# St-DRC: Stretchable DRAM Refresh Controller with No Parity-overhead Error Correction Scheme for Energy-efficient DNNs

**Duy-Thanh Nguyen[1]**, Nhut-Minh Ho[2], Ik-Joon Chang[1]

[1]Kyung Hee University, Republic of Korea

[2]National University of Singapore, Singapore

# Outline

- Motivation
  - The Effect of DRAM Refresh Relaxation
- Major Challenge
  - Floating-point IEEE754 under Retention Errors
  - Characteristic of DNN's Data
  - Large Sensitivity to Bit-errors of Some Exponent Bits
- Our Approach: Significant-Bit Protection
  - DRAM Controller with Stretchable Refresh Period
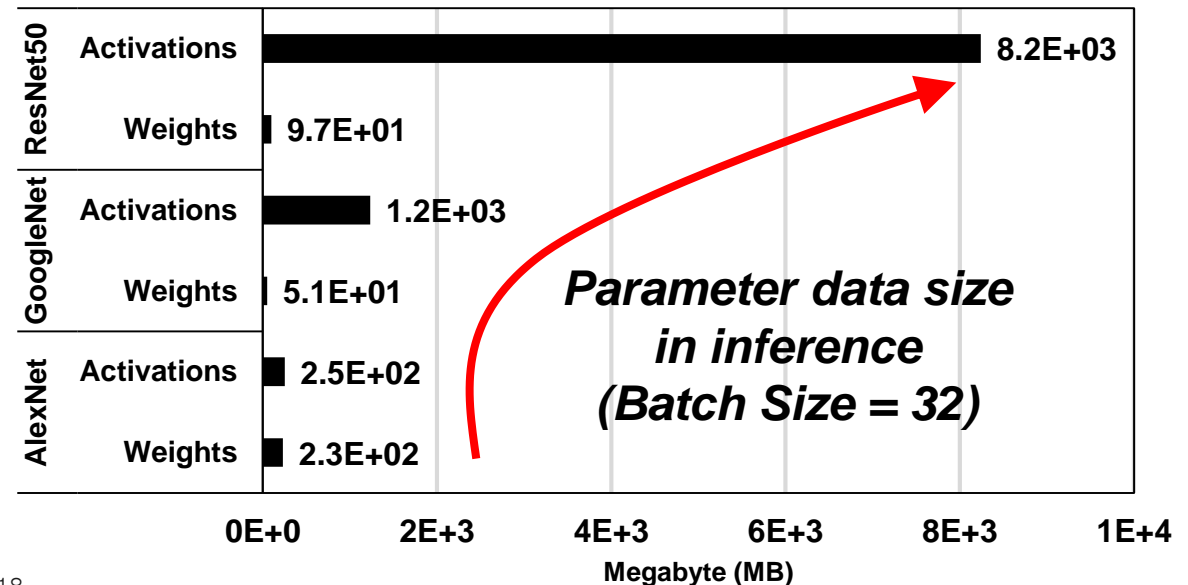  - Validation
- Energy and Performance Simulation Results
- Conclusion

- **Large processing time for the training of DNNs**
  - Training speed needs to be improved
- **Hybrid CPU-GPU platform is widely used for training DNNs**
  - Main Memory DRAM + GPU Memory DRAM → **DRAM power is very significant**

| Nvidia Server | Main Memory (DDR4) | GPU Memory (HBM2) |
|---|---|---|
| DGX-1 (8X Tesla V-100) | 512GB | 32GB per GPU  x  8 = 256GB |
| DGX-2 (16X Tesla V-100 /8X Tesla V-100) | 1.5TB | 32GB per GPU  x  16 = 512GB/ 32GB per GPU  x  8 = 256GB |

- **DNNs become deeper and wider**
  - The size of DNN Parameters tends to be larger
  - DRAM size should be larger
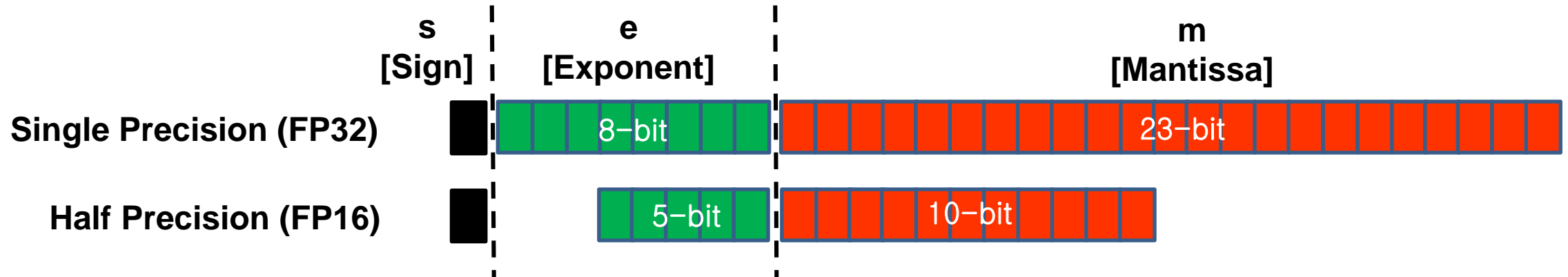  - **DRAM power would be more significant in data-center**



*Parameter data size in inference (Batch Size = 32)*

ResNet50 Activations: 8.2E+03
ResNet50 Weights: 9.7E+01
GoogleNet Activations: 1.2E+03
GoogleNet Weights: 5.1E+01
AlexNet Activations: 2.5E+02
AlexNet Weights: 2.3E+02

Megabyte (MB)

