
IWES 2018

Energy-Efficient Deep Learning

Valentino Peluso



**POLITECNICO
DI TORINO**

Department of Computer Engineering

Deep Learning at the edge of the IoT

Data-Analytics in-the-Cloud

High data volume

High communication latency

Huge capacity storage

In-situ Data-Analytics

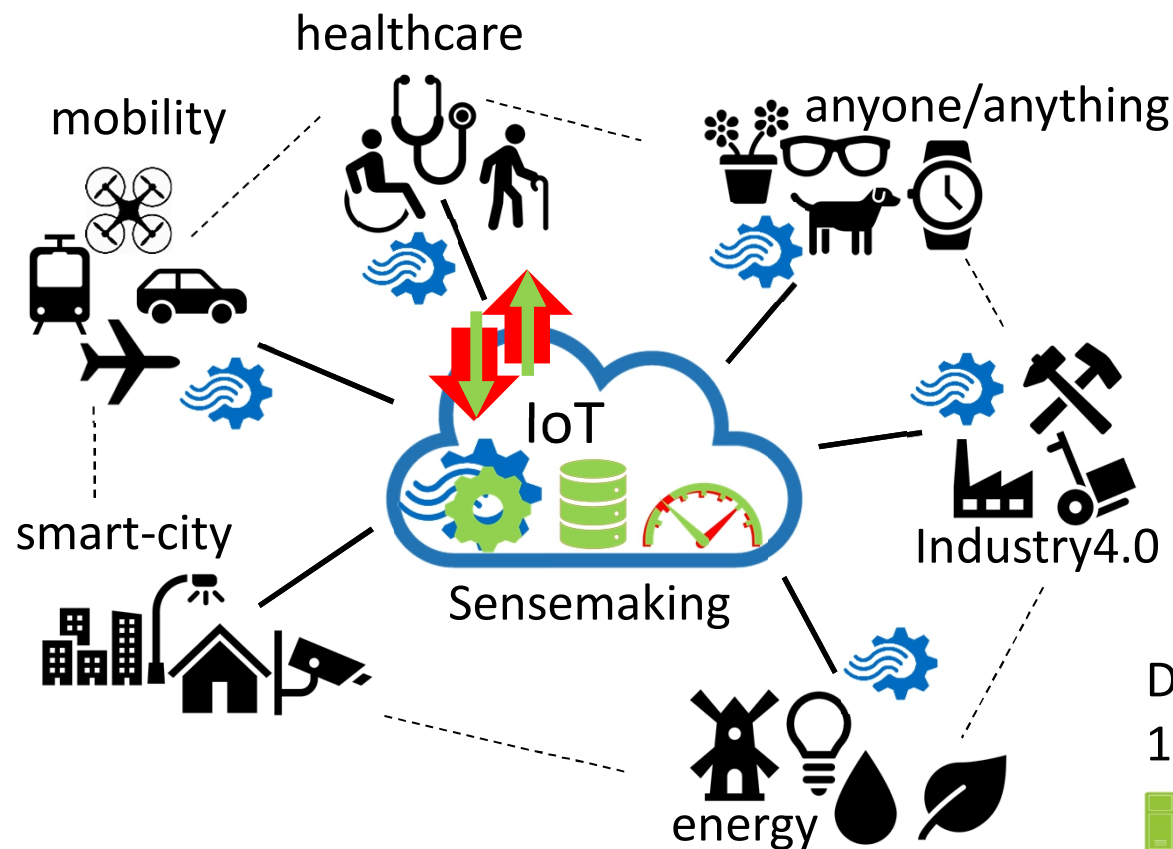
High information throughput

Low data volume

Low storage and network usage

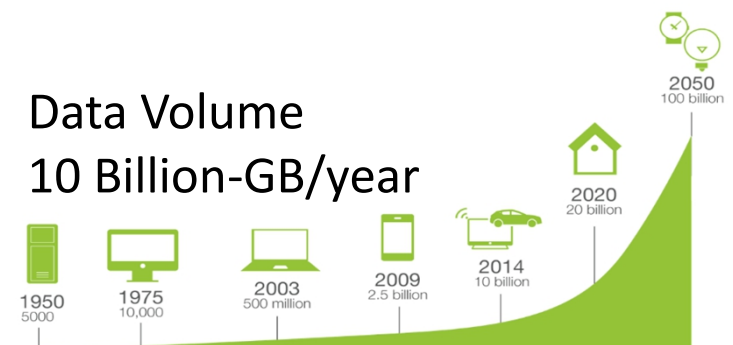
Real-Time M2M

Privacy



Scalability?

