

# HEALM: Hardware-Efficient Approximate Logarithmic Multiplier with Reduced Error

**Shuyuan Yu, Maliha Tasnim and Sheldon Tan**

VSCLAB

**Department of Electrical and Computer Engineering  
University of California, Riverside**



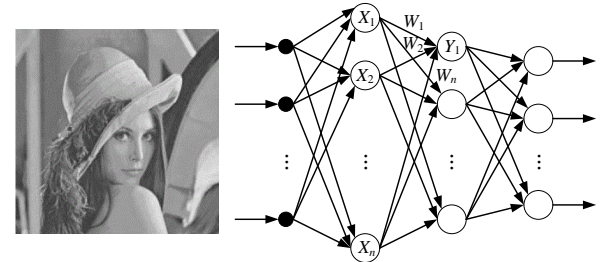
# Outline

- **Background**
- **Related Works & Motivation**
- **HEALM Method**
- **Results & Discussion**
  - **Error Evaluation**
  - **Hardware Performance**
  - **DCT application**
- **Conclusion**

# Background

## More emerging workloads and application are error tolerant

- Image process and computer vision (DCT, edge detection, contrast stretch, etc.)
- Machine learning (DNN, CNN, Transformer, Language model, etc.)
- Applications interfacing with human being do not need exact values



**Approximate computing, which allows the trade-off between area, delay, and power, is more efficient for error tolerant applications.**



# Background

## Approximate Multipliers

- Ad-hoc based
- **Inexact adder based**
- Smaller multiplication with reduced precision
- **Log-based**

## Approximate Log-based Multiplier (ALM)

**Concept:**

$$\mathbf{a} = 2^{ka} \cdot (1+x)$$

$$\mathbf{b} = 2^{kb} \cdot (1+y)$$

$$\begin{aligned} C_{Exact} &= 2^{k_a+k_b} \cdot (1+x) \cdot (1+y) \\ &= 2^{k_a+k_b} \cdot (1+x+y + \mathbf{X}) \end{aligned}$$

$$C_{ALM} = \begin{cases} 2^{k_a+k_b} \cdot (1+x+y), & x+y < 1, \\ 2^{k_a+k_b+1} \cdot (x+y), & x+y \geq 1 \end{cases}$$

## Advantages for ALM:

- Scalability (same Rel. error for 8-bit, 16-bit, 32-bit)
- Area/Power/Energy efficient
- Acceptable accuracy: **3.76%** (Mean Rel.) & **11.11%** (Peak Rel.)



# Related Work & Motivation

## State of art works

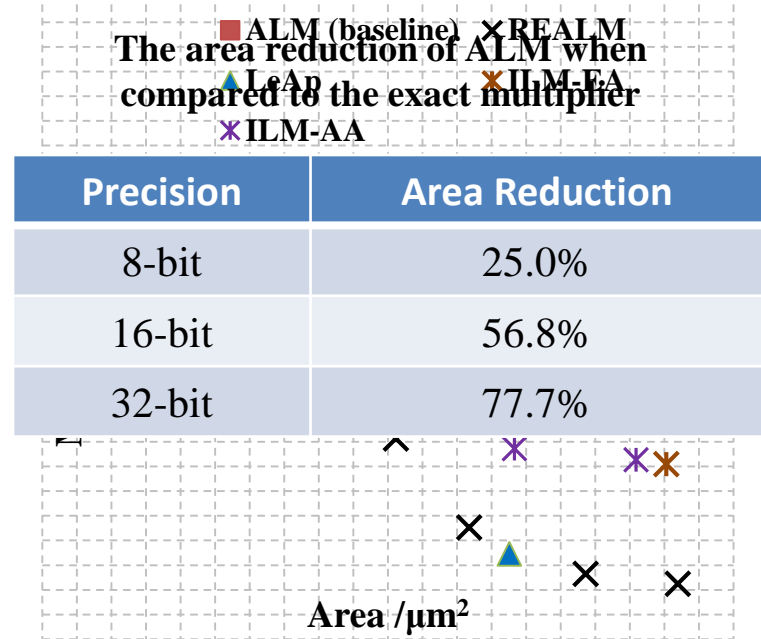
- LeAp [15]
- ILM-EA, ILM-AA [13]
- REALM [10]