# The Effect of DRAM Refresh Relaxation

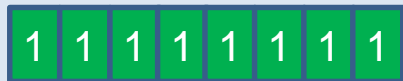| | RAIDR @ISCA12 | REFLEX @ISCA15 | Flikker @ASPLOS2011 | Quality-aware DRAM @DATE2015 |
|---|---|---|---|---|
| DRAM Power Saving Rate | 16.1% | 20% | 25~32% (standby), 20~25%(overall) | 73% (But some quality-loss) |
| System Performance Improvement Rate | 8.1% | 7~10% | Not discussed | Not discussed |
| Application | Generic | Generic | Video | Video |
| Precise/Approximate | Precise | Precise | Approximate | Approximate |
| Challenge | Retention time characterization of full rows for all chips (Large Verification Overhead) | Retention time characterization of full rows for all chips (Large Verification Overhead) | Not applicable to DNN (Large Accuracy Degradation) | 1. Retention time characterization of full rows for all chips (Large Verification Overhead) 2. Not applicable to DNN (Large Accuracy Degradation) |

➤ Our Contribution: Negligible accuracy degradation in DNN in spite of some retention errors, Reasonable verification effort, Significant power saving, System performance improvement

# Floating point IEEE 754 under Retention Errors

**s**
**[Sign]**
**e**
**[Exponent]**
**m**
**[Mantissa]**

**Single Precision (FP32)**
8–bit
23–bit

**Half Precision (FP16)**
5–bit
10–bit

**Conversion {s,e,m} $\rightarrow$ {$-1^s$ x M x $2^{(e-Bias)}$ | Bias = 127 for FP32 or 15 for FP16, M = 1.m}**

**Significant Challenge Due to Retention Errors**

1 1 1 1 1 1 1 1

When all 1's are in the exponent,
the data will become ±Infinity/NaN

**Operand1**
**±Infinity/NaN**

**ALU**

**Operand2**

**NaN**
**(Not-A-Number)**

**Error propagation due to +/- Infinity and NaN $\rightarrow$ Catastrophic system errors**

Weight parameter

Log₂ scale
- AlexNet
- GoogleNet
- ResNet50

$99\% \sim [2^{-13} : 2^{-2}]$

0 1 1

Force

(FP32)

3-bit can be freely used for another purpose

10~23−bit

Activation parameter

Log₂ Scale
- AlexNet
- GoogleNet
- ResNet50

$99\% \sim [2^{0} : 2^{11}]$

$99\% \sim [2^{-10} : 2^{11}]$

Weight and Activation are fully covered with only 5-bit exponent $\sim [2^{-15} : 2^{15}]$
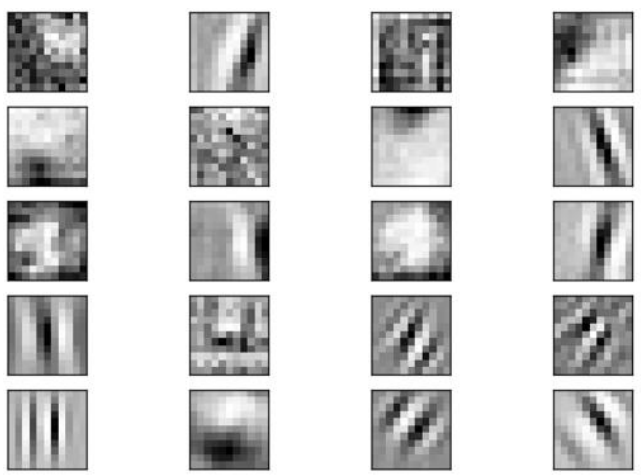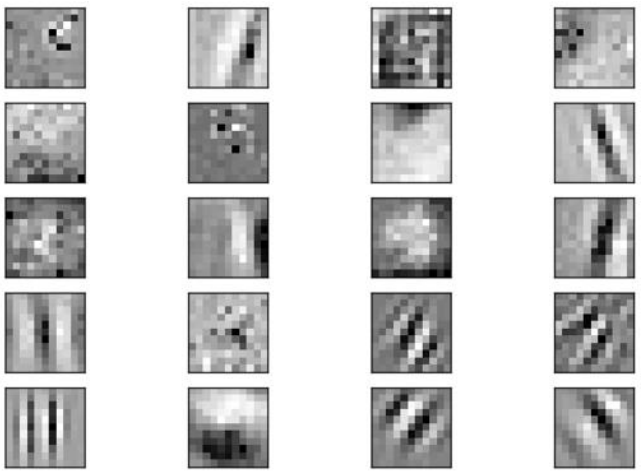
(FP32/FP16)   5−bit   10~23−bit

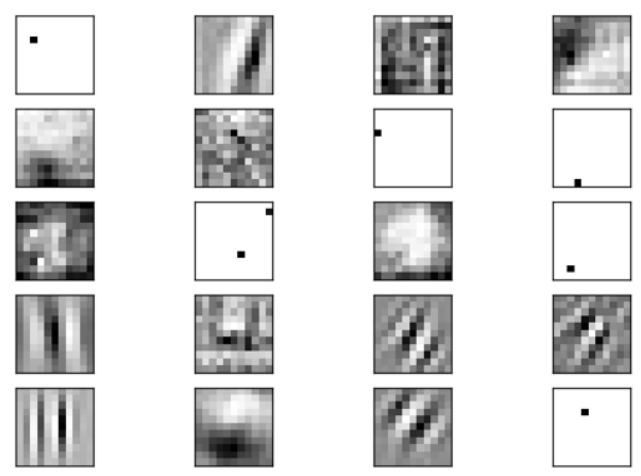# Large Sensitivity to Bit-errors of Some Exponent Bits (Inference)