Data Volume
10 Billion-GB/year



Bring CNNs at the edge of the IoT

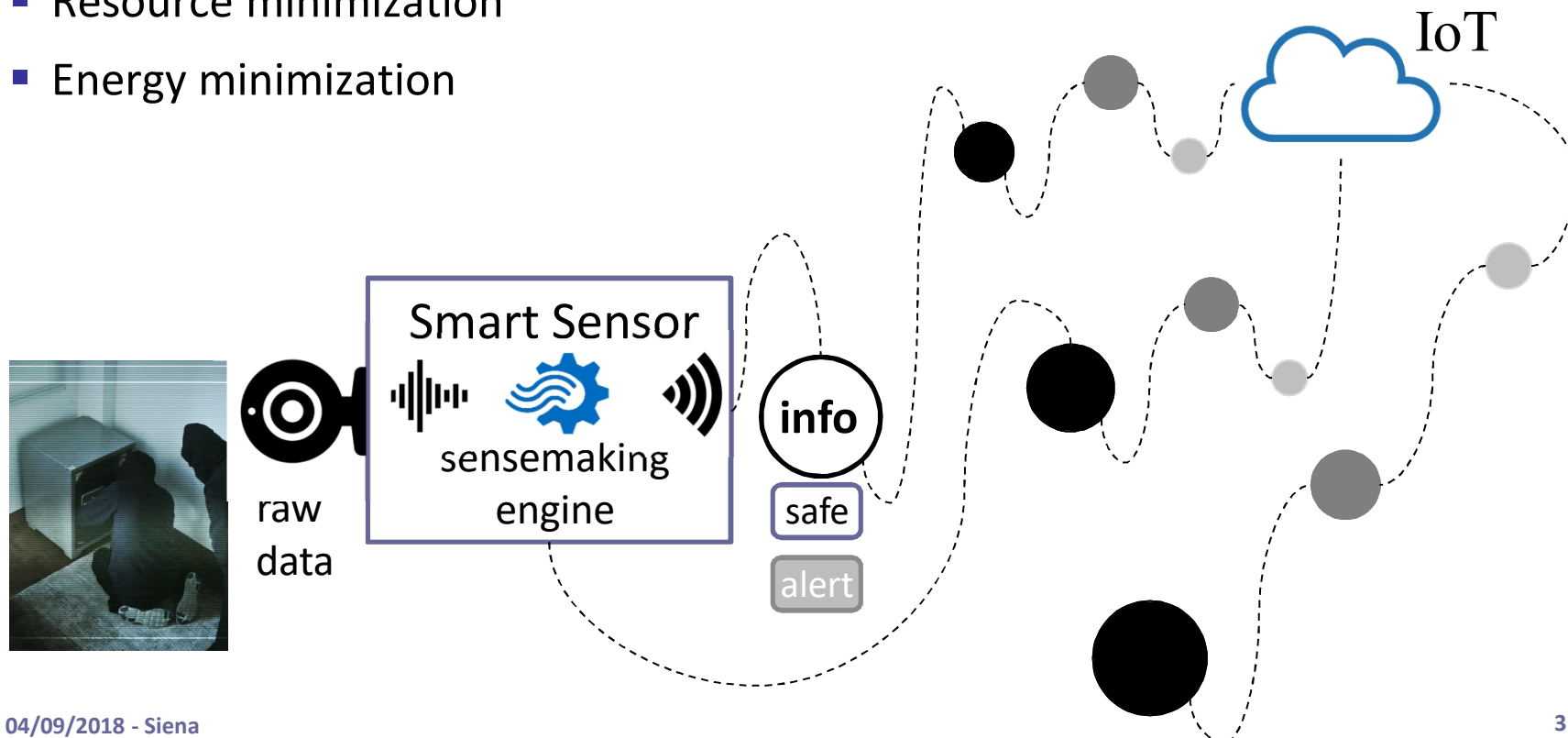
- **Data-flow computing on embedded platforms**

- Ultra-low power architecture
- Reconfigurable architecture (host different CNNs)