## Problems

- Improve accuracy with resource overhead
- Save resource with accuracy loss
- Lower precision, less resource saving

**Motivation:** propose 8-bit logarithmic-based multiplier with resource saving and accuracy improvement at the same time

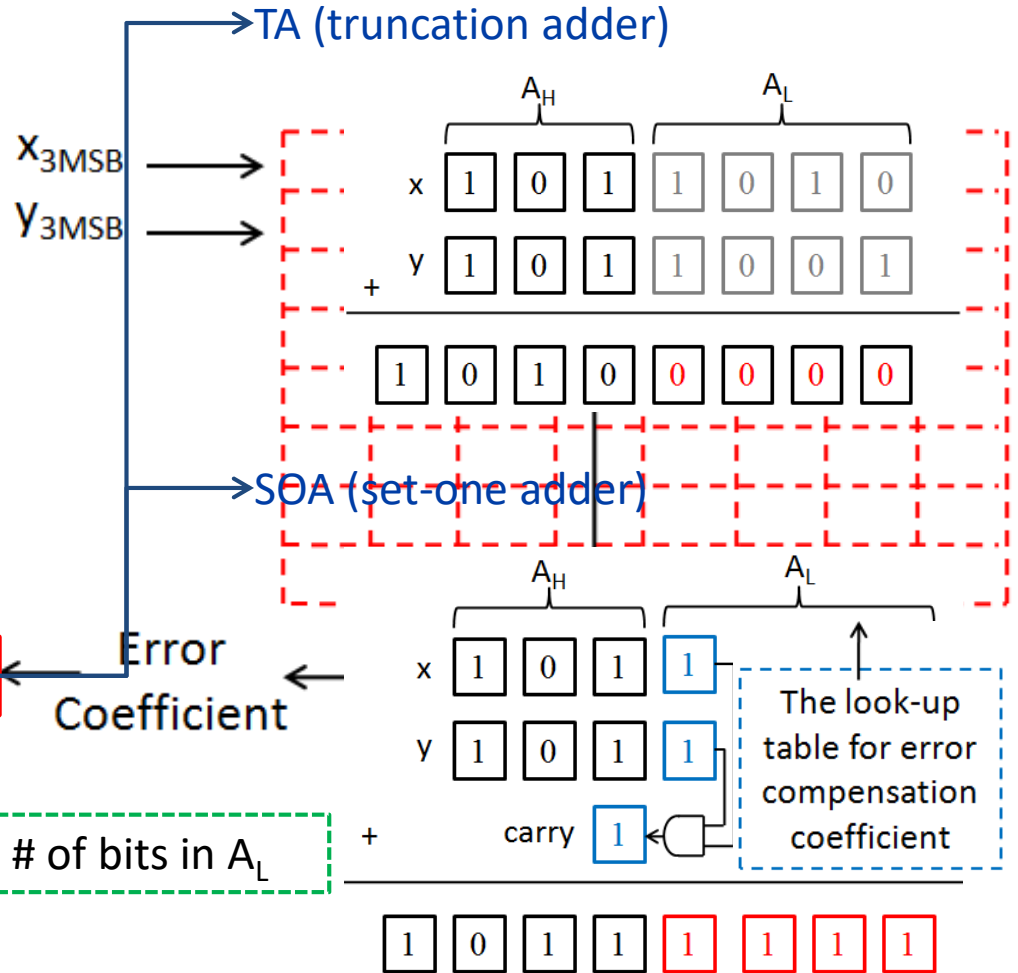
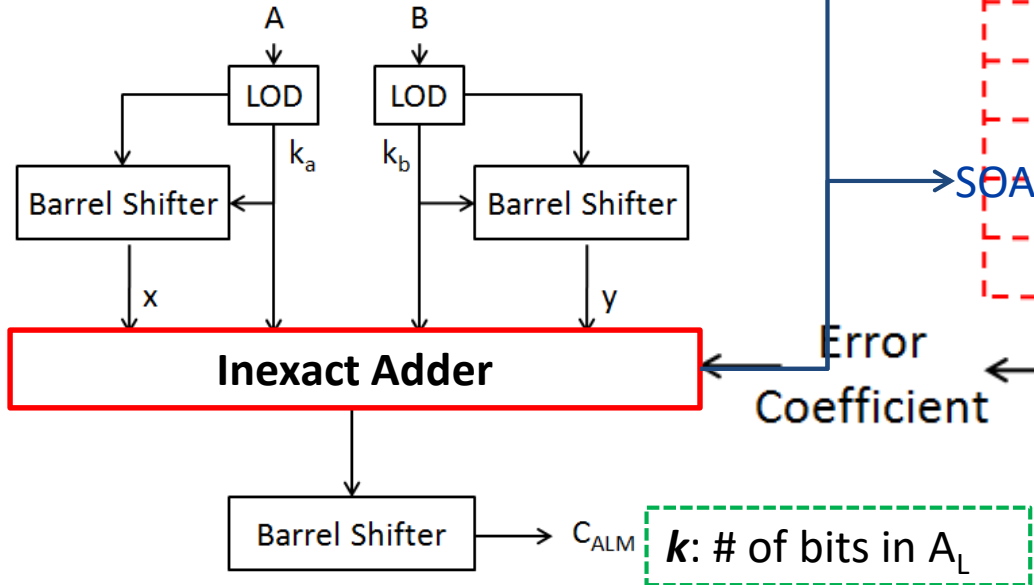
State of art works with 8-bit precision compared with ALM baseline (Area vs. Mean Rel.)





