



FPGA Implementation of In-Memory Search Macro with Range Encoding

Speaker : B11901027 王仁軒
Mentor : Rick Huang
Advisor: Prof. An-Yeu (Andy) Wu
Date : 2025/05/02

ACCESS IC LAB

1



Outline

- ❖ Introduction to few-shot learning
- ❖ BORE range encoding scheme
- ❖ IMS system architecture
- ❖ Implementation on FPGA
- ❖ Result

- ❖ Embedding dimension reduction
- ❖ Analysis of PCA (Principle Component Analysis)
- ❖ Analysis of quantization method
- ❖ Modification of AutoEncoder method

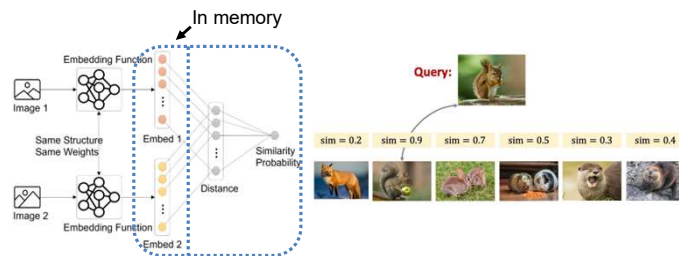
P2

2



Few-Shot Learning

- ❖ In few-shot learning, we compare similarity between query and supporting data in memory
 - ❖ Process data with NN, store embeddings into memory
 - ❖ Compute similarity between query and supporting embeddings



