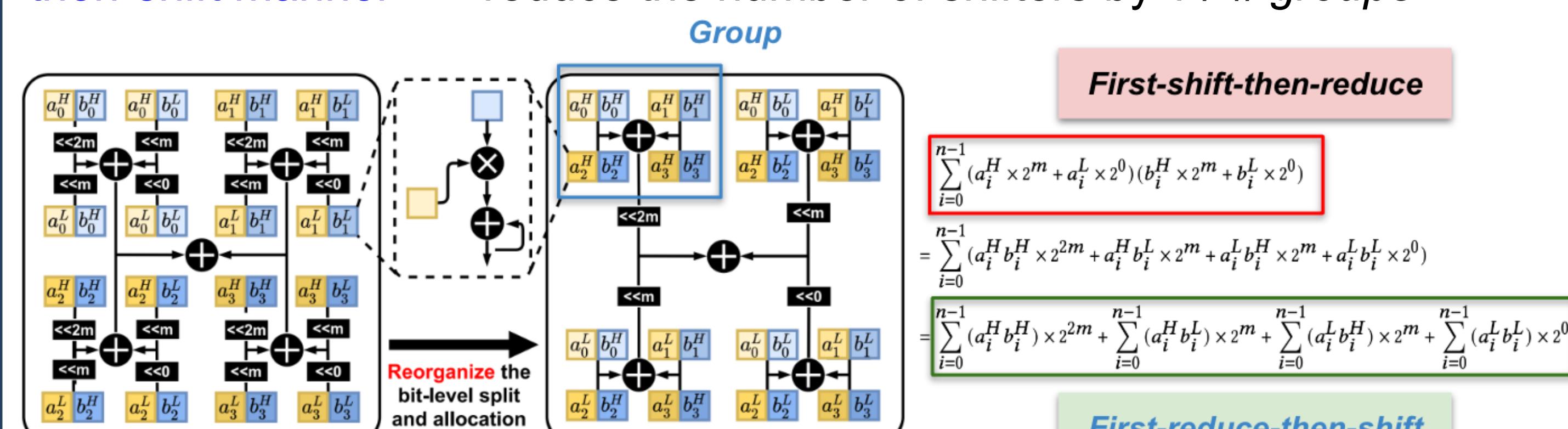
Better area-efficiency: Bit-serial units only need to support up to 4-bit