# HEALM Design

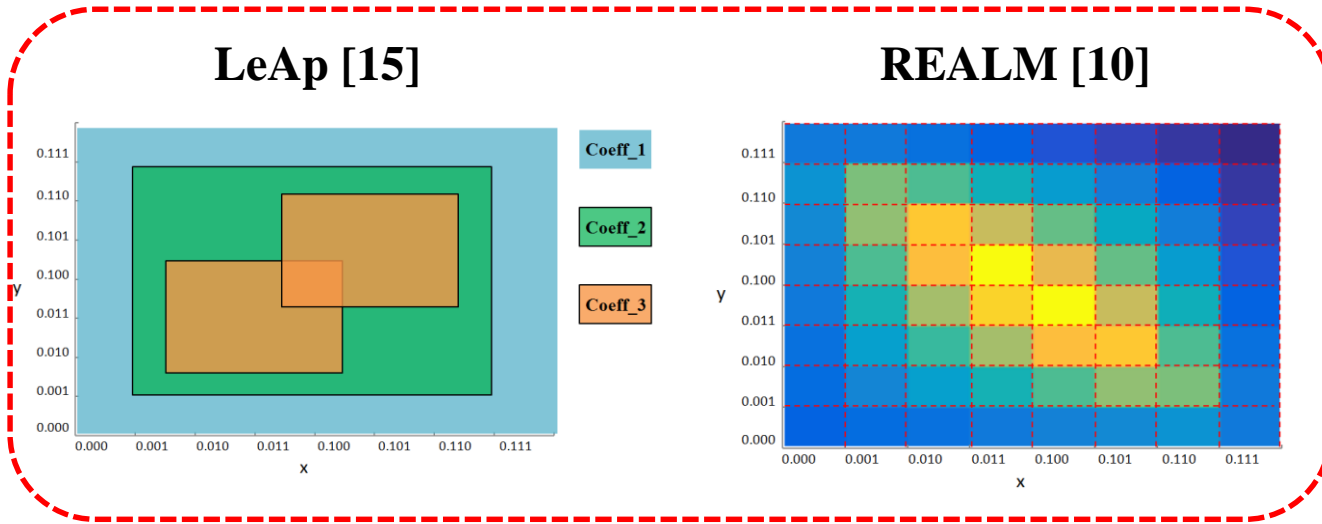
## ALM structure





# Error Compensation

State of art improved ALMs



**ONE  
Pattern**

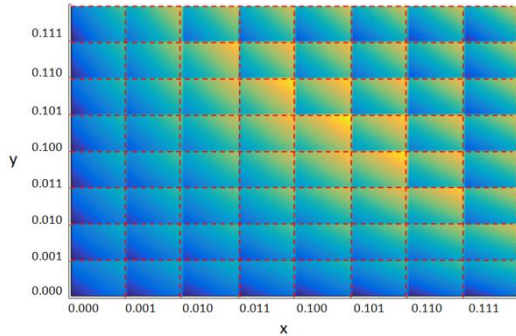
**NOT** suitable when we introduce inexact summation unit with **8-bit** precision.



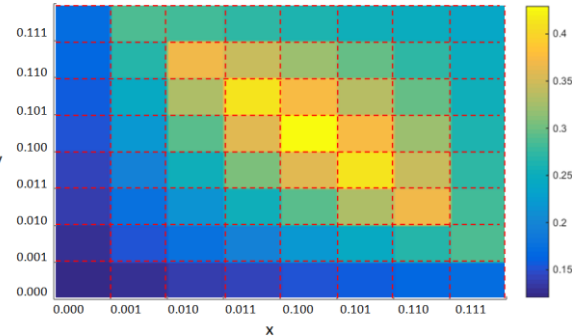
# HEALM Method

Generate error coefficient table based on different value of  $k$

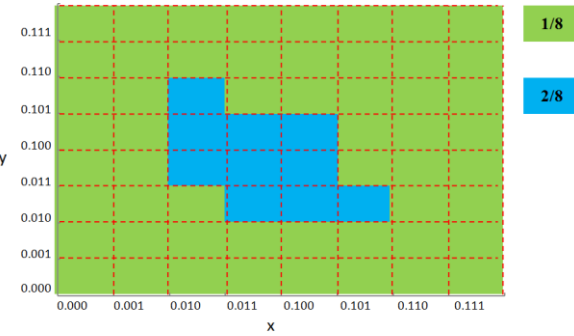
Error profile with inexact adder & error compensation