**Original Kernel**

**50 times Error Accumulation**
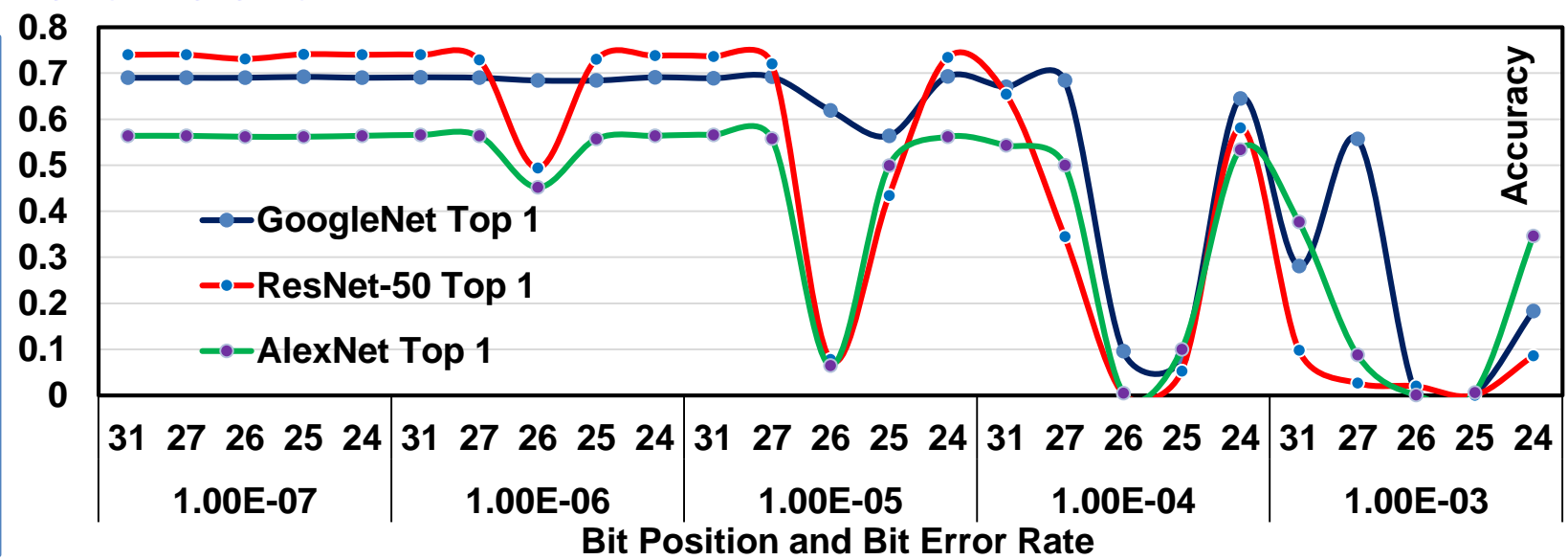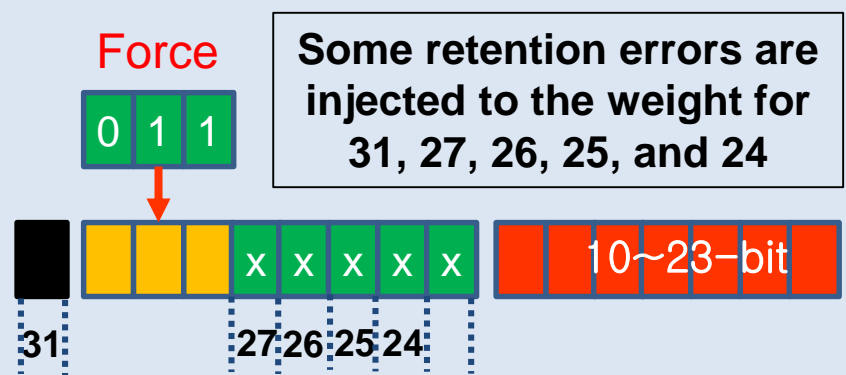(8 MSB: No Errors, Others: $10^{-5}$ BER)
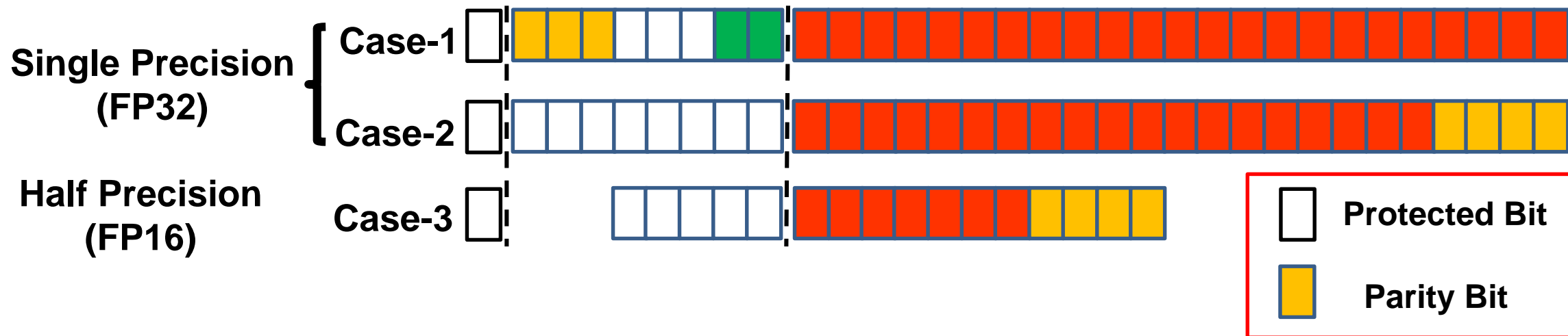
**50 times Error Accumulation**
(All Bits: $10^{-5}$ BER)



- Exponent bits are extremely sensitive to the error

**Simulation for FP32**

Force

| 0 | 1 | 1 |

Some retention errors are injected to the weight for 31, 27, 26, 25, and 24

| 31 | | | | x | x | x | x | x | | 10~23−bit |

31      27 26 25 24



GoogleNet Top 1
ResNet-50 Top 1
AlexNet Top 1

Accuracy

| 31 | 27 | 26 | 25 | 24 | 31 | 27 | 26 | 25 | 24 | 31 | 27 | 26 | 25 | 24 | 31 | 27 | 26 | 25 | 24 | 31 | 27 | 26 | 25 | 24 |

