

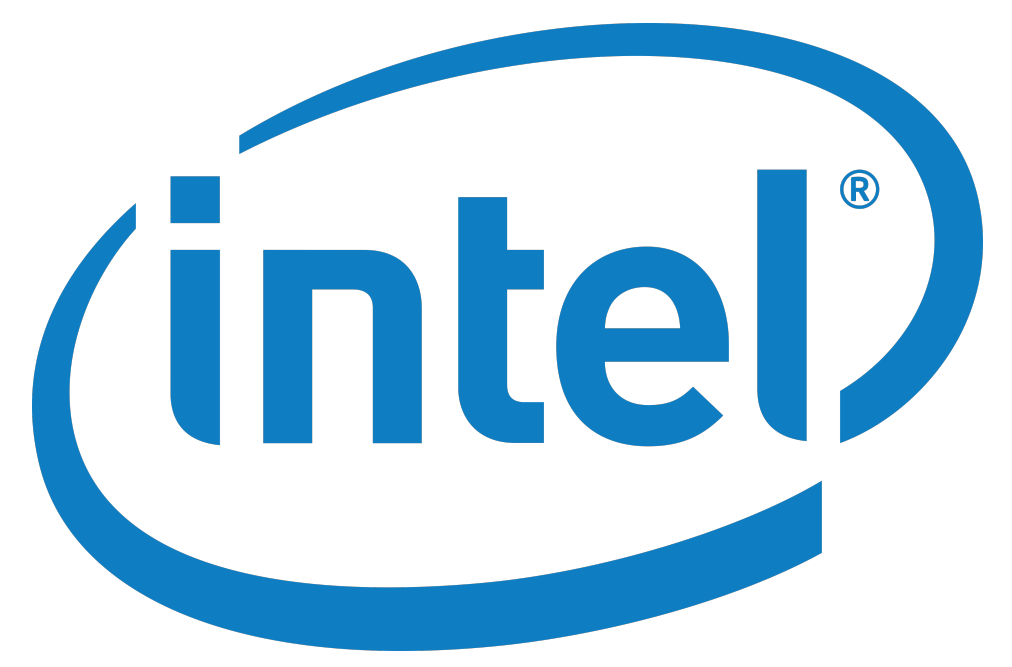


# Data-type Aware Arithmetic Intensity for Deep Neural Networks

Nandan Kumar Jha<sup>1</sup>, Sparsh Mittal<sup>1</sup>, Sasikanth Avancha<sup>2</sup>

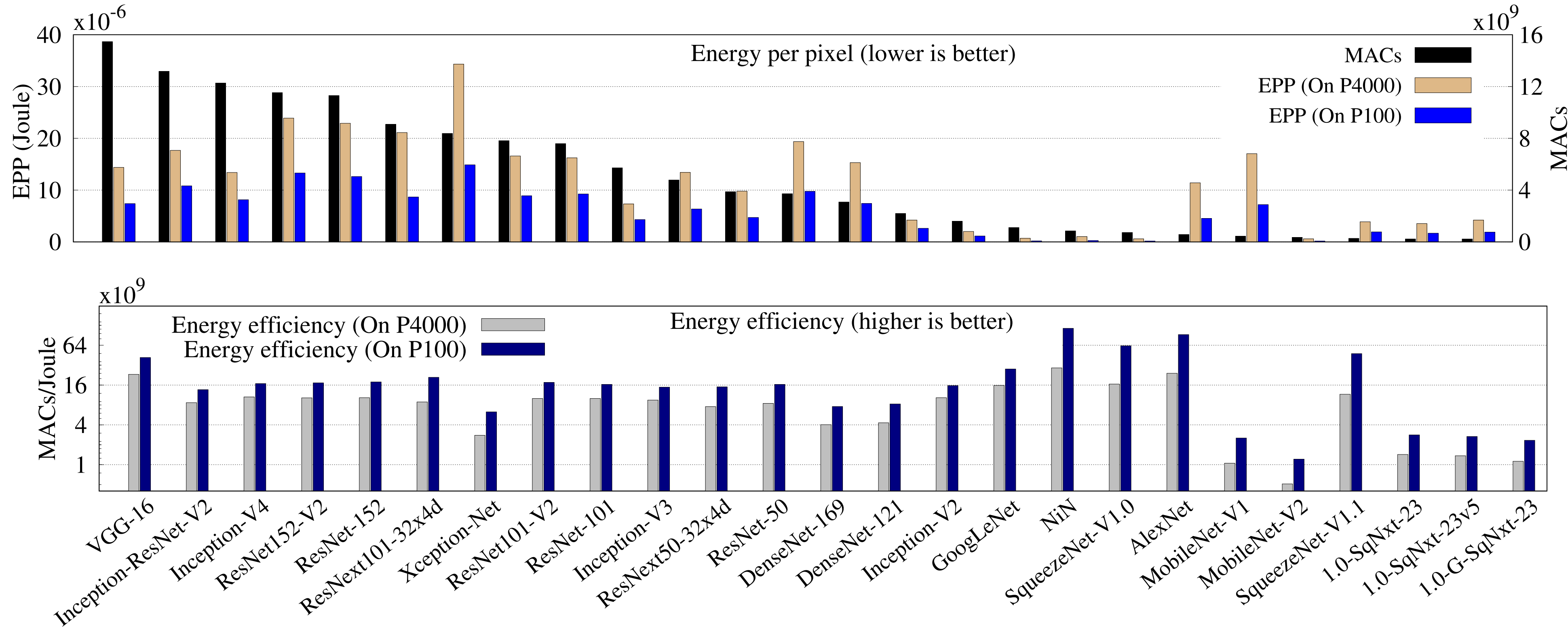
<sup>1</sup>Department of Computer Science and Engineering, IIT Hyderabad, India

