

ZEM: Zero-cycle Bit-masking Module for Deep Learning Refresh-less DRAM

Duy-Thanh Nguyen¹, Nhut-Minh Ho², Minh-Son Le¹,
Weng-Fai Wong², Ik-Joon Chang¹
¹{dtnguyen,sonlm,ichang}@khu.ac.kr, ²{minhnhn,wongf}@comp.nus.edu.sg



I. Abstract

In sub-20nm technologies, DRAM cells suffer from poor retention time. With the technology scaling, this problem tends to be worse, significantly increasing refresh power of DRAM. This is more problematic in memory heavy applications such as deep learning systems, where a large amount of DRAM is required, DRAM refresh power contributes to a considerable portion of total system power. With the growth in deep learning workloads, this is set to get worse. In this work, we present a *zero-cycle bit-masking* (ZEM) scheme that exploits the asymmetry of retention failures, to eliminate DRAM refresh in the inference of convolution neural networks, natural language processing, and the image generation based on generative adversarial network. Through careful analysis, we derive a bit-error-rate (BER) threshold that does not affect the accuracy of inference. Our proposed architecture, along with the techniques involved, are applicable to all types of DRAMs. Our results on 16Gb devices show that ZEM can improve the performance by up to 17.31% while reducing the total energy consumed by DRAM by up to 43.03%, dependent on the type of DRAM.