1.00E-07          1.00E-06          1.00E-05          1.00E-04          1.00E-03

**Bit Position and Bit Error Rate**

# Our Approach: Significant-Bit Protection

**Single Precision (FP32)**

Case-1

Case-2

**Half Precision (FP16)**

Case-3

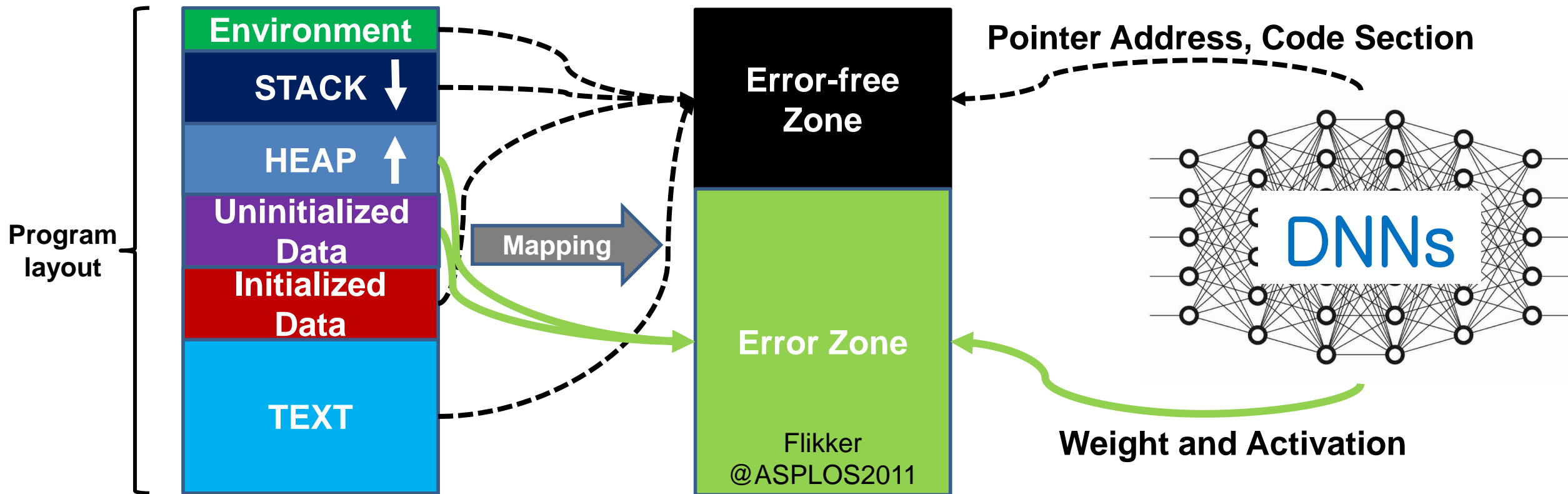| | Protected Bit |
|---|---|
| | Parity Bit |

- Utilize some non-critical bits as parity bits for ECC
  - No additional memory overhead
- Protect some critical bits from retention errors
- Hamming codes for ECC
  - Hamming(7,4) : Case-1
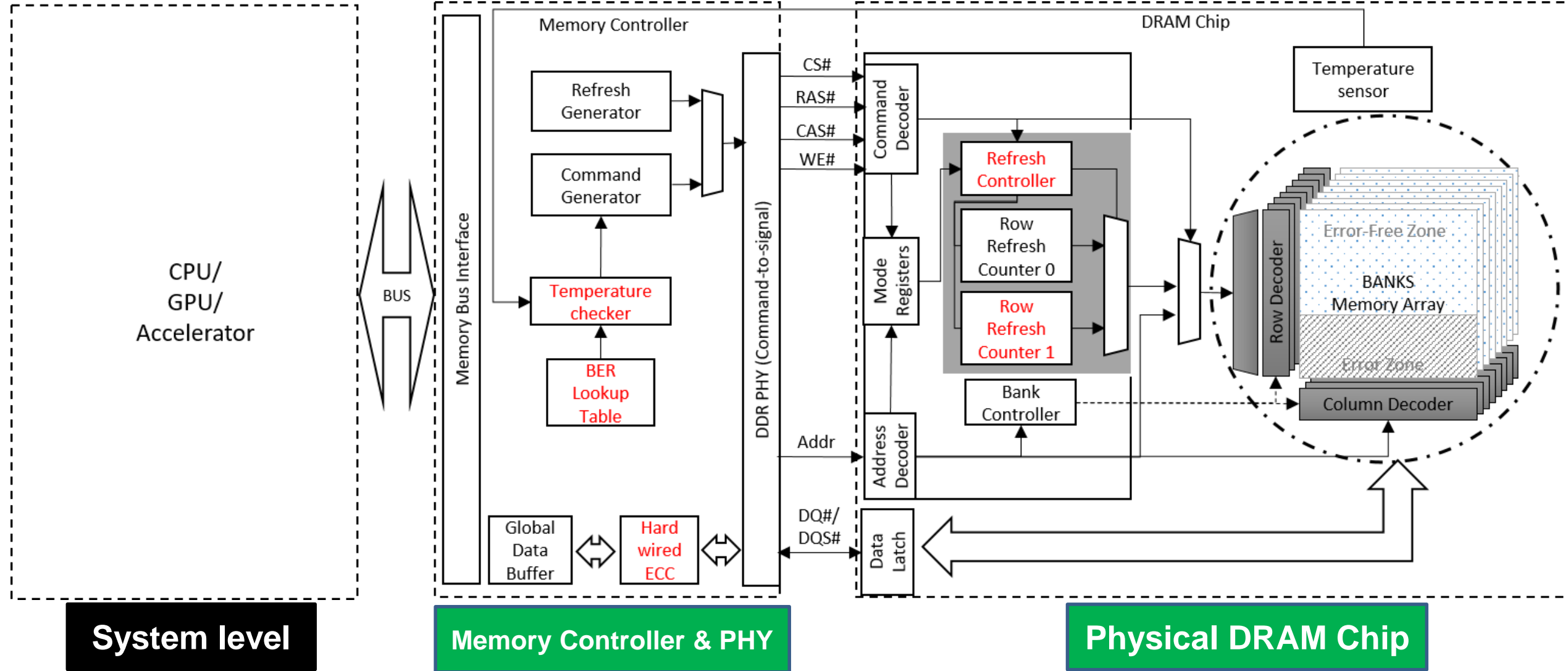  - Hamming(15,11) : Case-2 and Case-3