**Error profile**



**Error profile (avg. in block)**



**Error coefficient**

**# of bits of coeff**  
 $\leq$   
**# of bits of  $\Delta_H$**

**HEALM-TA**

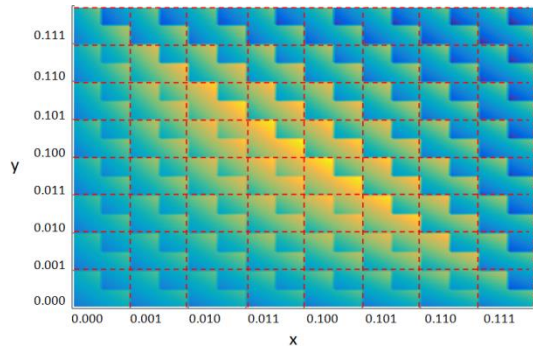




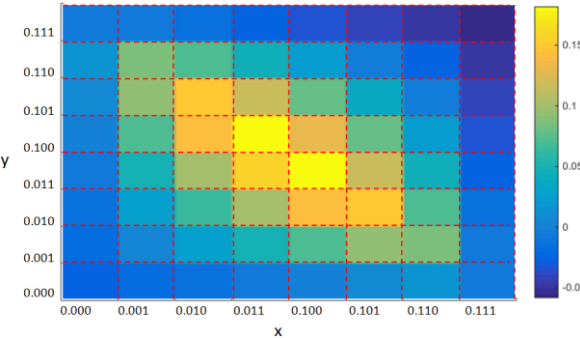
# HEALM Method

Generate error coefficient table based on different value of  $k$

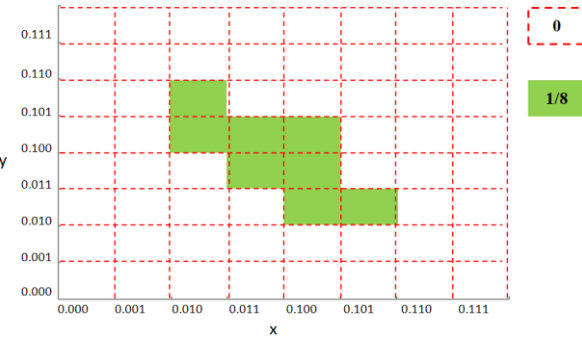
Error profile with inexact adder & error compensation



**Error profile**



**Error profile (avg. in block)**



**Error coefficient**

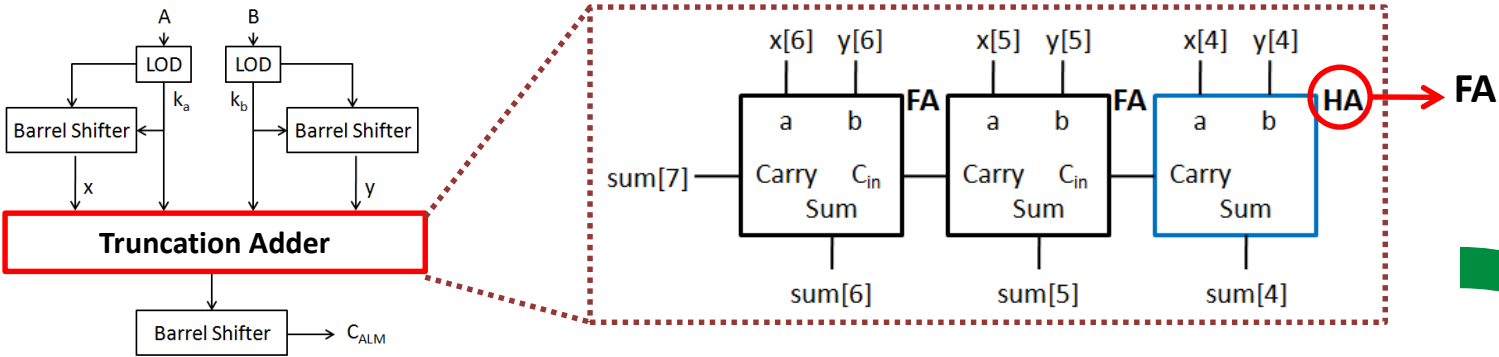
**Ex.  $k = 4$**

**HEALM-SOA**



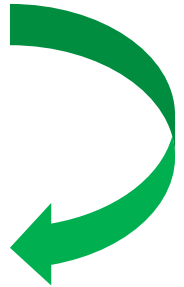
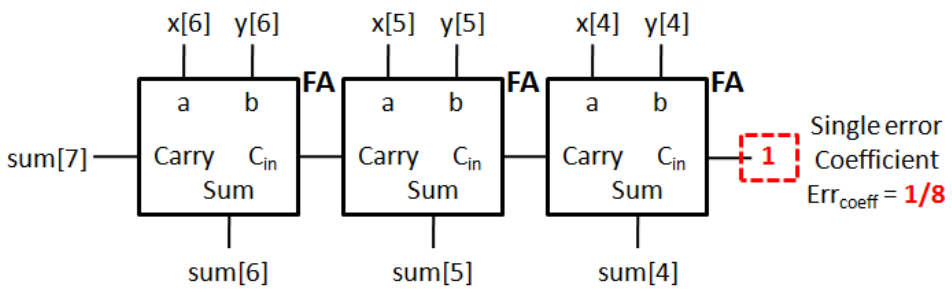
# Single Coefficient Mode (HEALM-TA-S)