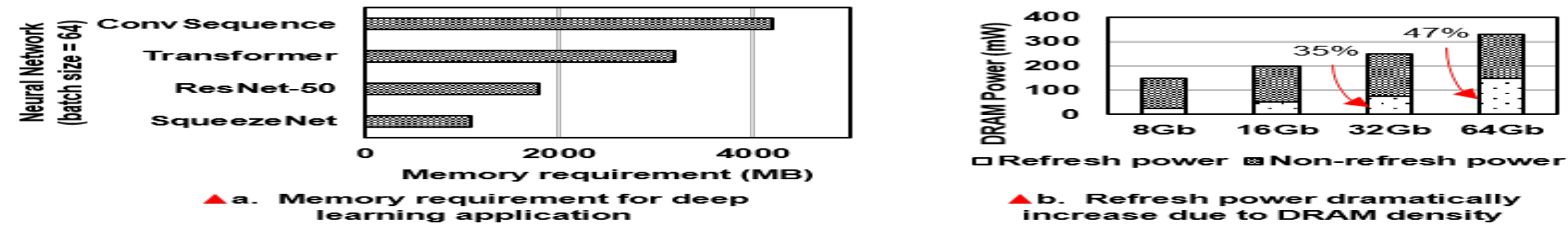
II. Introduction

Motivation

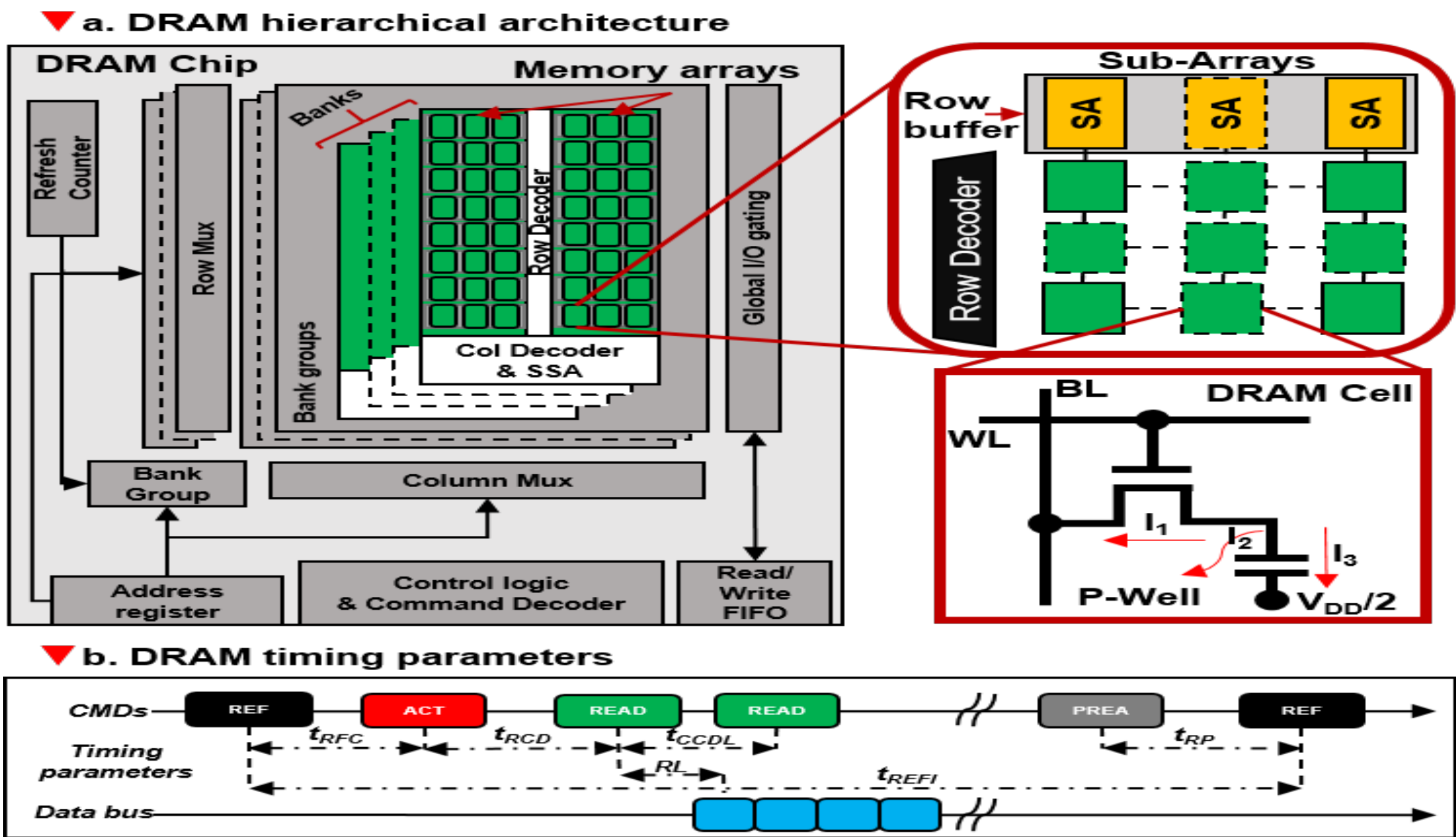
- DNNs become deeper and wider-> Large memory
- DNN can be tolerant to the errors
- DRAM refresh can consume up to 46% of total system power