<sup>2</sup>Parallel Computing Lab, Intel Corporation



## INTRODUCTION AND MOTIVATION

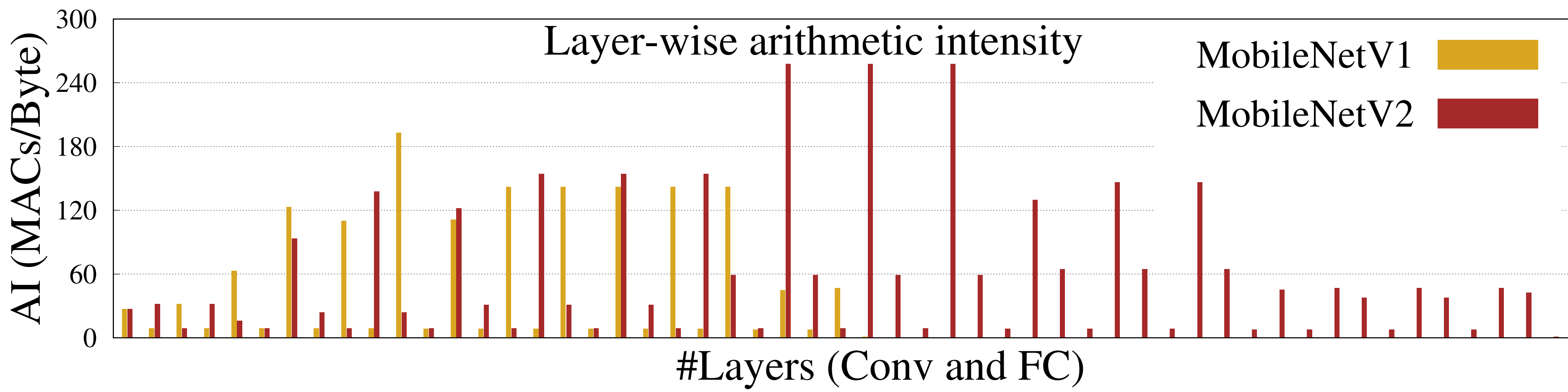
- Energy consumption of DNNs depend on computational complexity (#MACs) and energy efficiency of each MAC operation, which in turn depends on **data movement**.



## CHALLENGES

- DNN is a **special** type of workload which has computation phases with *radically* different computational intensity (OPs/Bytes). For example, depthwise, pointwise, groupwise convolution.