Ex. With largest value of  $k$  (8-bit,  $k=4$ )



For **ALL** the inputs, introduce just **ONE** single error coefficient

$x[6:0] + y[6:0]$   
 $k = 4$ , **truncate 4 bits**  
 $x[6:4] + y[6:4]$





# Results & Discussion

## Experimental setup

- Baseline: ALM

- EDK 32nm standard cell library (SNPS DC)

- 400MHz working frequency

- 1 million random input pairs (Matlab)

## Error Metrics & Hardware Performance

Logarithmic MUL	$k$	Mean /%	Peak /%	Area / $\mu\text{m}^2$	Power / $\mu\text{W}$
<b>ALM (baseline)</b>	<b>0</b>	<b>3.76</b>	<b>11.11</b>	<b>820.63</b>	<b>72.83</b>
<b>HELAM-TA</b>	1	1.12	4.86	1216.84	114.50
	2	1.40	8.25	938.05	92.17
	<b>3</b>	<b>2.17</b>	<b>9.75</b>	<b>743.63</b>	<b>74.14</b>
	4	3.66	13.77	595.46	51.41
<b>HEALM-SOA</b>	1	1.13	4.71	1175.41	113.96
	2	1.38	5.90	966.76	90.32
	<b>3</b>	<b>1.78</b>	<b>7.65</b>	<b>808.94</b>	<b>75.05</b>
	4	3.12	12.17	664.33	59.55



# Results & Discussion

## Experimental setup

- Baseline: ALM

- EDK 32nm standard cell library (SNPS DC)

- 400MHz working frequency

- 1 million random input pairs (Matlab)

## Error Metrics & Hardware Performance

Logarithmic MUL	$k$	Mean /%	Peak /%	Area / $\mu\text{m}^2$	Power / $\mu\text{W}$
<b>ALM (baseline)</b>	<b>0</b>	<b>3.76</b>	<b>11.11</b>	<b>820.63</b>	<b>72.83</b>
<b>ALM-TA</b>	1	4.02	12.03	763.45	69.41
	2	4.79	13.83	702.71	65.07
	3	6.58	17.25	612.74	56.20
	<b>4</b>	<b>10.29</b>	<b>23.53</b>	<b>533.19</b>	<b>41.73</b>
<b>HEALM-TA-S</b>	1	3.21	11.11	763.70	71.19
	2	3.17	12.02	716.69	65.68
	3	3.39	13.79	614.01	56.62
	<b>4</b>	<b>4.26</b>	<b>17.12</b>	<b>534.72</b>	<b>42.32</b>