→ DRAM refresh power can be reduced

DRAM Power And DNN memory requirements



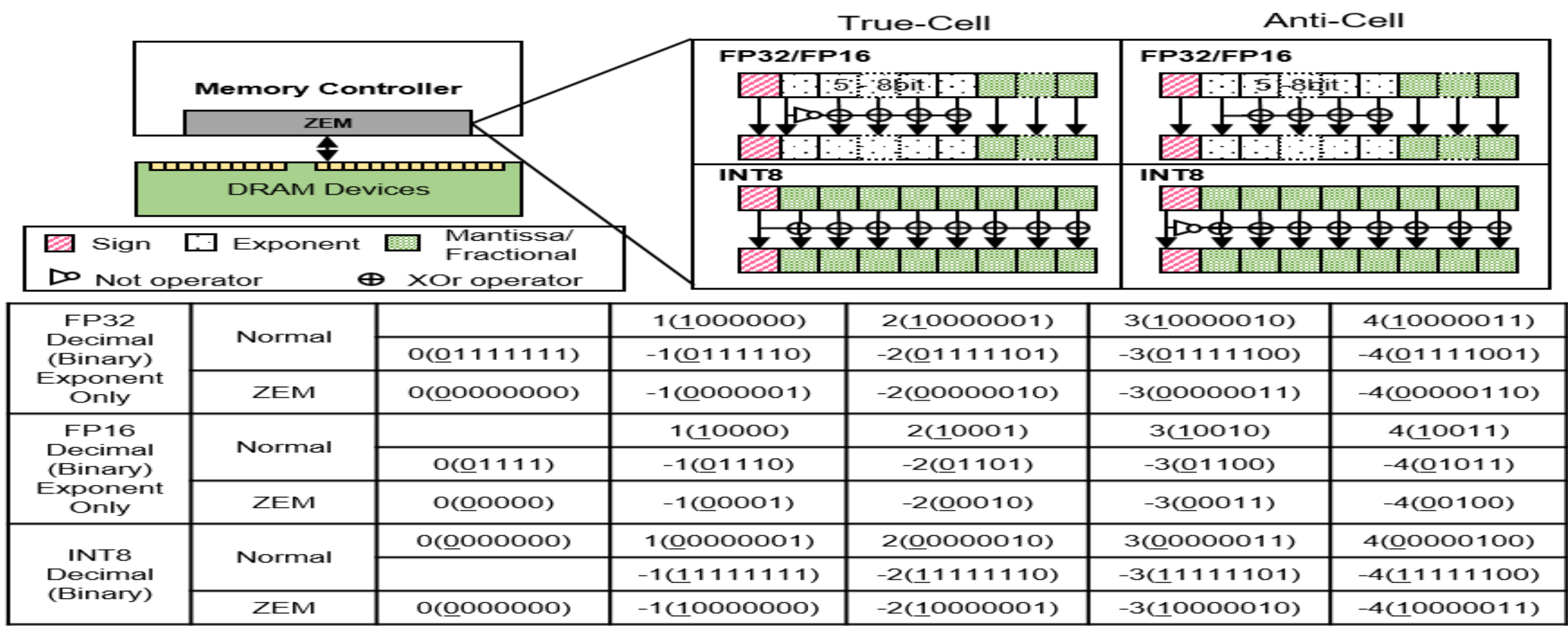
DRAM architecture



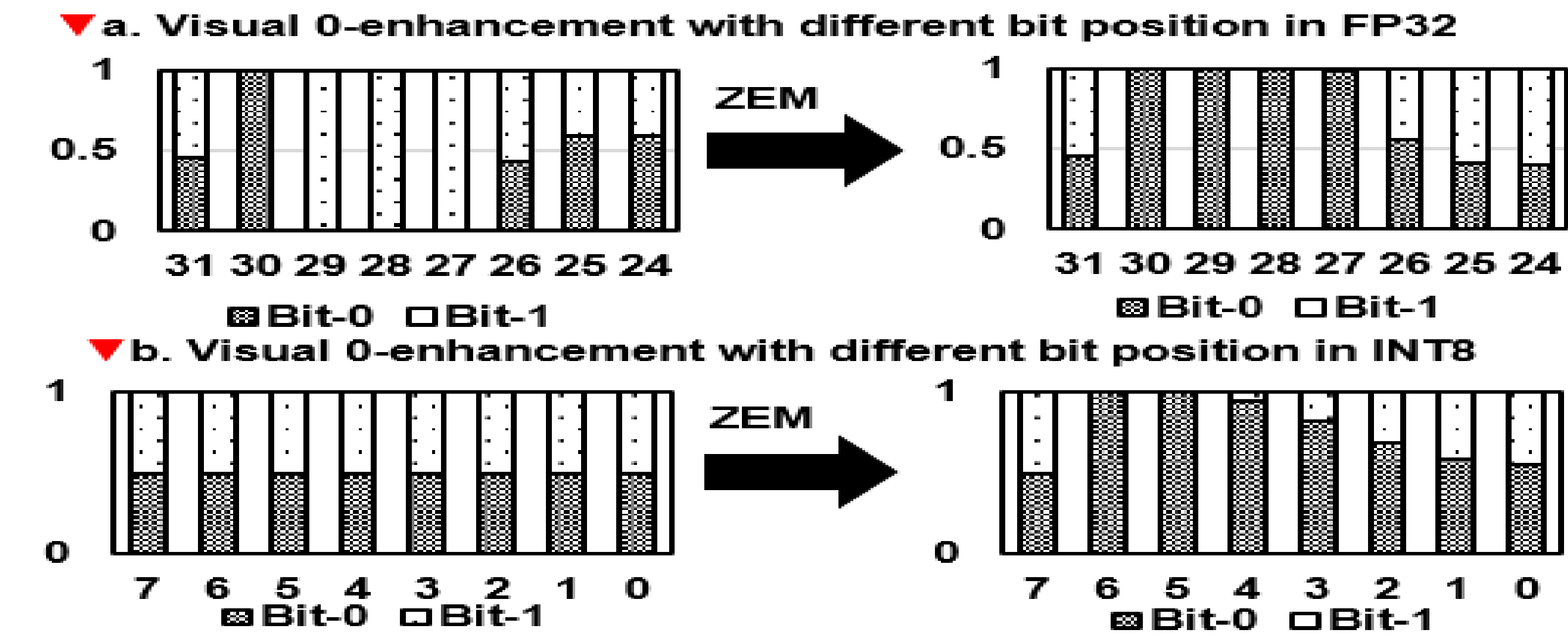