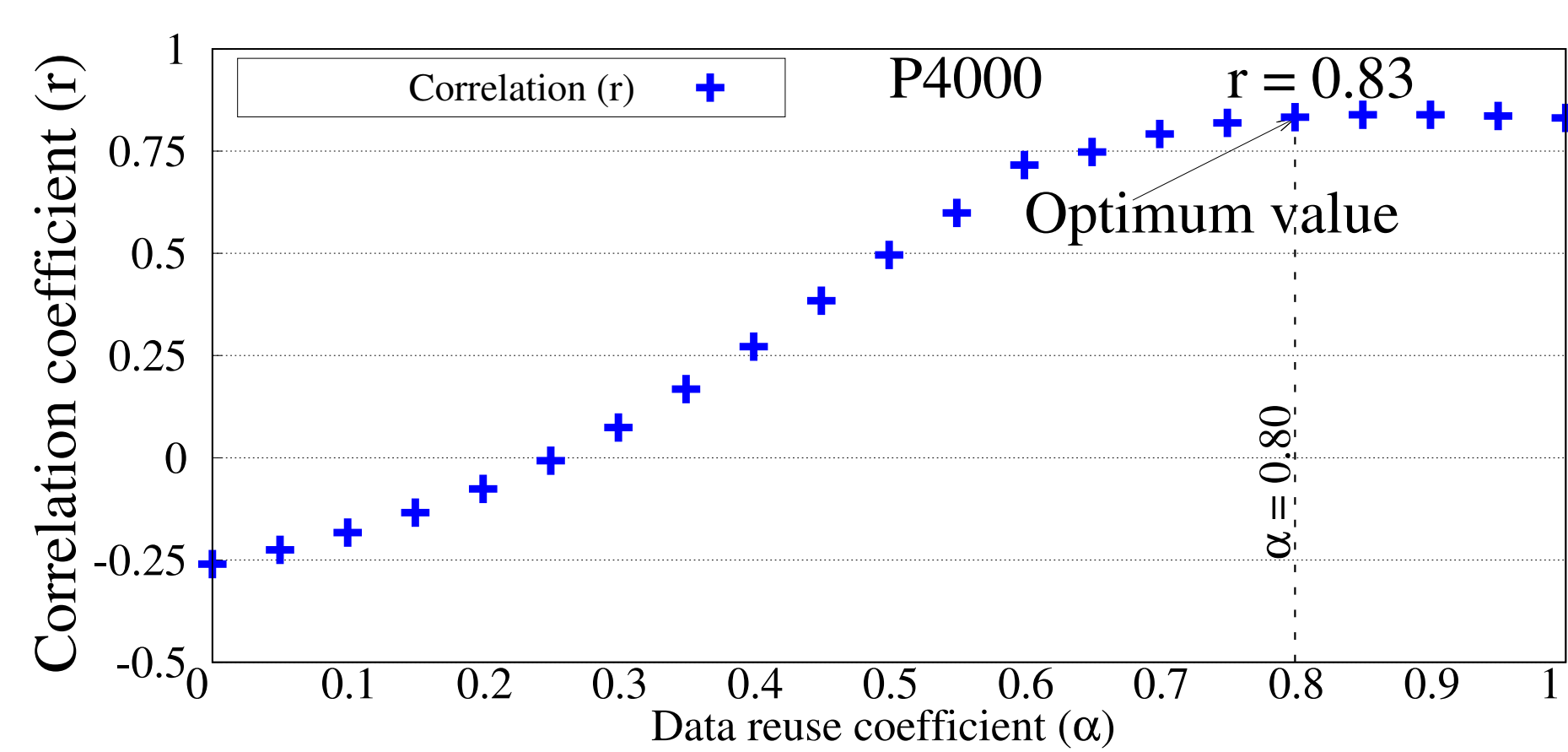
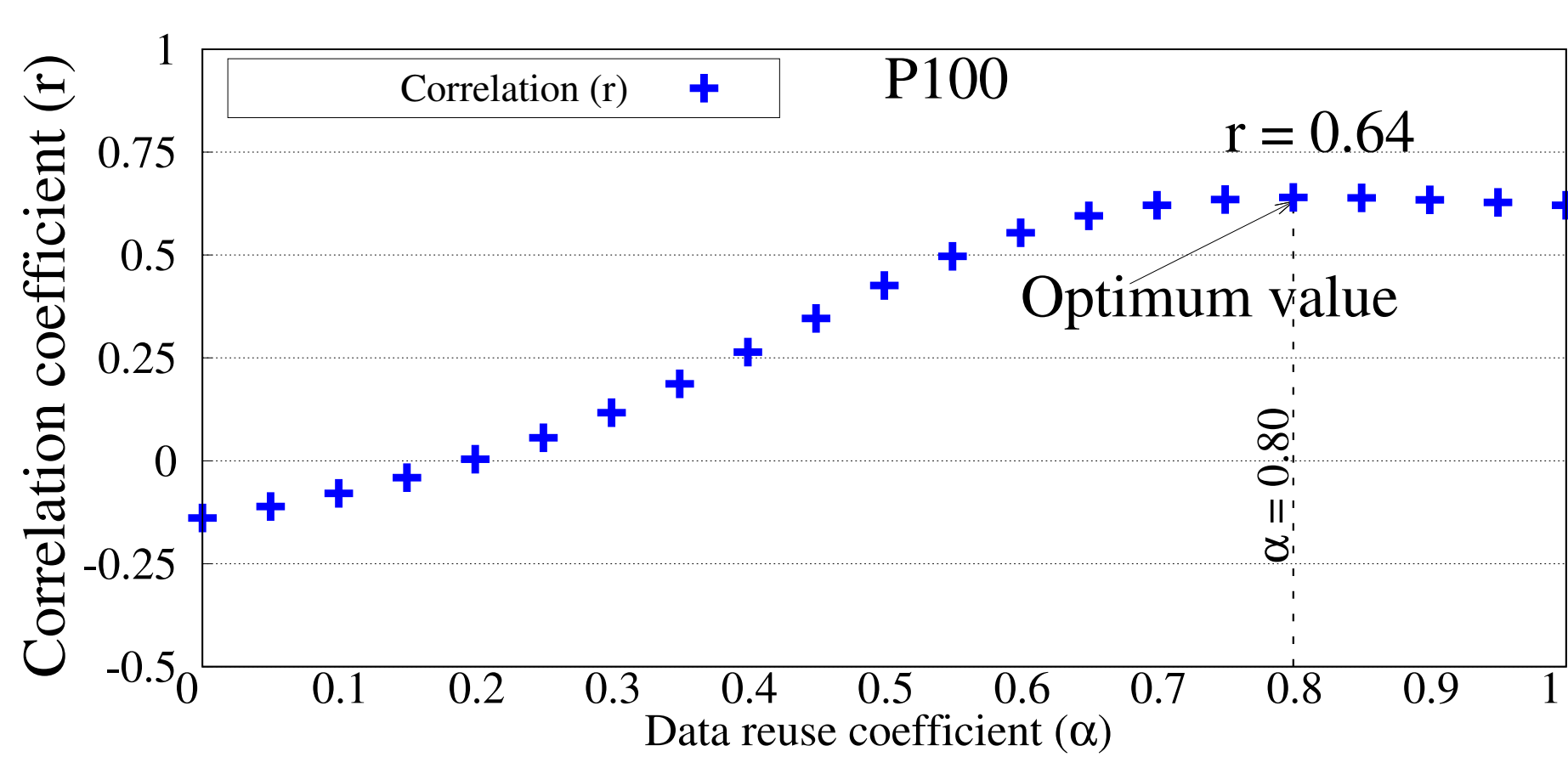
Convolution	Arithmetic intensity (AI)	AI_r
Standard	$\frac{M \times N \times S_k^2 \times S_o^2}{M \times N \times S_k^2 + (M+N) \times S_o^2}$	1.00
Pointwise	$\frac{M \times N \times S_o^2}{M \times N + (M+N) \times S_o^2}$	0.24
Group	$\frac{M \times N \times S_k^2 \times S_o^2}{M \times N \times S_k^2 + g \times (M+N) \times S_o^2}$	0.45
Depthwise	$\frac{M \times S_k^2 \times S_o^2}{M \times S_k^2 + (M+N) \times S_o^2}$	0.01