- **Optimal software-to-silicon mapping (design automation)**

- Resource minimization
- Energy minimization



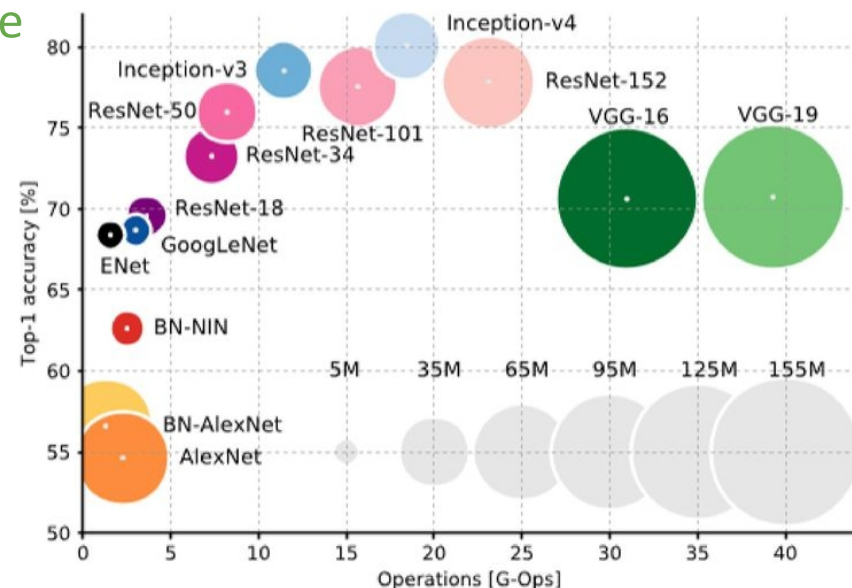
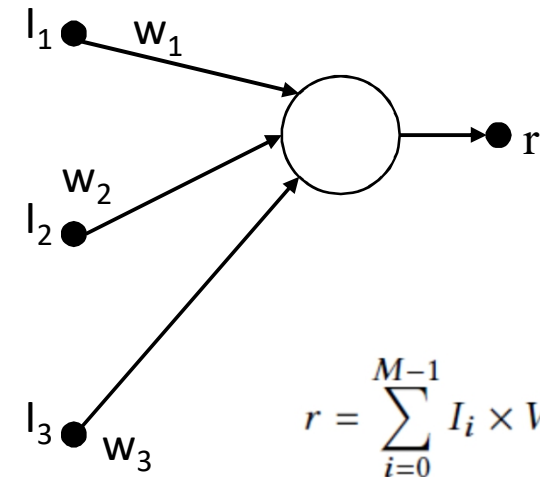
Fixed-Point CNNs

- Fixed-point (INT) arithmetic instead of floating-point (FP32) is a well established strategy for inference

- FP32: $\pm S \cdot 1.M \times 2^{\text{EXP}}$ [S=1b, M=23b, Exp=8b]
- INT16: I.Q [8b.8b]

- Fixed-point arithmetic

- require less memory
 - fit in resource-constrained HW
- minimize memory bandwidth usage
 - leverage SIMD data-paths
- achieve lower accuracy
 - reduced range
 - reduced precision



How to choose optimal precision

- **State-of-the-art: accuracy driven**

→ Need to consider also energy

1. Input-dependant

- Dynamic bit-width reconfiguration

2. Network-dependant

- Different DNNs show different tolerance

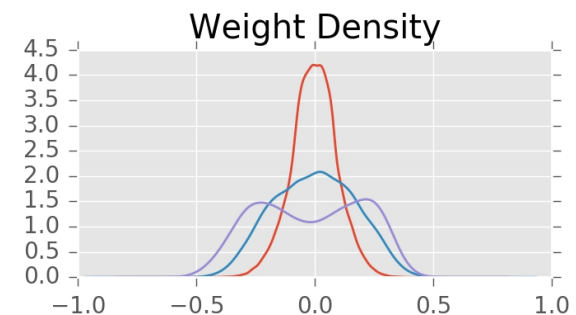
- Topology (Size, Depth, #Kernels)
- Trainable parameters (sparsity)