DRAM errors rarely occur on one bits

III. Our Approaches

Zero-cycle bit-Masking

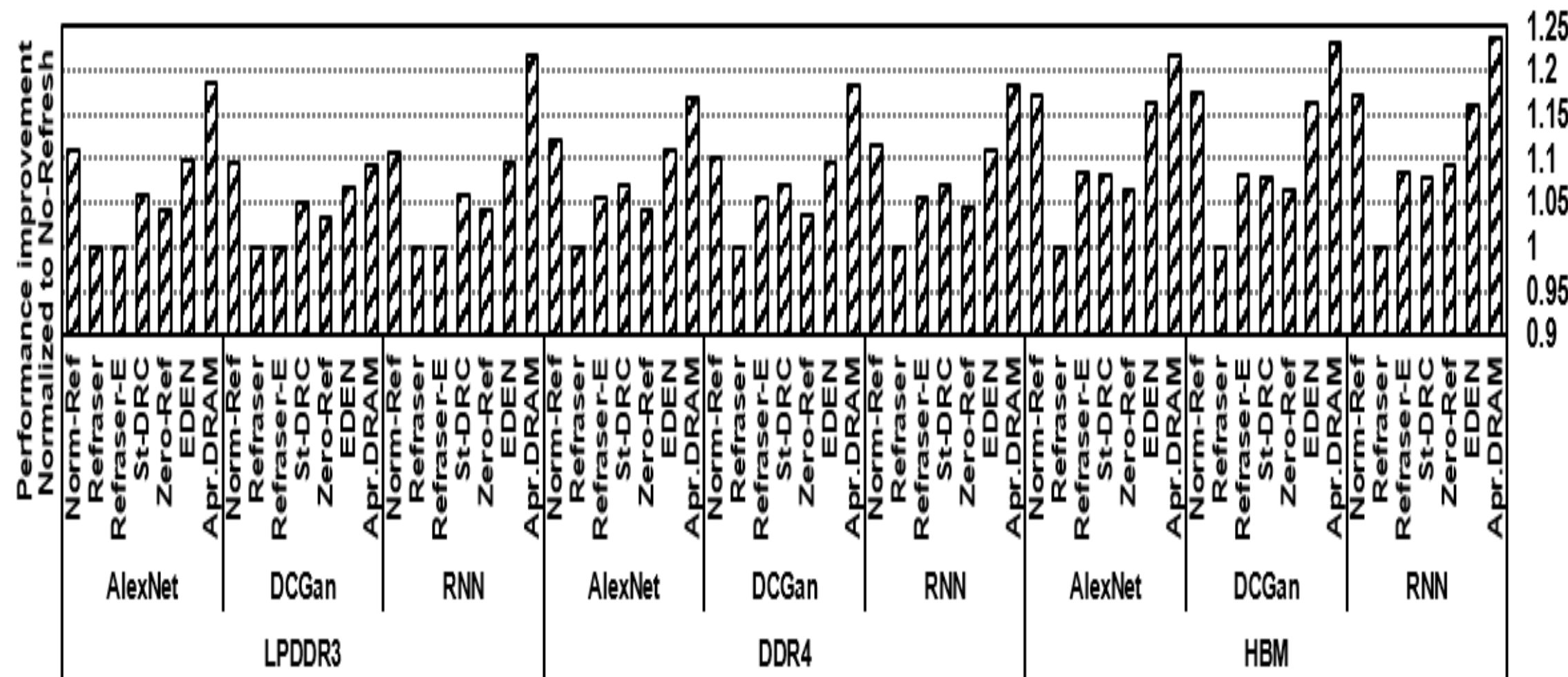


Asymmetry-enhancement with ZEM

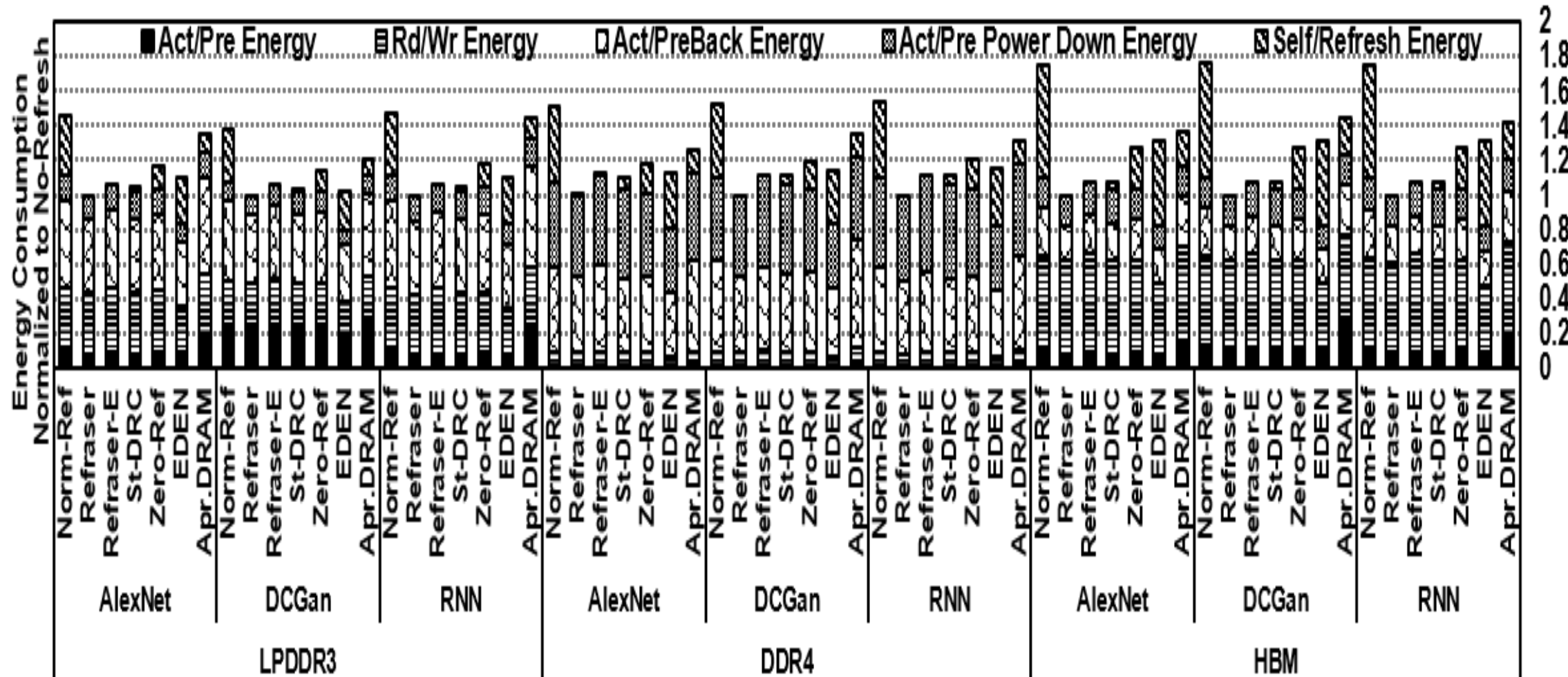