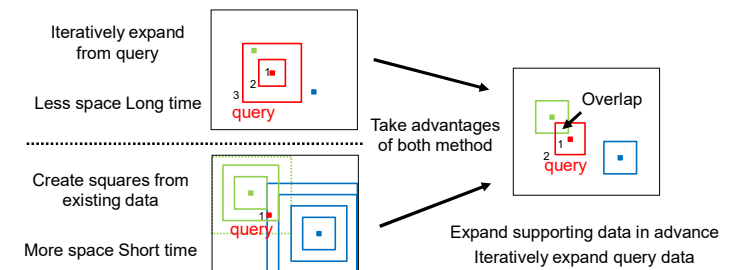
P3

3



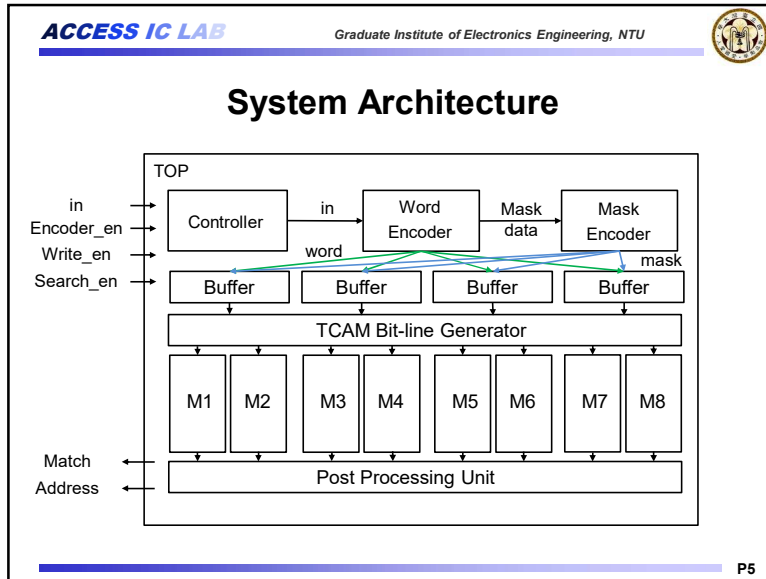
Similarity Search

- ❖ Which class does query belong to?
 - ❖ Iterative vs one-shot method

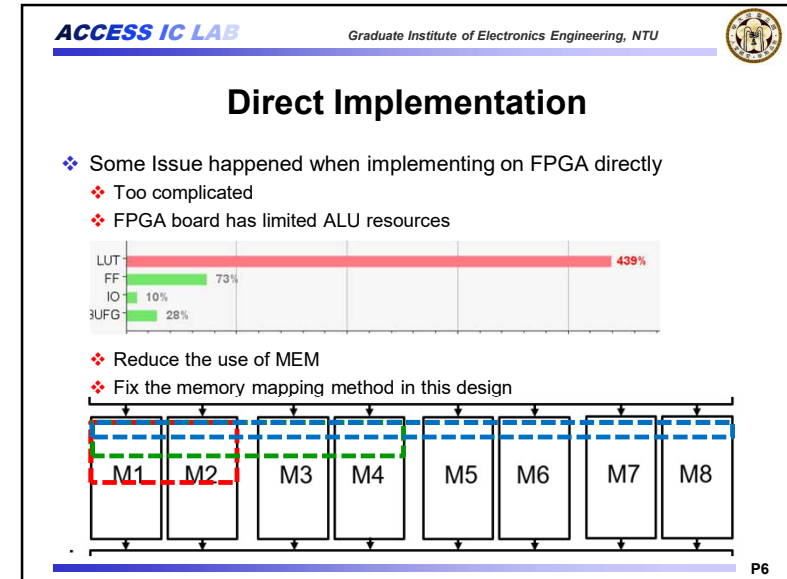


P4

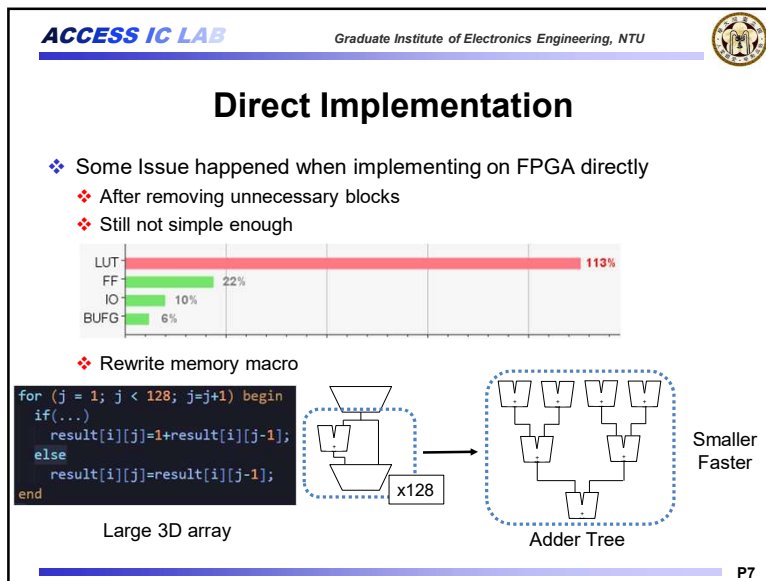
4



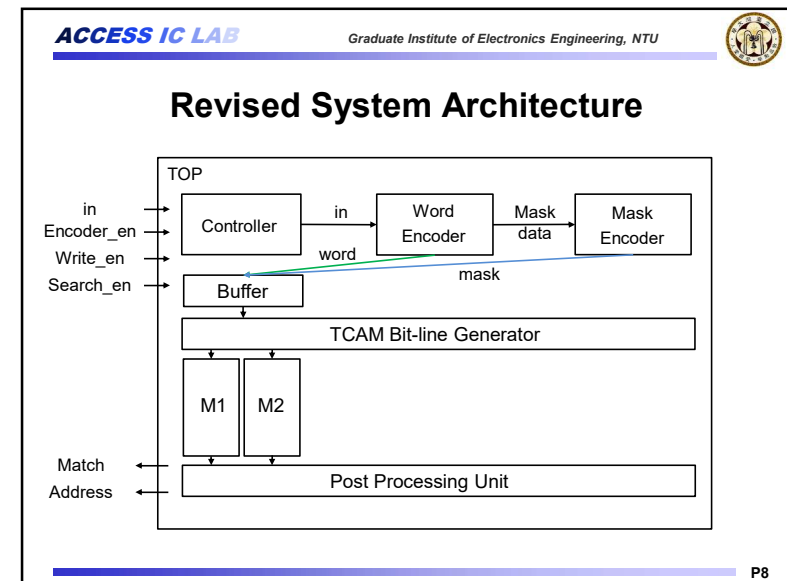
5



6



7



8

ACCESS IC LAB Graduate Institute of Electronics Engineering, NTU

Input Stimuli Generation

- ❖ Input stimuli were provided by testbench → Non-synthesizable
- ❖ Rewrite testbench into a synthesizable module
- ❖ Using RAM on FPGA board to preload data

P9

9

ACCESS IC LAB Graduate Institute of Electronics Engineering, NTU

Result

- ❖ After start signal is sent, TB receives MATCH signal from TOP once Search SW is triggered (active low)
- ❖ LED shows the ADDRESS signal

P10

10

ACCESS IC LAB Graduate Institute of Electronics Engineering, NTU

Curse of Dimensionality

- ❖ As dimensions increases, data become increasingly sparse, causing traditional algorithms to struggle