- Per-net quantization

- Same precision for each and every layer

- Per-layer quantization

- Exploit the dynamic range of activations and weights across different layers



Energy-Efficient Precision Scaling for ConvNets

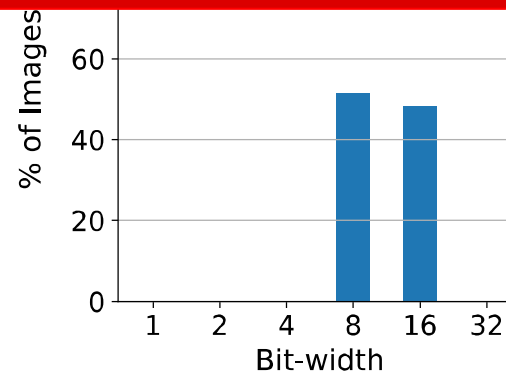
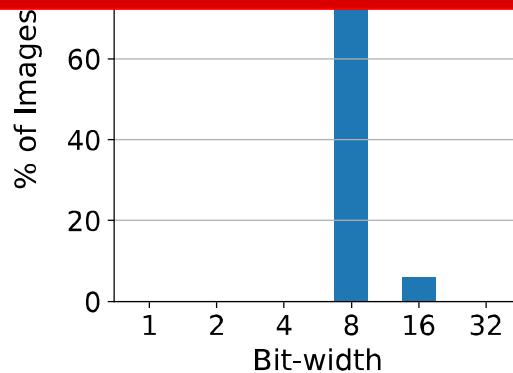
- **Broad Objective:**
 - Enable Dynamic Energy-Accuracy Scaling
 - Run-time adaptation according
 - application requirements and/or the context
 - available energy budget
- **Knobs:**
 - Arithmetic precision
- **Key feature:**
 - No Retraining
 - Costly
 - Training data may be not always available
 - Make use of a single weight-model rather than multiple models

- **Dynamic Bit-width Reconfiguration**
- **Energy-Aware Precision Scaling for ConvNets**

Motivation

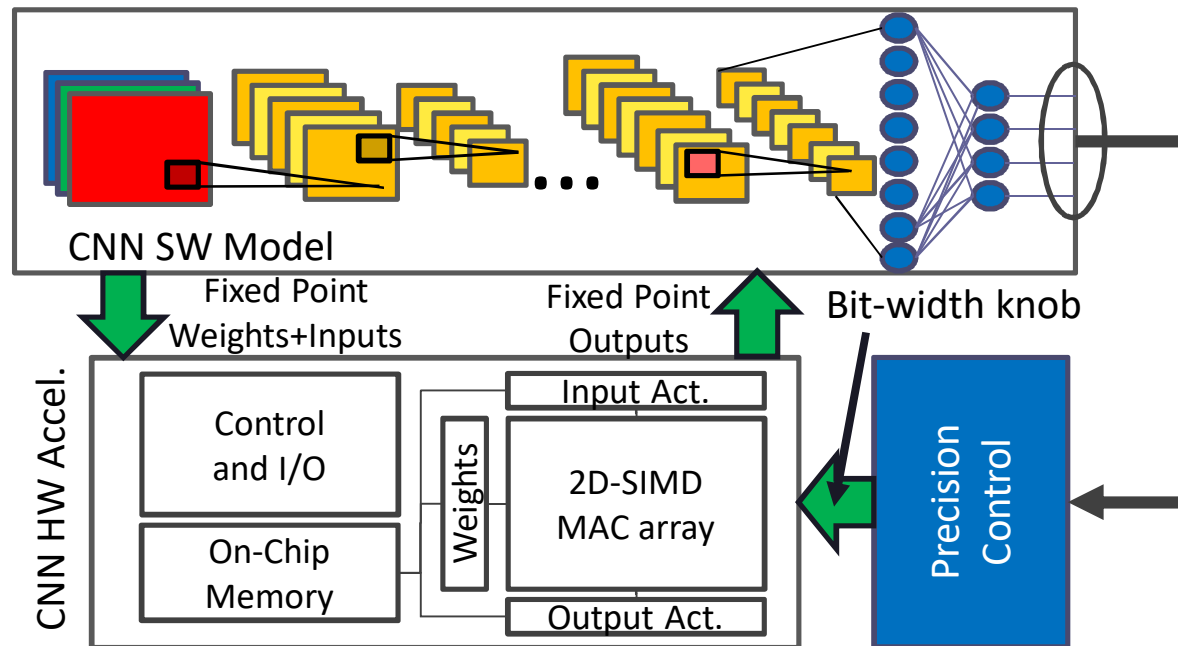
- Idea: not all inputs are **equally hard to classify**.
- Using the **same fixed-point bit-width** for all inputs may be sub-optimal.
- **Example:**

Consequence: static bit-width solutions require a costly retraining to (partially or totally) restore accuracy!