## PROPOSED METHOD

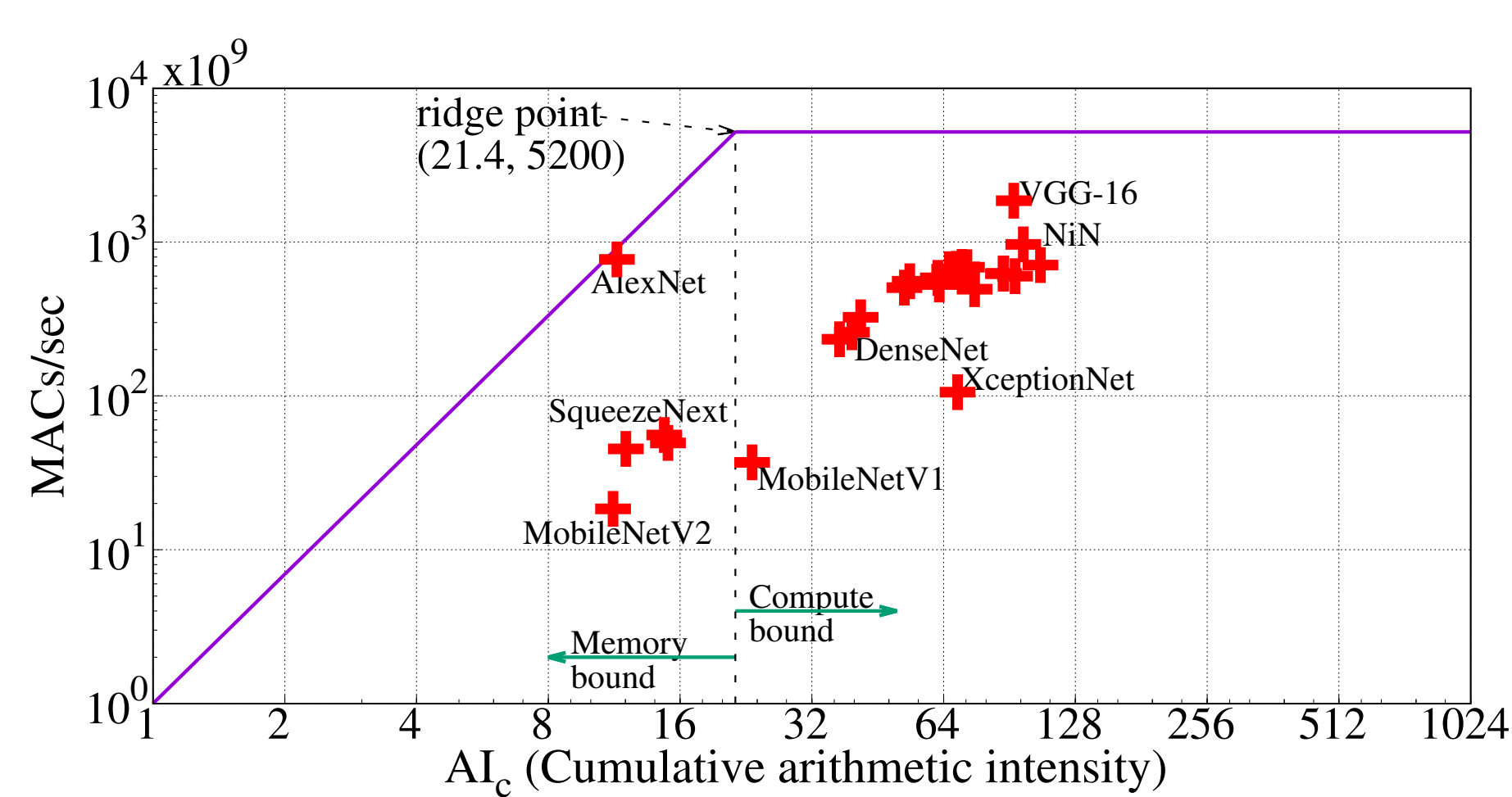
$$\frac{W+A}{2} \geq \frac{2 \times W \times A}{W+A} \Rightarrow \frac{2 \times M_c}{W+A} \leq \frac{M_c \times (W+A)}{2 \times W \times A} \Rightarrow \frac{M_c}{W+A} \leq \frac{1}{4} \times \left[ \frac{M_c}{A} + \frac{M_c}{W} \right]$$

$$\Rightarrow AI_c \leq \frac{1}{4} \times [\text{ActivationReuse} + \text{WeightReuse}] \quad \text{Now, } DI = \frac{1}{4} \times [\alpha \times \text{ActivationReuse} + (1-\alpha) \times \text{WeightReuse}]$$