Legend:
- Hamming(7,4)
- Hamming(15,11)

Y-axis: BER After ECC (1E+00, 1E-02, 1E-04, 1E-06, 1E-08, 1E-10, 1E-12, 1E-14, 1E-16)

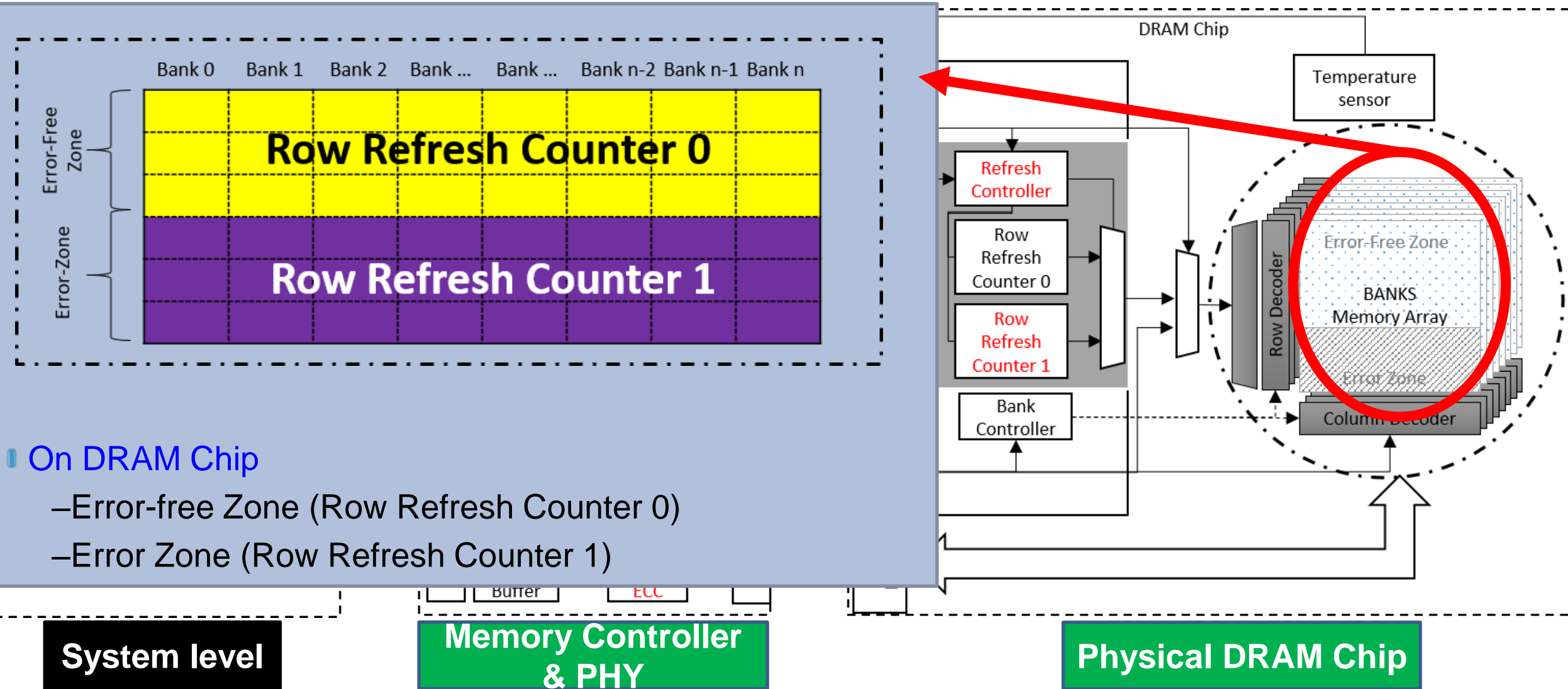X-axis: BER Before ECC (1E-8, 1E-7, 1E-6, 1E-5, 1E-4, 1E-3, 1E-2, 1E-1)

