# Enabling Energy-Efficient, High-Performance **DNN Inference Using Approximate DRAM**

Skanda Koppula A. Giray Yağlıkçı Lois Orosa Roknoddin Azizi Taha Shahroodi **Onur Mutlu** Konstantinos Kanellopoulos

#### **1: Summary**

**<u>Problem</u>**: Challenges of DNN inference:

- High DRAM energy consumption and latency
- **Goal:** Reduce **DRAM oltage/timing scaled DRAM** for DNN inference to exploit error tolerant DNN workloads

**EDEN:** DNN Inference Using Approximate DRAM

- Techniques to maintain accuracy through **error tolerance** boosting, DNN characterization, DNN to DRAM mapping, and **DRAM error modeling** 

#### **Results**:

- Average **21% DRAM energy savings**, **8% speedup** on CPU
- Average 37% DRAM energy savings on GPU
- Average **31% DRAM energy savings** on DNN accelerators

EDEN is applicable to other parameters and technologies

## **2: DNN Basics and DRAM Parameters**



- Modern DNNs can have **100+ layers**
- **3 main data types** in a DNN layer: Weights, Input Feature Maps (IFMs), **Output Feature Maps (OFMs)**
- Large # weights/IFMs enable **high** learning capacity
- If weights/IFMs have **bit errors**, the DNN can still maintain accuracy



CPU, GPU, or

**DNN Accelerator** 





DRAM operates at a **standard** 

voltage (e.g., DDR3 at 1.35V)

**ETH** zürich

SAFARI

#### **4: DNN Error Tolerance Boosting**

#### **3: Overview of EDEN**



**Problem:** Retraining is **not always feasible** on the approx. DRAM device • **8 DNN workloads** with int4, int8, int16, FP32 quantizations ResNet101, MobileNetV2, VGG16, DenseNet201, SqueezeNet1.1, AlexNet, YOLO, YOLOTiny

**Example:** Boosting Error Tolerance of **ResNet101** 

DNN tolerance boosting

improves

a DNN's **bit error** 

tolerance by 5-10x



**Goal:** Perform retraining and error characterization without use of the approximate DRAM device

Model 0: Uniform Random Model 1: Bit Value Dependent Model 3: Bitline Correlated Model 4: Wordline Correlated



**Modification** to Retraining: Forward Pass with the DRAM **Error Model** 

- Custom **PyTorch**-based DNN framework to run DNN inference with error models
- **SoftMC** framework<sup>1</sup> to run inference data accesses on **real DDR3 DRAM modules**
- Ramulator, ZSim, GPGPUSim, and SCALE-Sim used for DRAM, CPU, GPU, Eyeriss, and TPU simulation
  - Full configuration can be found in the paper Ο
- Inference libraries from **DarkNet**, Intel OpenVINO, TVM

	FP32			int8		
Model	BER	$\Delta V_{DD}$	$\Delta t_{RCD}$	BER	$\Delta V_{DD}$	$\Delta t_{RCD}$
ResNet101	4.0%	-0.30V	-5.5ns	4.0%	-0.30V	-5.5ns
MobileNetV2	1.0%	-0.25V	-1.0ns	0.5%	-0.10V	-1.0ns
Sample of DNNs and their tolerable BERs						

#### **DRAM Energy Reductions on CPU Performance Improvements on CPU:**



**Average 21% DRAM energy reduction** maintaining accuracy within 1% of original

1.20
1.15 EDEN EDEN tRCD = 0
FP32 int8
YOLO-T YOLO ResNet VGG SqueezeNet DenseNet Gmean
YOLO-T YOLO ResNet VGG SqueezeNet DenseNet Gmean

Average 8% system speedup, with some achieving **up to 17% speedup** 

#### **DRAM Energy Reduction** on GPU and Accelerators

- Average **31% on Eyeriss**
- Average **32% on TPU**
- Average **37% on Titan X**



1.35

2.5

5.0

tRCD (ns)

1.05

1.20

Voltage (V)



**Example:** DNN Accuracy of **LeNeT** on **SoftMC** 

**Boosting with error models** helps maintain accuracy while reducing voltage and latency on real DRAM modules

## **Other Results in the Paper**

acy (%)

- Error resiliencies across different DNNs/ quantizations
- Validation of the boosting mechanism

7.5 10.0 12.5

- Support for error models using real DRAM modules
- Comparison of different DRAM error models
- Breakdown of energy savings on different workloads for **GPU** and **TPU**

#### skoppula.github.io/pdfs/eden.pdf



GPUs and accelerators are **effective at hiding DRAM latency** due to (1) *effective pre-fetching* and (2) large register banks and SRAM buffers (exploiting the fixed memory access patterns on DNN inference), so we find little performance improvement reducing tRCD on GPU and DNN accelerators