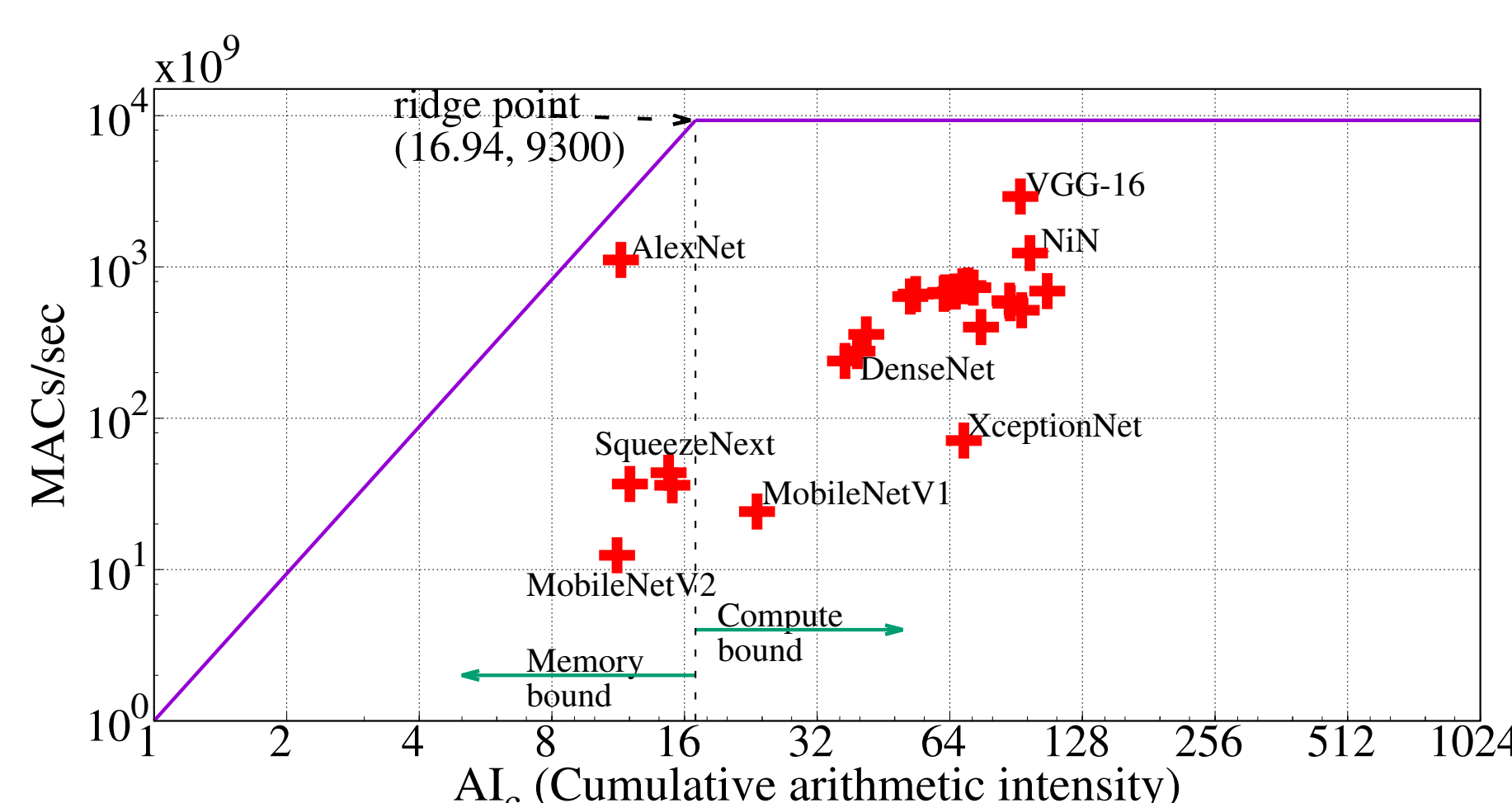


$$DI = \frac{1}{4} \times [0.8 \times \text{ActivationReuse} + 0.2 \times \text{WeightReuse}]$$