- Text, initialized data, environment, and stack sections are stored in the 'Error-free Zone'.
  - **Prevent retention errors for control information**
- Weight and activation parameters (significantly dominant in DNN) are stored in the 'Error Zone'.
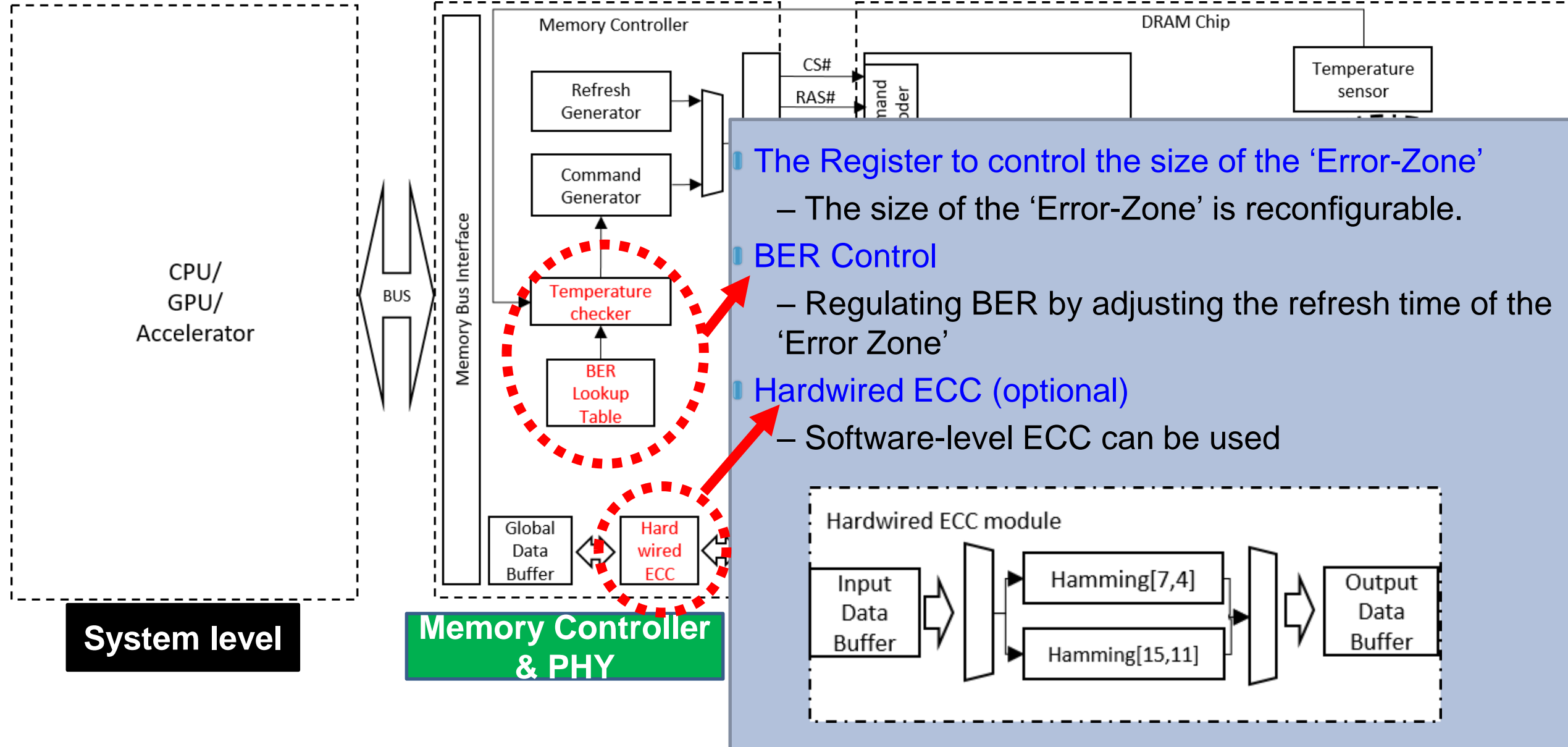
# DRAM Controller with Stretchable Refresh Period



**System level**

**Memory Controller & PHY**

**Physical DRAM Chip**
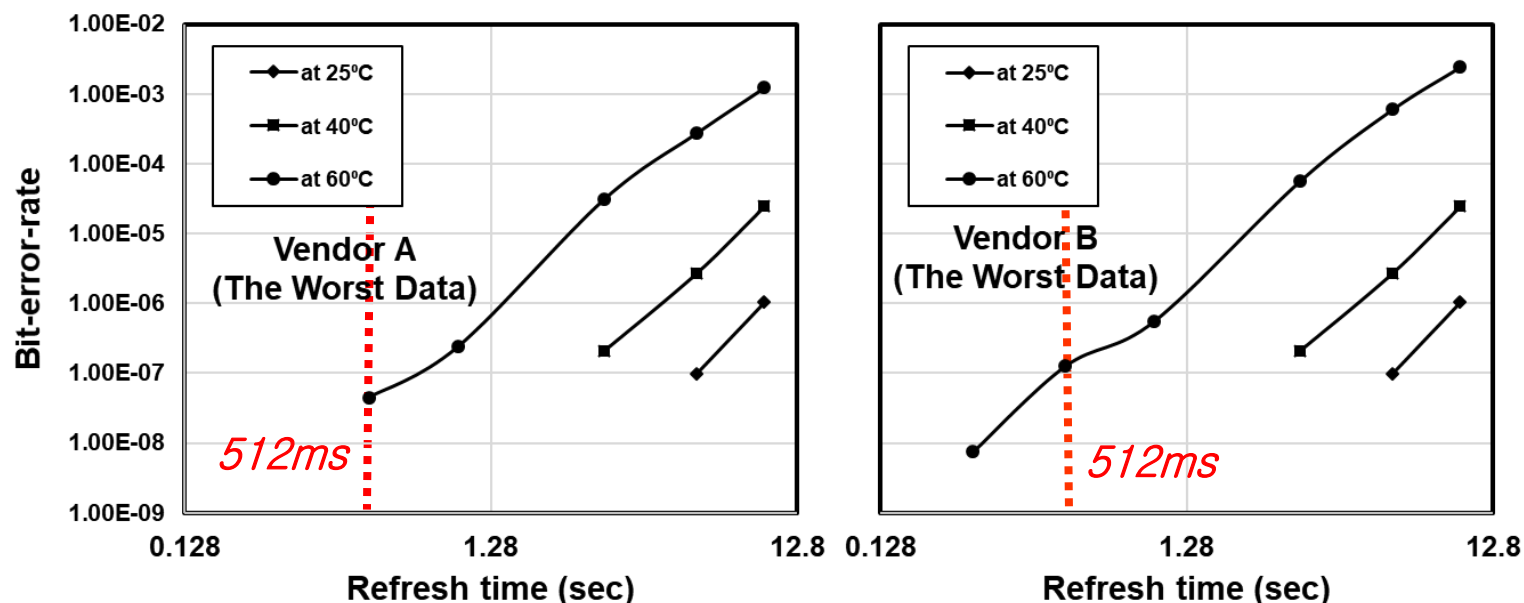
**We slightly modify the Memory Controller and the Physical DRAM Chip**

On DRAM Chip

- Error-free Zone (Row Refresh Counter 0)
- Error Zone (Row Refresh Counter 1)