Proposed Approach

- Automatically **adapt the bit-width to the current input**.
- **Reduce energy consumption** compared to a *conservative* static bit-width (e.g. 16-bit in the previous two examples).
- **Increase accuracy** compared to lower precision (e.g. 8-bit)
- **General Framework:**

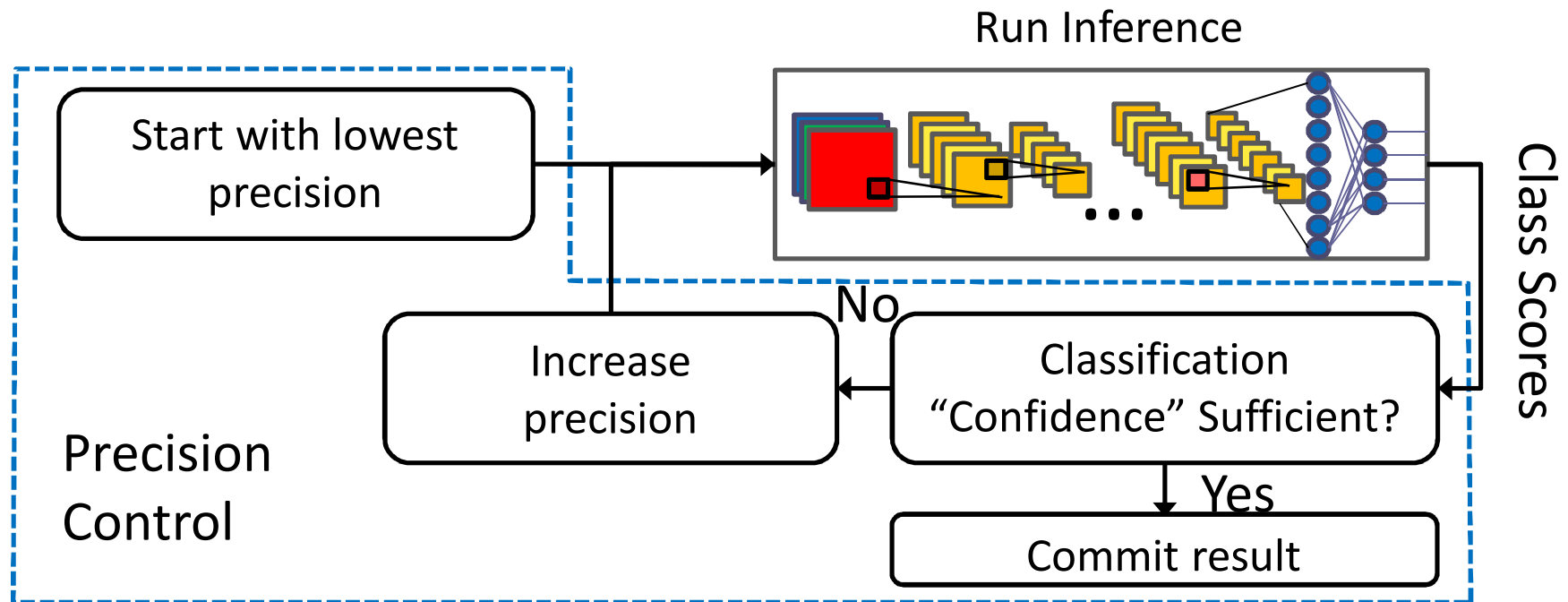


Implementation Details

- **Data format: Dynamic Fixed Point**
 - **Integer bits:** fixed, determined from variation ranges of activations and weights in each layer
 - **Fractional bits:** change depending on overall bit-width (may be negative)

- **Advantage: No storage duplication.**
 - Weights store at maximum precision, precision is reduced by LSB-truncation.