**6.0%**

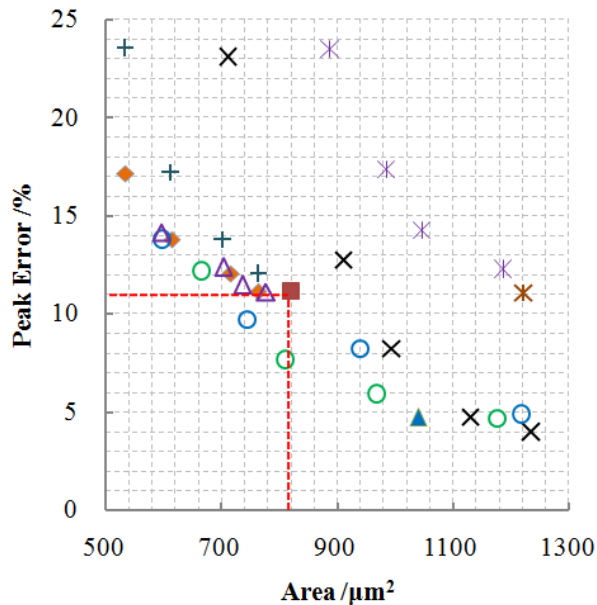
**6.4%**

**Almost NO overheads**

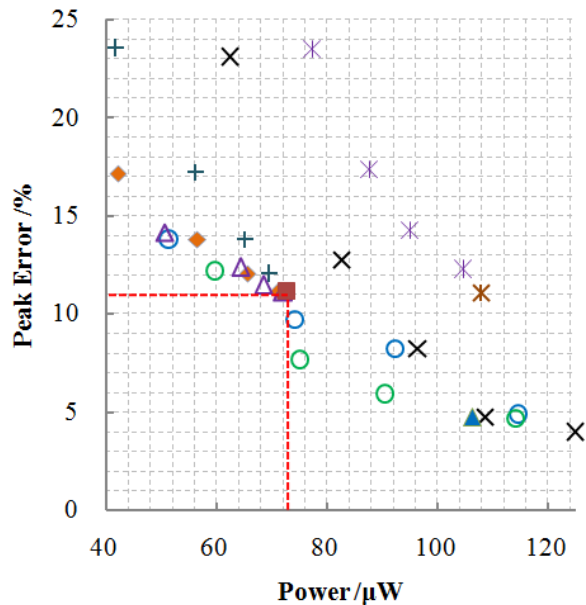


# Comparison with State of Art Works

(a) Area vs. Peak Error



(b) Power vs. Peak Error



■ ALM

▲ LeAp

× REALM

✱ ILM-EA

✱ ILM-AA

+ ALM-TA

◆ HEALM-TA-S

○ HEALM-TA

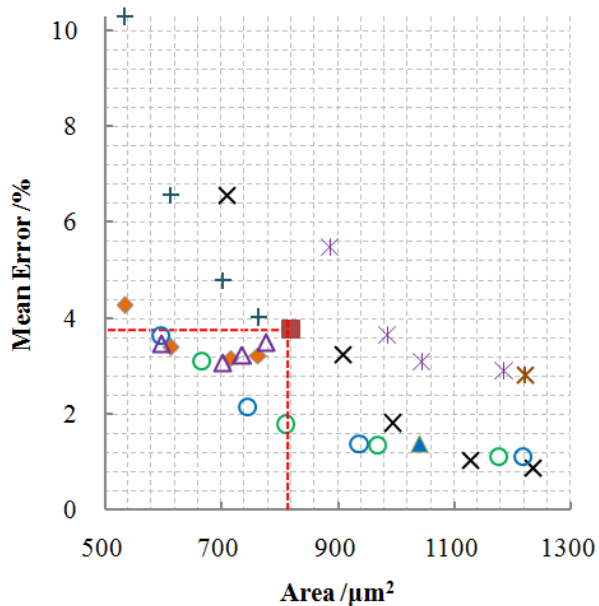
△ ALM-SOA

○ HEALM-SOA

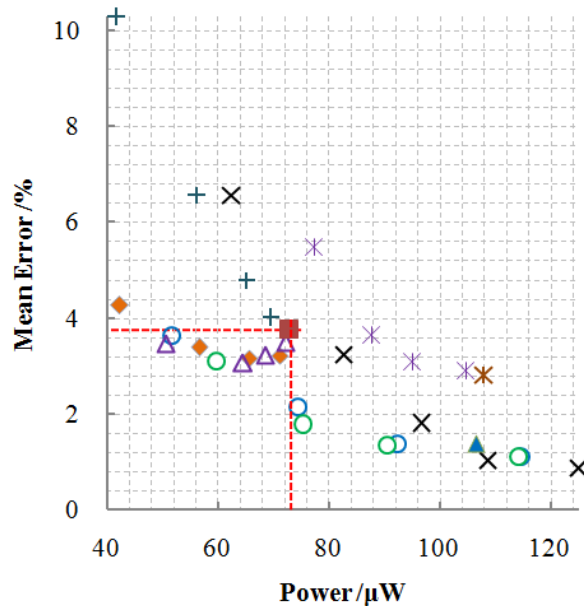


# Comparison with State of Art Works

(c) Area vs. Mean Error



(d) Power vs. Mean Error



■ ALM

▲ LeAp

× REALM

✖ ILM-EA

✖ ILM-AA

+ ALM-TA

◆ HEALM-TA-S

○ HEALM-TA

△ ALM-SOA

○ HEALM-SOA



## Complementary Results for **16-bit** Cases

### Experimental setup

### Error Metrics & Hardware Performance

- Baseline: ALM

- EDK 32nm standard cell  
library (SNPS DC)

- 200MHz working  
frequency

- 1 million random input  
pairs (Matlab)

Logarithmic MUL	$k$	Mean /%	Peak /%	Area / $\mu\text{m}^2$	Power / $\mu\text{W}$
<b>ALM (baseline)</b>	<b>0</b>	<b>3.76</b>	<b>11.11</b>	<b>1825.52</b>	<b>110.91</b>
REALM	0	0.75	3.70	2383.36	164.50
	<b>9</b>	<b>1.06</b>	<b>5.27</b>	<b>1572.90</b>	<b>94.07</b>
LeAp	0	0.98	4.76	1990.71	128.20
ALM-TA	9	4.88	12.93	1263.86	66.26
<b>HEALM-TA-S</b>	<b>9</b>	<b>4.87</b>	<b>12.02</b>	<b>1267.16</b>	<b>66.45</b>
<b>HEALM-TA</b>	<b>9</b>	<b>1.64</b>	<b>5.83</b>	<b>1511.39</b>	<b>87.62</b>
ALM-SOA	9	3.07	12.03	1383.56	74.10
<b>HEALM-SOA</b>	<b>9</b>	<b>1.38</b>	<b>5.15</b>	<b>1577.47</b>	<b>91.89</b>