**System level**

**Memory Controller & PHY**

**Physical DRAM Chip**

# DRAM Controller with Stretchable Refresh Period (Cont..)



- **The Register to control the size of the 'Error-Zone'**
  – The size of the 'Error-Zone' is reconfigurable.
- **BER Control**
  – Regulating BER by adjusting the refresh time of the 'Error Zone'
- **Hardwired ECC (optional)**
  – Software-level ECC can be used

**System level**

**Memory Controller & PHY**

# Validation - Setup



**Retention Errors vs. Refresh Time (Measurement Results)**



512ms based line

Batch Size 16/32/64

- The working temperature is less than 60ºC in the data-center (**)
    - BER ~ $10^{-7}$ at 60°C for two major vendors (DDR3)
    - Inject $10^{-7}$ BER to weights and activations during forward/backward phase
- (Forward time+ backward time)/iteration << 512ms  (Measurement on Torch and Caffe)
- JEDEC requires $10^{-15}$ BER @ 32/64ms in the normal working temperature

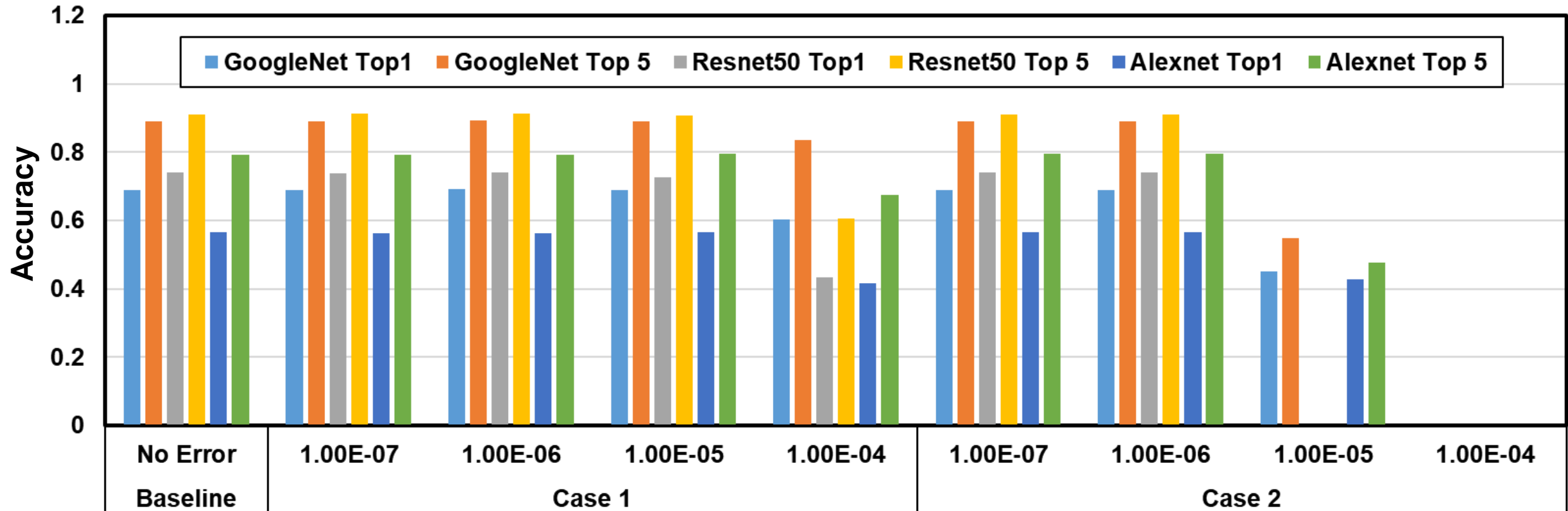*(**)Donghyuk Lee, Adaptive-latency DRAM: Optimizing DRAM timing for the common-case. HPCA 2015*

**ResNet-50 @ $10^{-7}$ BER**

Legend:
- Original Accuracy Top 1
- Case 1 Accuracy Top 1
- Case 2 Accuracy Top 1
- Case 3 Accuracy Top 1
- Without ECC Accuracy Top 1

Works accurately and more efficiently than JEDEC standard, which is $10^{-15}$ BER