High flexibility: Support irregular precision choices like 3/6-bit

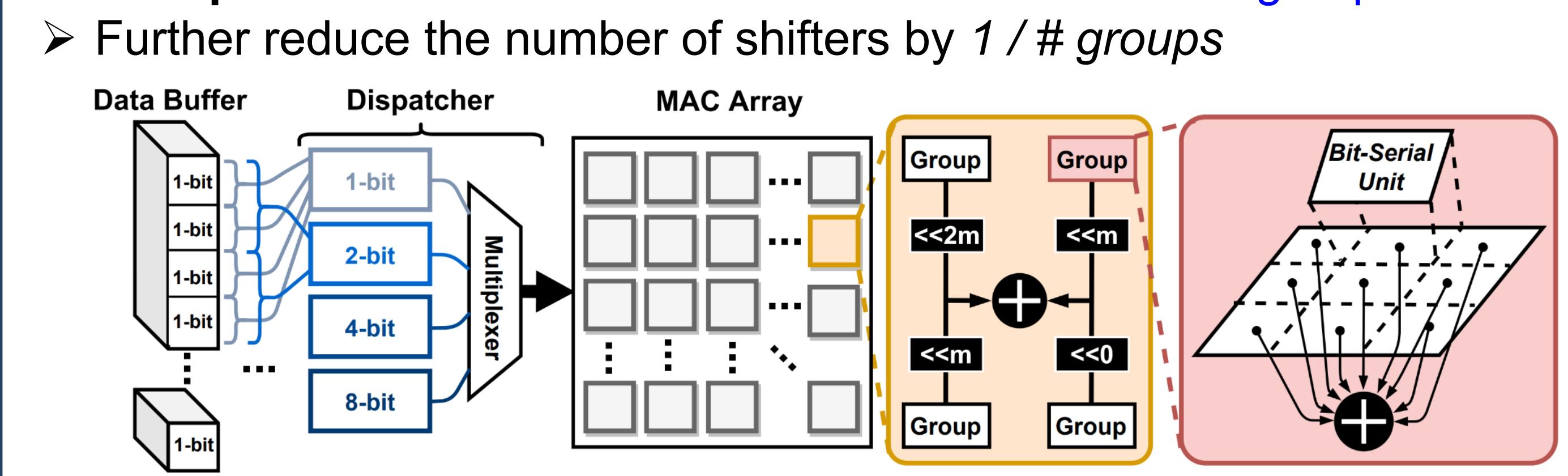


2-in-1 Accelerator's Architecture: Bit-level Reorganization & Fusion

Technique 1: Reorganize the bit-level split and allocation in a first-reduce-then-shift manner → reduce the number of shifters by 1 / # groups



Technique 2: Fuse the shift-add of bit-serial units in one group



Our proposed architecture: Achieve 2.3x throughput/area and 4.88x energy-efficiency/operation compared with Bit Fusion

2-in-1 Accelerator: Evaluation

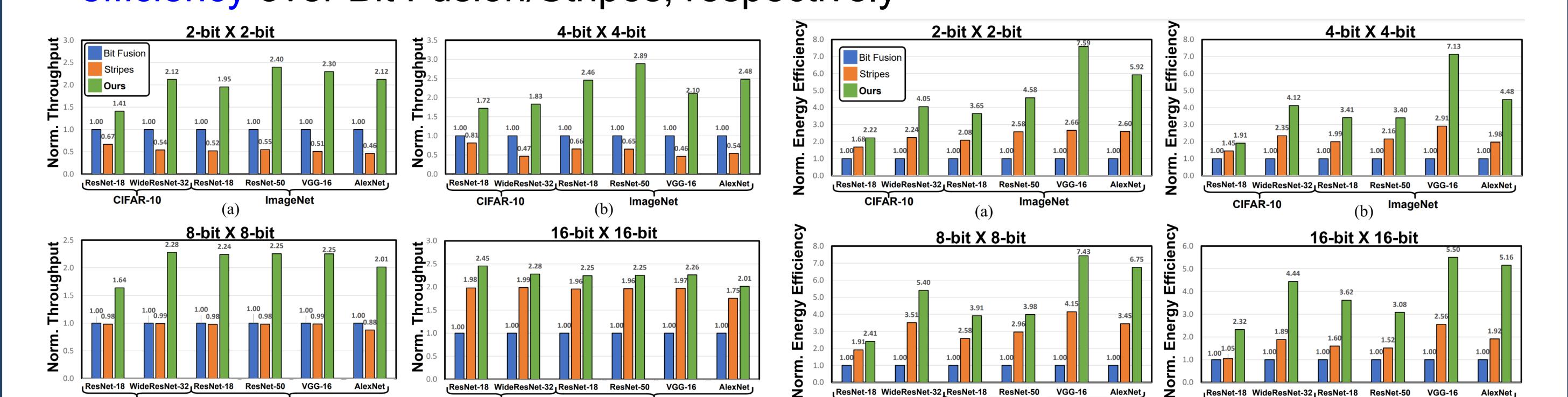
Evaluation 2-in-1 Accelerator's Algorithm

- ✓ Consistently boost adversarial robustness, e.g., +13.98%/+12.14% robust accuracy on two networks with PGD-7 training against PGD-20 attacks

Adversarial Training Method	PreActResNet-18			WideResNet-32		
	Natural (%)	PGD-7 (%)	PGD-100 (%)	Natural (%)	PGD-20 (%)	PGD-100 (%)
FGSM	67.04	41.48	41.37	66.76	40.78	40.55
FGSM + RPS	80.58	64.08	63.56	64.09	50.70	48.72
FGSM-RS	86.08	41.76	41.13	89.95	45.33	44.77
FGSM-RS + RPS	82.11	59.33	59.32	87.87	60.07	59.12
PGD-7	82.02	51.17	50.93	85.25	54.61	54.36
PGD-7 + RPS	82.16	65.15	64.88	81.52	66.75	66.28

Evaluation 2-in-1 Accelerator's Architecture

- ✓ 1.41x~2.88x/1.15x~4.59x throughput and 1.91x~7.58x/1.25x~2.85x energy efficiency over Bit Fusion/Stripes, respectively



Throughput evolution with precisions

- ✓ Superior performance: Outperform both baselines
- ✓ High Flexibility: Achieve consistent improvements in throughput as the precision decreases