- ❖ The number of dimensions should meet the HW provided by MXIC
- ❖ 64-dimensional INT3 L2-norm search

P11

11

ACCESS IC LAB Graduate Institute of Electronics Engineering, NTU

Methods of Dimension Reduction

- ❖ Linear
 - ❖ Pooling
 - Average Pooling
 - Magnitude Pooling
 - ❖ Principle Component Analysis
- ❖ Non-Linear
 - ❖ UMAP
 - ❖ Autoencoder (AE)

Closeer points → exponential decay to maintain local distance
Farther points → inverse-polynomial to maintain global distance

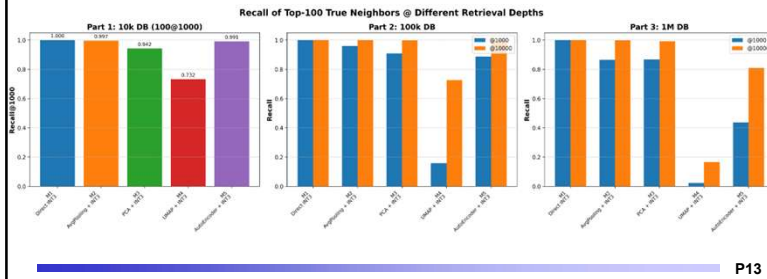
P12

12



Experimental Results (SIFT-1M)

- ❖ Using SIFT-1M dataset
 - ❖ Vector dimension: 128 -> 64
 - ❖ Quantization method: min-max quantization
 - ❖ Avg. Pooling > PCA > AE > UMAP

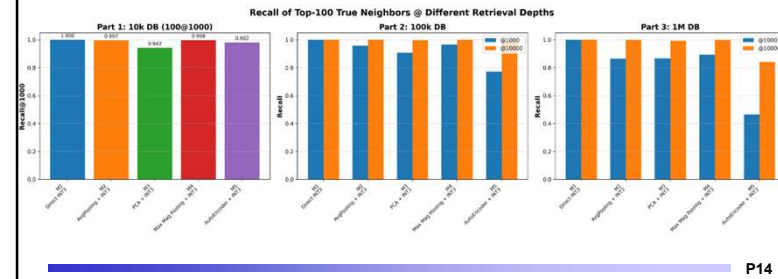


13



Experimental Results (SIFT-1M)

- ❖ Using SIFT-1M dataset
 - ❖ Vector dimension: 128 -> 64
 - ❖ Quantization method: min-max quantization
 - ❖ Since the performance of UMAP is worst, it is replaced by magnitude pooling
 - ❖ Avg. Pooling = Mag. Pooling >= PCA >= AE



14



Experimental Results (Omniglot)

- ❖ Using Omniglot 20-way-5-shot experiment
 - ❖ Embedding dimension: 128 -> 64
 - ❖ Quantization method: min-max quantization
 - ❖ AE = Avg. Pooling = PCA >= Mag. Pooling

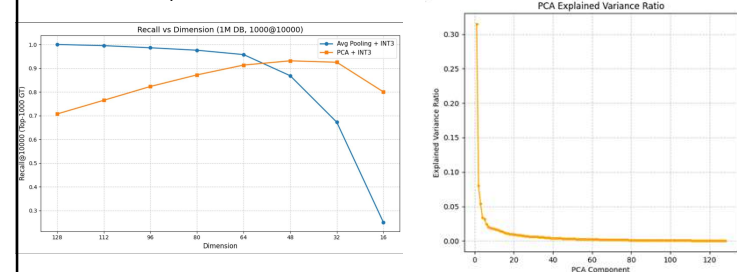


15



Analysis of PCA

- ❖ PCA seems cannot win against average pooling
 - ❖ Is this true for all dimensions?
 - ❖ Why PCA achieves better recall rate when dim < 60
 - ❖ Too much noise from less representative dimensions
 - ❖ Low explained var. from dim 60 to 128

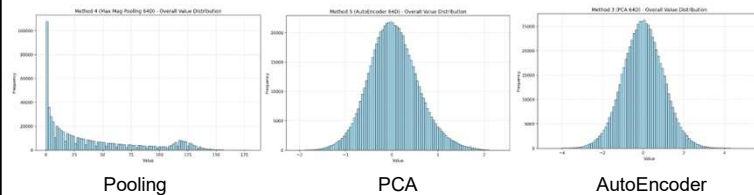


16



Data Distribution under Different Methods

- ❖ What does the distributions look like after different dimension reduction methods
- ❖ Some methods produces long-tail distributions
- ❖ Others produces bell-shaped distributions
- ❖ Can our quantization method fit all of these distributions?