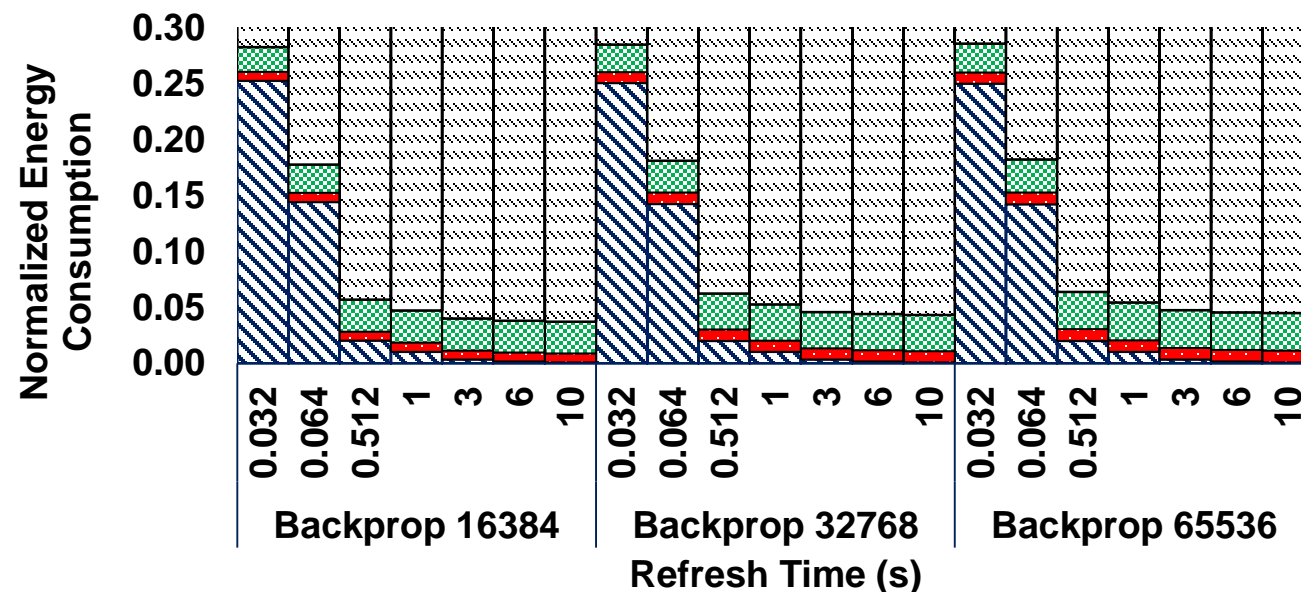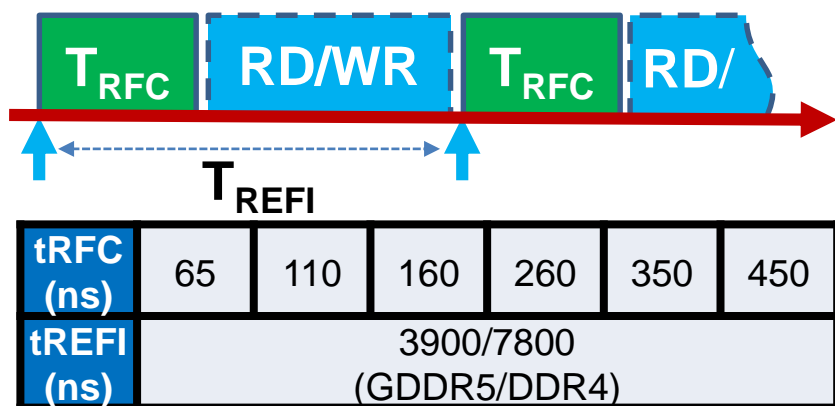| CNNs | Last Trained Epoch | Original Acc. Top-1 | Case-1 Acc. Top-1 | Case-2 Acc. Top-1 | Case-3 Acc. Top-1 |
|---|---|---|---|---|---|
| LeNet | 10 | 99.07 | 99.12 | 99.06 | 99.02 |
| ConvNet | 10 | 75.69 | 74.66 | 76.56 | 75.83 |
| SqueezeNet | 7 | 58.50 | 58.83 | 58.34 | 58.21 |
| GoogleNet | 22 | 70.05 | 69.55 | 69.73 | 70.07 |

- We figure out the safe BER threshold for the inference of each DNN
    - → The Case-1 is a little robust than The Case-2
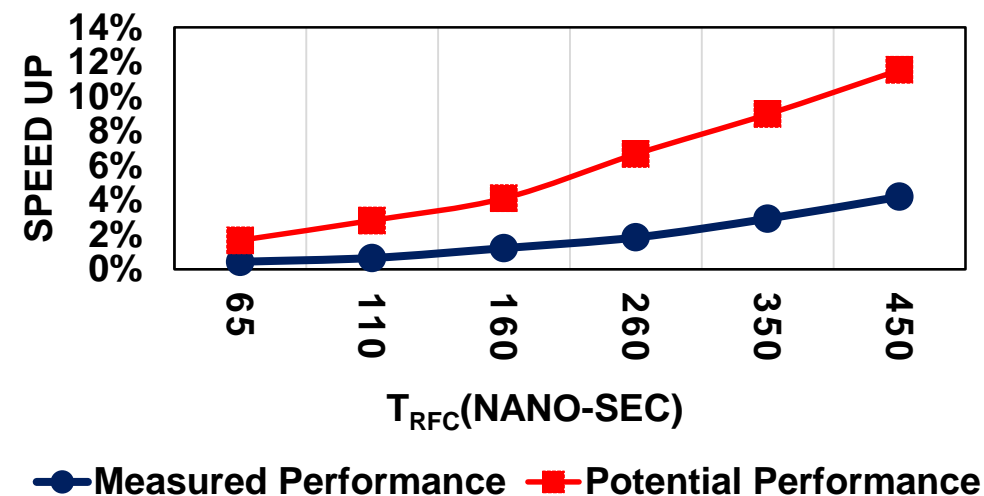    - → The Case-3 is almost similar to the Case-2

# Energy and Performance Simulation Results

| $T_{RFC}$ | RD/WR | $T_{RFC}$ | RD/ |
|---|---|---|---|

$T_{REFI}$

| tRFC (ns) | 65 | 110 | 160 | 260 | 350 | 450 |
|---|---|---|---|---|---|---|
| tREFI (ns) | 3900/7800 (GDDR5/DDR4) | | | | | |

Normalized Energy Consumption vs Refresh Time (s)

X-axis values: 0.032, 0.064, 0.512, 1, 3, 6, 10 for each of: Backprop 16384, Backprop 32768, Backprop 65536

Legend: RefE, Act/PreE, Rd/WrE, ActBack/PreBackE

- Back-propagation test results with different tRFC/tREFI
(Benchmark Set: Rodinia from IISWC 2009)
  - Hybrid CPU-GPU Platform (GPGPU + GEM5)
  - Refresh time is 512ms ($10^{-7}$ BER), DRAM energy reduces:
    - 23% on graphic memories
    - 12% on main memories
  - Performance improves 0.43~4.12%

SPEED UP vs $T_{RFC}$(NANO-SEC): 65, 110, 160, 260, 350, 450

Legend: Measured Performance, Potential Performance

# Conclusion

- Present the stretchable DRAM Refresh controller to control the BER according to
  - Temperature
  - User desired BER
- The proposed Error Correction Schemes can safeguard the important bit
- With only 512ms ($10^{-7}$ BER), our proposed system can potentially help:
  - Speed up the training process up to 4.12%
  - Reduce 12% and 23% DRAM energy in main and graphic memories

# Acknowledgment & Thank you

- Any question?