Runtime Bit-width Tuning



Classification “Confidence” Estimation

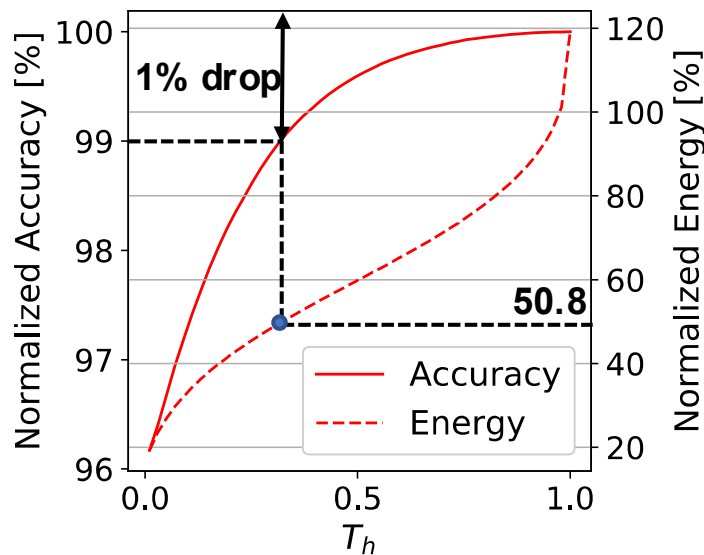
- Uses the **Score Margin (SM)**

$$SM = P(y_i|x) - P(y_j|x)$$

- $P(y_i|x)$ = output probability of the selected class
- $P(y_j|x)$ = second largest probability.
- **Idea:**
 - Large difference → only one class is probable
 - Small difference → “uncertain” between at least two classes.
- **Decision Strategy:**
 - Compare SM with a **Threshold (Th)**.
 - Changing the threshold allows exploring the **energy** (i.e. number of inferences per image) **versus accuracy trade-off**.

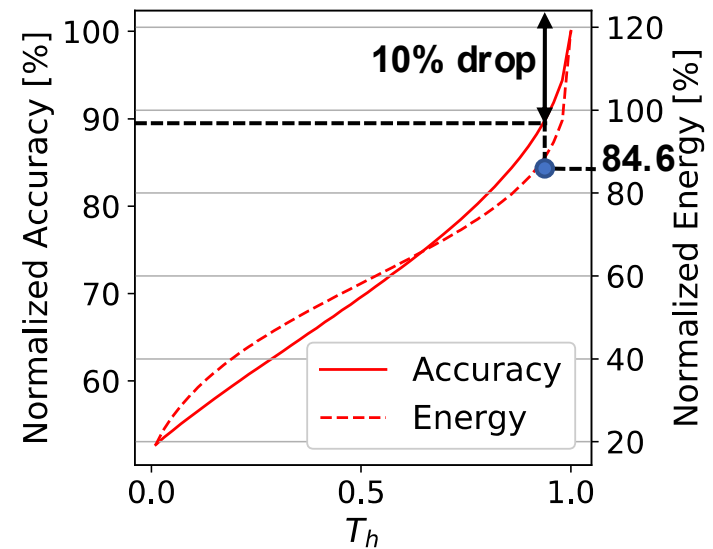
Results

- **Target CNNs:** CaffeNet (i.e. Caffe's AlexNet, Krizhevsky et al., NIPS2012) and CNN-M (Chatfield et al, arXiv2014)
 - Same task: ImageNet classification.
- **Target Accelerator:** Envisions (Moons et al., ISSCC2017)
- Different networks yield very different tradeoffs:



CaffeNet

**Baseline: Static 16-bit
Top-1 Accuracy**



CNN-M

- **Dynamic Bit-width Reconfiguration**
- **Energy-Aware Precision Scaling for ConvNets**

Per-net assignment

- **Four options:**
 - 16b activations, 16b weights (16x16)
 - 8b activations, 16b weights (8x16)
 - 16b activations, 8b weights (16x8)
 - 8b activations, 8b weights (8x8)
- **Unconstrained assignment**

	CNN-1	CNN-2	CNN-3
Network Architecture	input 32×32×3 conv 3×3×32 conv 3×3×64 maxpool dense 128 dense 10 softmax	input 32×32×3 conv 3×3×32 conv 3×3×32 maxpool conv 3×3×64 conv 3×3×64 maxpool dense 512 dense 10 softmax	input 32×32×3 conv 3×3×32 conv 3×3×32 maxpool conv 3×3×64 conv 3×3×64 maxpool conv 3×3×128 conv 3×3×128 maxpool dense 1024 dense 512 dense 10 softmax
	CIFAR-10	CIFAR-100	
	20 828 426	25 718 538	