## EXPERIMENTAL RESULTS (1)

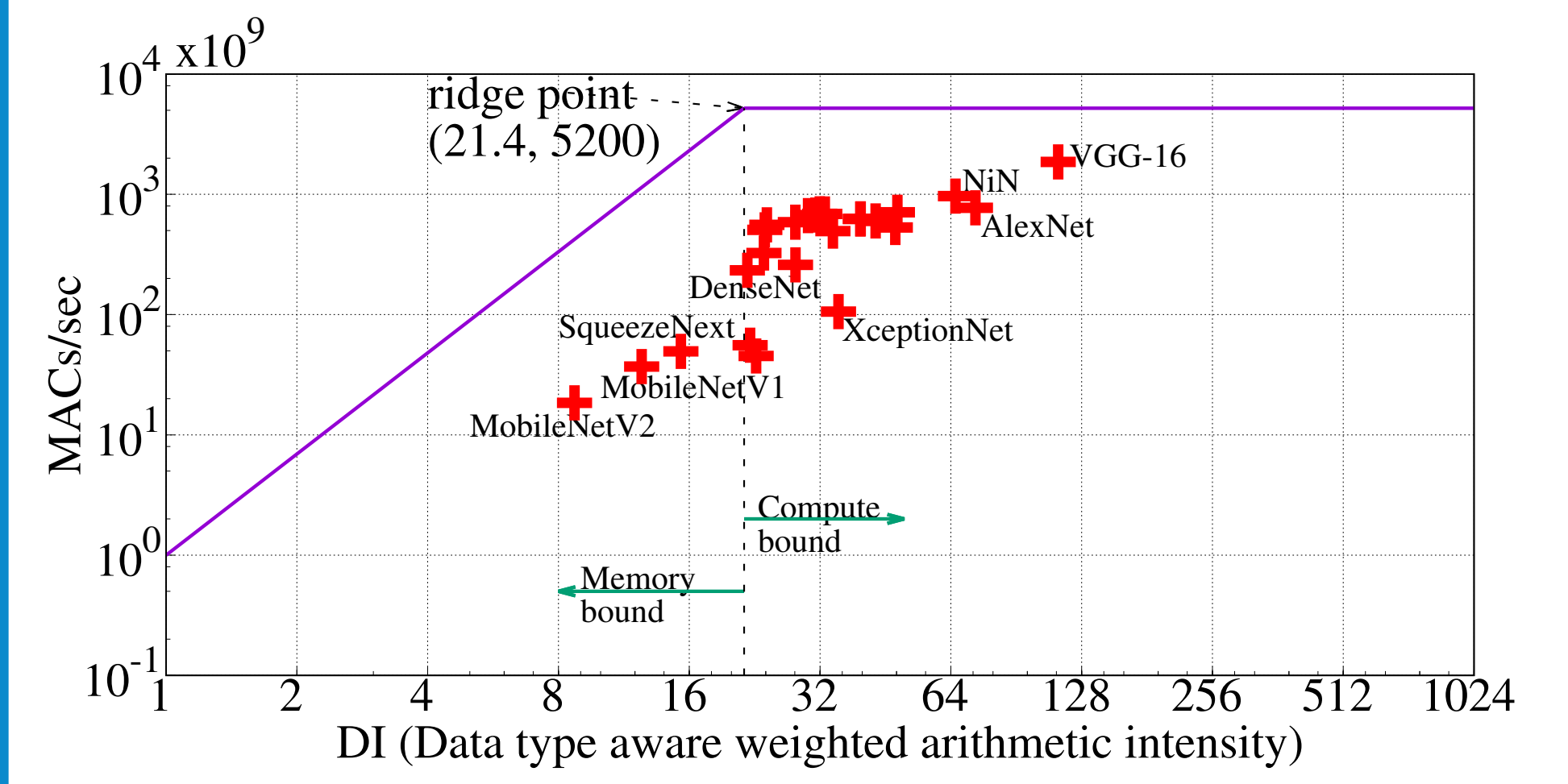


(a) Roofline model (with AI<sub>c</sub>) on P4000 GPU