# Application: Discrete Cosine Transform

Image quality: PSNR (dB)

MUL	Lena	Boat	Barbara	House	Pepper	Avg.
ALM (baseline)	19.1	18.7	19.3	18.4	18.7	18.8
<b><math>k = 4</math></b>						
ILM-AA	20.2	18.9	20.1	18.5	19.5	19.5
ALM-TA	14.1	13.7	14.1	13.4	13.7	13.8
ALM-SOA	18.9	18.9	19.6	16.7	18.9	18.6
REALM	19.8	20.0	19.8	17.8	19.5	19.4
<b>HEALM-TA</b>	<b>29.1</b>	<b>28.7</b>	<b>29.1</b>	<b>25.9</b>	<b>28.5</b>	<b>28.3</b>
HEALM-TA-S	23.1	22.0	22.6	24.3	22.3	22.9
HEALM-SOA	22.2	22.1	22.6	19.4	22.4	21.7

Up to  
**17.2 dB**

At least  
**2.9 dB**





# Conclusion

- Propose **HEALM** approach: Hardware-Efficient Approximate Logarithmic Multiplier, with reduced error.
- Do improvements upon the ALM baseline in **error metrics** and **hardware performance** aspects **at the same time**.
- Outperform the state of art design with up to **2.9%, 9.3%, 16.1%, 17.6%** improvement in mean error, peak error, area, power consumption, respectively for 8-bit precision.
- The single coefficient mode **HEALM-TA-S** could improve up to **6.0%** and **6.4%** in mean and peak error respectively when compared to ALM with simple TA, with almost **NO** resource overheads.
- For DCT workloads, with  $k=1$ , **HEALM design** could achieve up to **17.2dB** improvement upon ALM baseline; with  $k=4$ , **HEALM design** could achieve at least **2.9dB** improvement in image quality.

Questions???