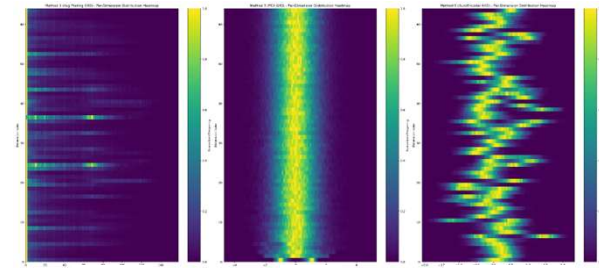
P17

17



Data Distribution under Different Methods

- ❖ What does the distributions look like after different dimension reduction methods
- ❖ Some methods produces long-tail distributions
- ❖ Others produces bell-shaped distributions
- ❖ Can our quantization method fit all of these distributions?



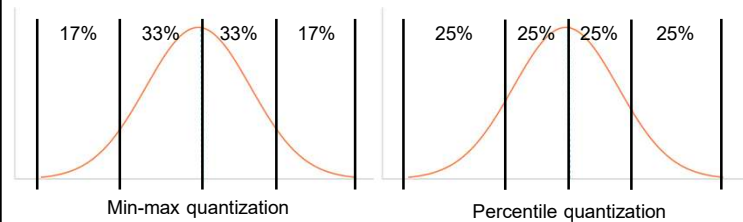
P18

18



Quantization Makes Big Differences

- ❖ New quantization method : Percentile Quantization
- ❖ Quantize data according to its percentile among all data
- ❖ Non-linear quantization
- ❖ Better to distinguish top-k nearest neighbors



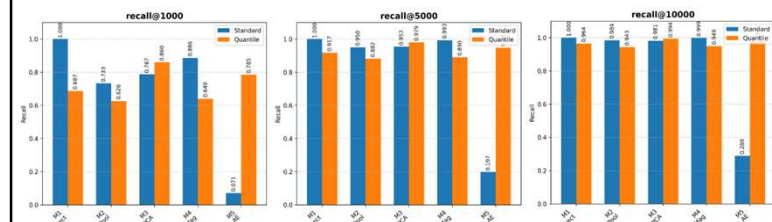
P19

19



Experimental Results

- ❖ Percentile quantization boosts the recall rate of AutoEncoder method
- ❖ Embedding dimension: 128 -> 64
- ❖ Blue: Min-max Quantization. Orange: Percentile Quantization
- ❖ Min-max -> better for long-tail distribution
- ❖ Percentile -> better for bell-shaped distribution



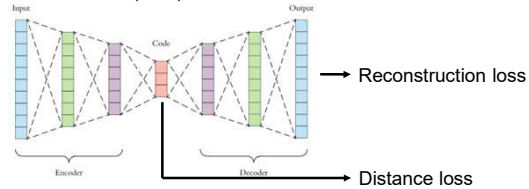
P20

20



Modification of AutoEncoder

- ❖ Original AutoEncoders consider reconstruction loss as loss function
 - ❖ Construct a low dimension space
 - ❖ Can it preserve distance information?
- ❖ New loss function
 - ❖ Distance loss = MSE(distance before and after dimension reduction)
 - ❖ $k \times \text{reconstruction loss} + (1 - k) \times \text{distance loss}$



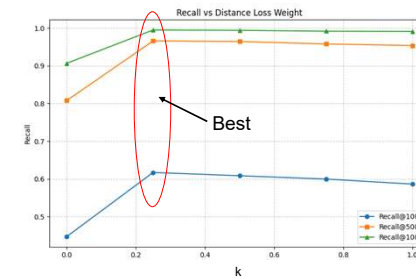
P21

21



Experimental Results

- ❖ New loss function
 - ❖ $k \times \text{reconstruction loss} + (1 - k) \times \text{distance loss}$
 - ❖ Recall rate increased significantly by adding distance loss into loss function
 - ❖ Considering both reconstruction and distance info achieves best recall rate



P22

22



Conclusion

- ❖ Embedding dimension reduction
 - ❖ Average pooling is simple but effective
- ❖ Analysis of PCA (Principle Component Analysis)
 - ❖ PCA outperforms average pooling when dimensions < 60
- ❖ Analysis of quantization method
 - ❖ Min-max quantization is better for pooling methods
 - ❖ Quantile is better for PCA and AutoEncoders
- ❖ Modification of AutoEncoder method
 - ❖ Preserving distance information increases recall rate
 - ❖ Best recall @ reconstruction : distance = 3 : 1

P23

23