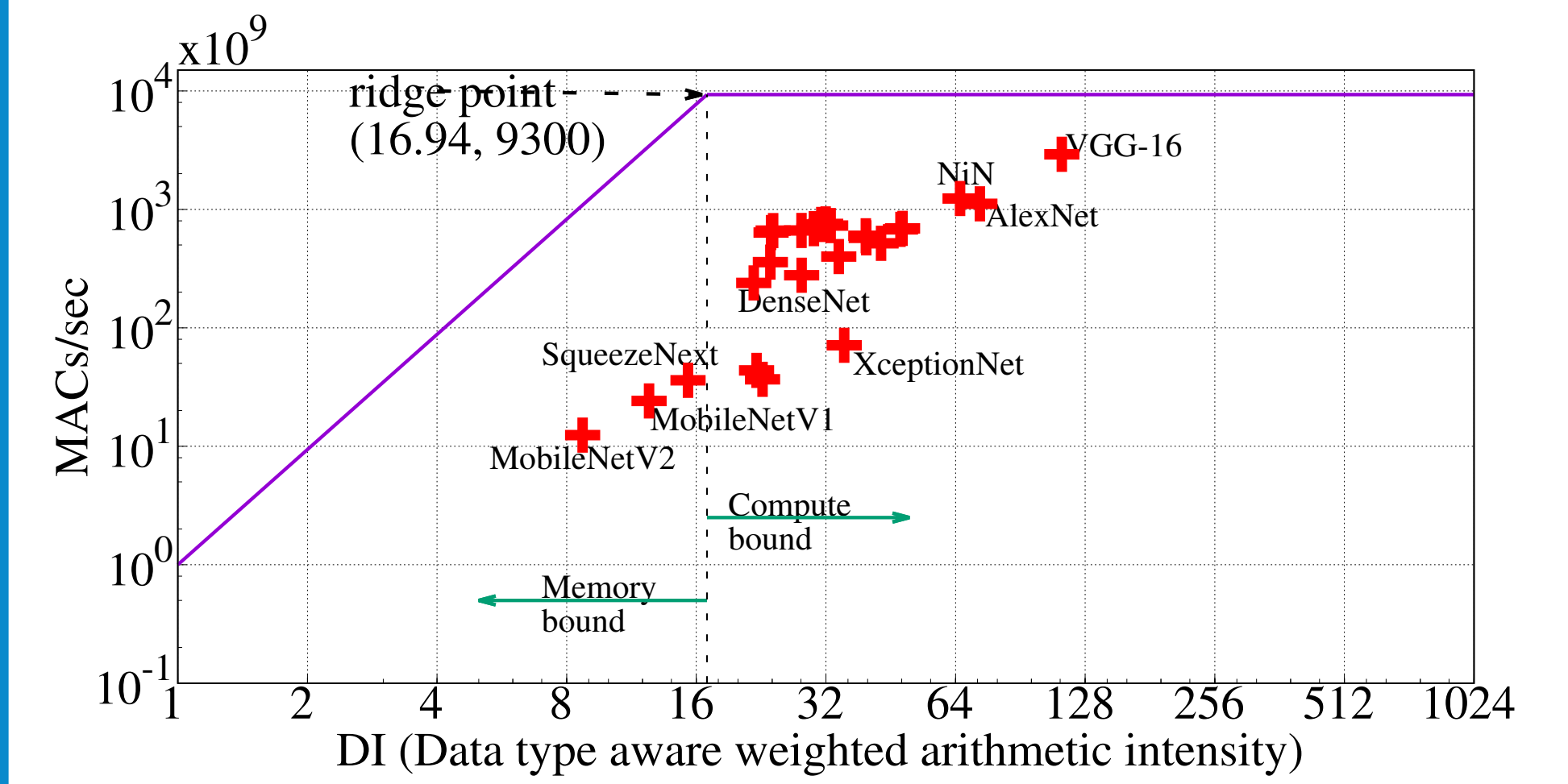


(b) Roofline model (with AI<sub>c</sub>) on P100 GPU

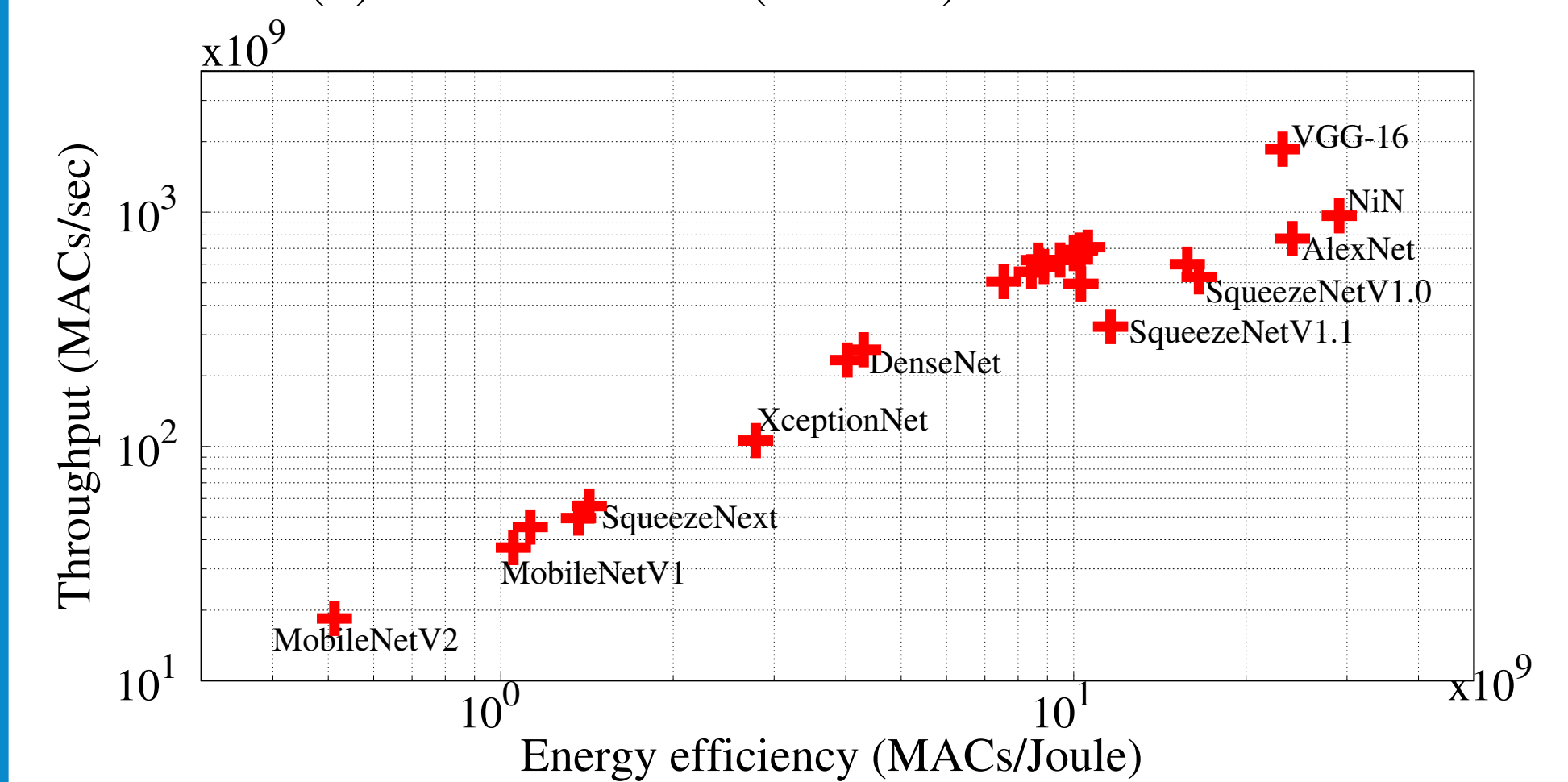
## EXPERIMENTAL RESULTS (2)