**Need for a finer precision assignment
→ cross-layer optimization**

	32-bit FP	16×16 Fix	8×16 Fix	16×8 Fix	8×8 Fix
CNN-1	90.02%	90.01%	89.66%	76.64%	74.54%
CNN-2	81.77%	81.75%	80.99%	77.80%	72.17%
CNN-3	55.80%	55.82%	52.79%	22.27%	14.69%

Optimization Flow

- FP-32 pretrained model
- Calibration Set (5000k samples)

- HW energy model
- Accuracy Constraint
- Calibration Set (5000k samples)

- Test Set

Fixed-Point Conversion (INT16)

Per-Layer Precision Assignment

Model Evaluation

- Per-layer radix point scaling based on range analysis
- GPU-accelerated emulation tool fully integrated in TensorFlow
- Energy Aware
- Based on Simulated Annealing

Results

■ Accuracy Results

Max. Drop λ_{max}	1%	5%	10%
CNN-1	89.69% (0.35%)	86.82% (3.55%)	86.82% (3.55%)
CNN-2	80.76% (1.21%)	77.77% (4.87%)	75.57% (7.56%)
CNN-3	55.20% (1.11%)	54.01% (3.24%)	51.54% (7.67%)

■ Per-layer Precision Scaling

CNN-1			CNN-2			CNN-3		
1%	5%	10%	1%	5%	10%	1%	5%	10%
8×16	8×8	8×8	8×16	8×8	8×8	16×16	8×16	8×16
8×16	8×8	8×8	8×16	8×16	8×8	16×16	8×16	8×16
8×16	8×16	8×16	8×8	8×8	8×8	16×16	8×16	8×16
8×8	8×8	8×8	8×8	8×8	8×8	8×16	8×16	8×16
8	8	8	8×16	8×8	8×8	8×16	8×16	8×16
			16×8	8×16	8×8	16×16	8×16	8×8
			16	16	16	8×16	8×16	8×8
						8×8	8×8	8×8
						16×8	8×8	8×8
						16	16	16