V. Energy/Performance Simulation Results

Performance Improvement

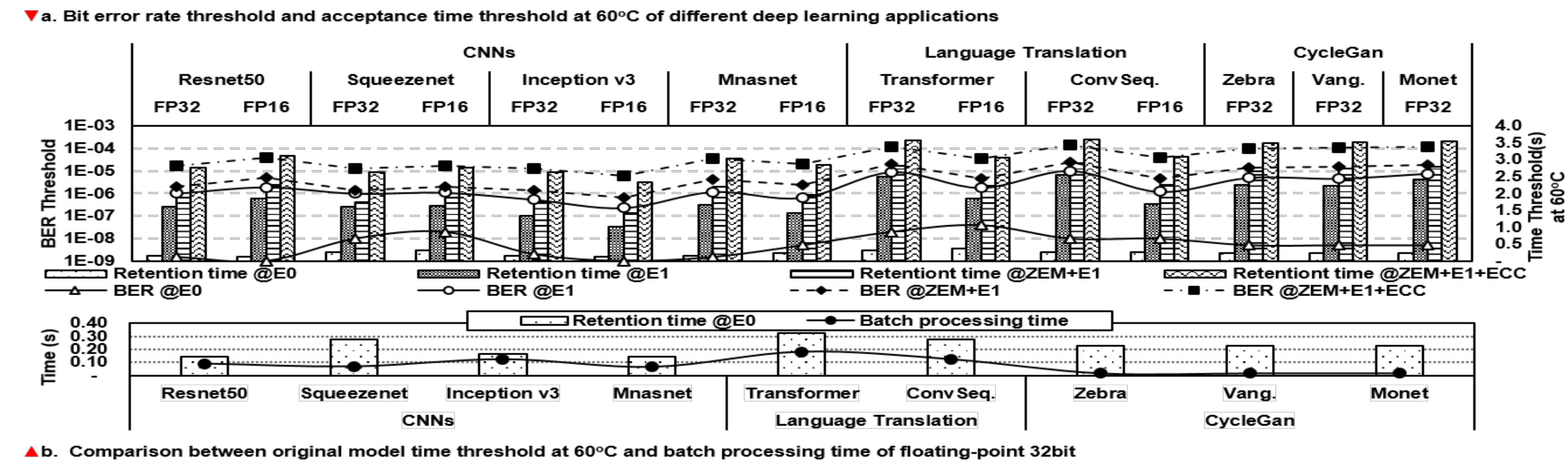


Energy Saving



IV. Validation

Floating-point 32/16bit



Integer 8bit