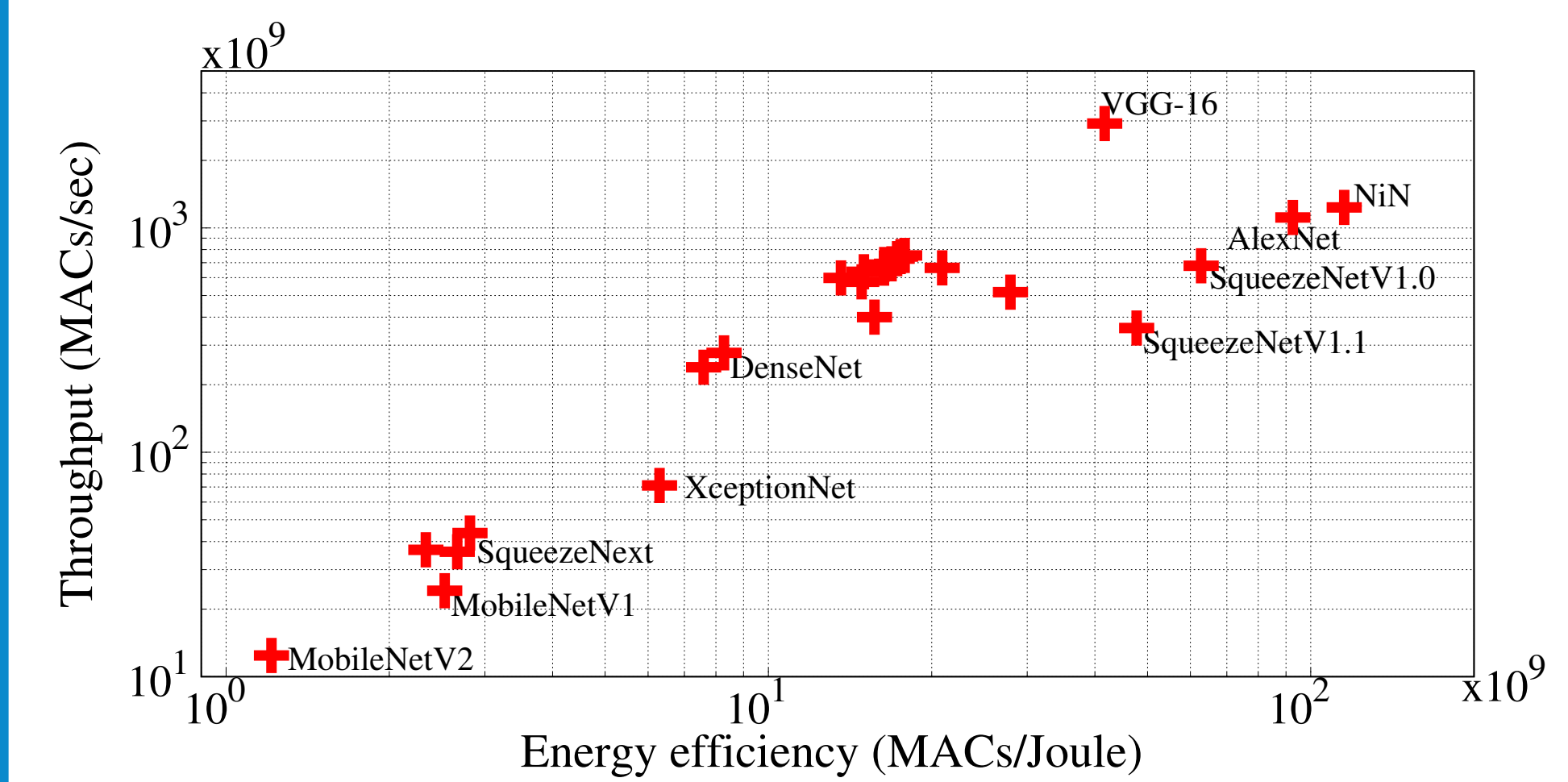
(c) Roofline model (with DI) on P4000 GPU



(d) Roofline model (with DI) on P100 GPU



(e) Throughput and energy efficiency on P4000 GPU



(f) Throughput and energy efficiency on P100 GPU

## DISCUSSION

GPU	X	Y	r = correlation (X,Y)
P4000	AI <sub>c</sub>	Energy efficiency	0.52
	DI (Ours)	Energy efficiency	<b>0.83</b>
P100	AI <sub>c</sub>	Energy efficiency	0.23
	DI (Ours)	Energy efficiency	<b>0.64</b>

Relative disparity ( $d_f$ ) between AI<sub>c</sub> and DI

$$d_f = \left( \frac{AI_c - DI}{AI_c} \right) \times 100 = 75 - 6.25 \times \left[ \frac{A}{W} + 3 \times \frac{W}{A} \right]$$

	Case 1: $A \ll W$	Case 2: $A \approx W$	Case 3: $A \gg W$
AI <sub>c</sub>	$\approx M_c/W$	$\approx 0.5 \times M_c/A$	$\approx M_c/A$
DI	$\approx 0.2 \times M_c/A$	$\approx 0.25 \times M_c/A$	$\approx 0.06 \times M_c/W$
$d_f$	$\approx 75 - 18.75 \times \frac{W}{A}$	$\approx 50$	$\approx 75 - 6.25 \times \frac{A}{W}$

## CONCLUSION AND FUTURE WORK

- AI<sub>c</sub> is a representative of energy-efficiency for mainly DNNs with  $A \approx W$ .
- AI<sub>c</sub> fails to estimate the data reuse in DNNs when  $|d_f|$  is high.
- DI is a better representative of energy-efficiency in DNNs.
- Evaluation of DI on DNN accelerators and FPGA.

## CONTACT

cs17mtech11010 [at] iith [dot] ac [dot] in