Comparison and Energy Savings

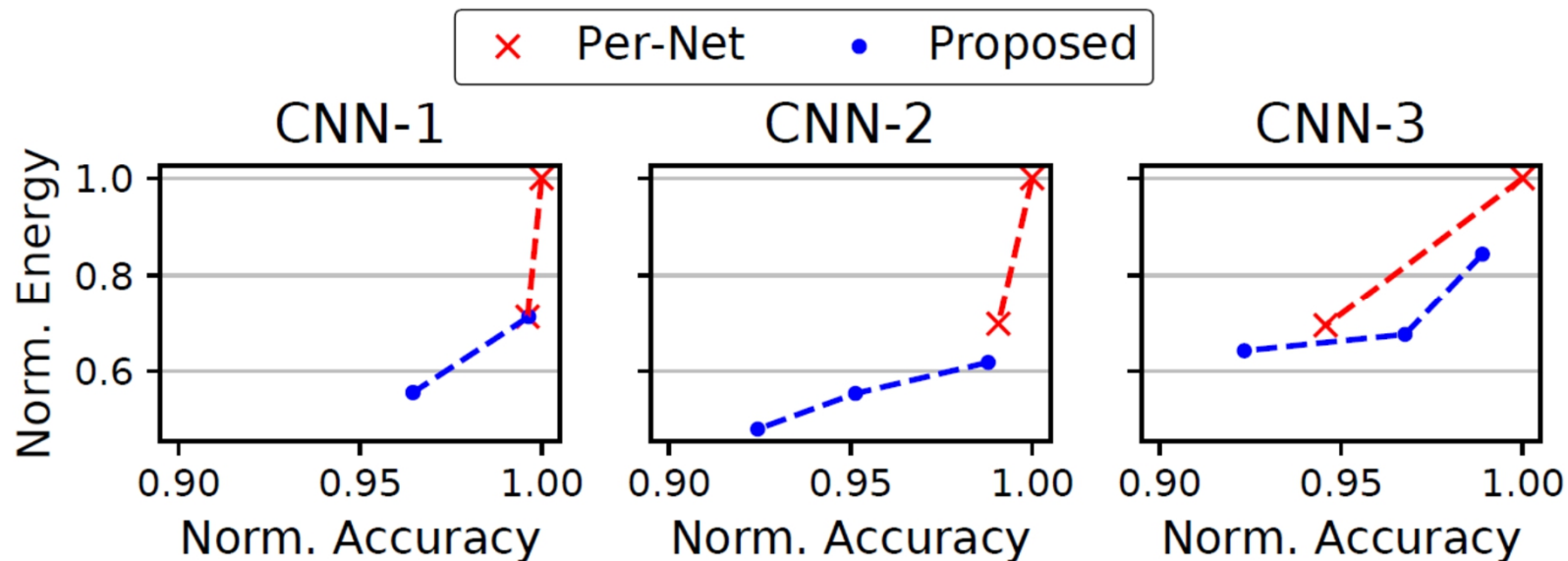
Per-Net

	32-bit FP	16×16 Fix	8×16 Fix	16×8 Fix	8×8 Fix
CNN-1	90.02%	90.01%	89.66%	76.64%	74.54%
CNN-2	81.77%	81.75%	80.99%	77.80%	72.17%
CNN-3	77.7%	77.7%	77.7%	77.7%	14.69%

10%
86.82% (3.55%)
75.57% (7.56%)
51.54% (7.67%)

Extended Pareto-points

→ dynamic adaptation to accuracy constraint or energy budget



- **Dynamic Energy Accuracy Scaling**
 - Allows to weigh the importance of accuracy vs energy depending on application/context
 - Especially interesting when retraining is not an option

- **Future works**
 - Combination of two techniques

- **Question?**

Dyn. Bit-width Reconfig. Results (2)

- **Comparison with other input-dependent strategies (on CaffeNet):**

Method	Energy Saving	Top-1 Drop	Multiple CNNs	Multiple Trainings
<i>Park et al, CODES2015</i>	53.7%	0.9%	Yes	Yes
<i>Tann et al, CODES2016</i>	32.61%	0.29%	No	Yes
This work	49.2%	0.89%	No	No

- **Warning:** different underlying hardware!

HW energy model

- **Computation Energy**

$$E^{comp} = \sum_{i=1}^L E^{MAC} N_{cycles}(x_i) \cdot N_i^{MAC} + E^{zero} \cdot N_i^{zero}$$

8x8 MAC depends on energy precision zero-skipping logic

- **Memory Energy**

$$E^{mem} = \sum_{i=1}^{2 \cdot L + 1} E^{MAC} \cdot [\alpha_i(x_i) + \beta_i(x_i) + \gamma_i(x_i)]$$

input activations weights output activations

Simulated Annealing

Algorithm 2: Simulated Annealing

Input: T_0 , T_f , X_0 , *cooling*, K_b , *iter*, λ_{max} , *valset*

Output: X

```
1  $T = T_0$ 
2  $E = \text{energy}(X_0)$ 
3  $E_{\max} = \text{energy}(\text{ones}(2L + 1)); E_{\min} = \text{energy}(\text{zeros}(2L + 1))$ 
4 while ( $T \geq T_f$ ) do
5     for  $i = 0; i < \text{iter}; i = i+1$  do
6          $\text{next\_state} = \text{move}(\text{current\_state})$ 
7         if  $\text{accuracy\_drop}(\text{next\_state}, \text{valset}, \text{tested}) < \lambda_{max}$  then
8              $E_{\text{next}} = \text{energy}(\text{next\_state})$ 
9              $\Delta E = (E_{\text{next}} - E_{\text{current}}) / (E_{\max} - E_{\min})$ 
10            if ( $dE < 0$ ) or ( $\exp[-\Delta E / K_b \cdot T] > \text{random}(0, 1)$ ) then
11                 $\text{current\_state} = \text{next\_state}$ 
12                 $E_{\text{current}} = E_{\text{new}}$ 
13                if  $E_{\text{current}} < E_{\text{best}}$  then
14                     $E_{\text{best}} = E_{\text{current}}$ 
15                     $\text{best\_state} = \text{current\_state}$ 
16                end
17             $\text{update}(\text{tested})$ 
18        end
19         $T = T \cdot \text{cooling}$ 
20 end
21 return  $\text{best\_state}$ 
```
