**Challenges and Opportunities in Accelerating** Large-Scale Search Using NVM-based In-Memory Computing X. Sharon Hu (shu@nd.edu) **Department of Computer Science and Engineering University of Notre Dame** 







JUMP

# Search as an important class of computation kernels



# Growing use of associative search in AI/ML workloads



# Al and memory wall



'The evolution of the number of parameters (410x) of SOTA AI models over the years, along with the AI accelerator memory capacity (2x)'

# Memory wall in conventional computing platforms



### **Reduce memory transfer overhead is critical!**

In-memory search via content addressable memory

# In-memory search: Content-addressable memory (CAM)

### □ CAM => Associative memory (AM)

- Support parallel search
- >Widely used in high-associativity cache, network router, etc.
- >Growing use in data analytics, machine learning, vector processing



# Choices of memory device technologies for CAMs











#### **FeFETs**

Silicide

Poly-Si

TiN

8nm HfO

SiO<sub>2</sub>

p-Si



RRAM

What memory devices should be used to construct CAM cells?

#### **Desired properties:**

- Large memory windows
- Multi-level data
- **Tunable resistance levels**
- Robustness
- Endurance
- Retention
- Speed, power, area igodol
- ...

NVM opens doors for designing more sophisticated CAM functions

# **CAM** data representations

#### **Generic CAM array structure**



# CAMs based on different memory technologies



[1] Haitong Li, et al., IEEE TED, 2021. [2] Seungchul Jung, et al., Nature, 2022. [3] Jing Li, et al., IEEE JSSC 2013. [4] Viacheslav V. Fedorov, et al., IEEE TC, 2016. [5] Arman Kazemi, et al., IEEE ISLEP, 2021. [6] Giacomo Pedretti, et al., Nat. Comm, 2021. [7] Xunzhao Yin et al., IEEE TED, 2020

# **Different types of match for T/MCAM**



#### Same CAM cell design but different sensing circuits



CAM Sub-Array EN Enable Circuit Delay Element Latch - C<sub>1,d</sub> C<sub>1.1</sub> - C<sub>1,2</sub> Pre-charger ML  $D \rightarrow D$ Latch c<sub>2,2</sub> C2.1 ------ C<sub>2,d</sub> ... ... ML Latch C32.2 C<sub>32</sub> C<sub>32.1</sub> •••• Dataline decoder Dataline driver



#### Arman Kazemi, et al., Sci. Report, 2022

# 2-FeFET MCAM design implement other distances

**Conventonal B/TCAM uses Hamming distance; can other distances be** implemented?

≻Yes!

## $\Box V_T$ and $V_q$ assignment problem:

 $\succ$  Find the optimal  $V_T$  and  $V_Q$  for each state

> Minimize the error between each CAM distance and the desired distance

> Define necessary constraints

CAM distance desired distance  $\min_{V_{Q},V_{T},\alpha} I(V_{Qi},V_{Tj}) - \alpha \times d_{ij} - \beta$ 

Subject to  $0 \leq V_{Qi} \leq V_{DD}$ ,  $V_{Tmin} \leq V_{Tj} \leq V_{Tmax}$ ,  $V_{Ti} \leq V_{Qi}$ ,  $\forall i, j$ 

# CAM for complex and large-scale search

# How big can a single CAM array be?

### □ Sub-array size, particularly # of columns, cannot be too big!



□ How to support high dimensional vectors with small CAM subarrays?

□ How do architectural solutions impact accuracy, energy, latency?

## **Hierarchical CAM architecture**



# Scatter and gather



#### Horizontal Gather When N>C How to combine partial results on different dimensions together?

### Vertical Merge When K>R

How to find the match results across different groups of entries?

Match Type	Horizontal only	Vertical only	Horizontal & Vertical
Exact	AND	<mark>Merge</mark>	AND and Merge
Best	Voting	<b>Comparator</b>	N/A
Threshold	N/A	Merge	N/A

- Merge without error
- Gather with approximation error!
- No effective gather schemes proposed yet!!

## **Cross-layer design choices in CAM-based accelerators**



# **CAMASim: CAM-based accelerator simulator**



https://github.com/ND-IMC/CAMASim-V1.0

# **Case study and look forward**

# Case study: Memory-augmented neural network (MANN)



- A CNN for generating highdimension embeddings
  - Embedding dimensions ranges from 64 to 512
- □ A **CAM**-based accelerator for classification

- High-dimension embedding → Horizonal gather is required
  - > Majority voting used here
- What is the best design configuration considering gather errors?

# Impact of embedding and subarray size



## Impact of device variation and sensing limit



# Looking forward



Use of associative search in retrieval augmented generation for Large Language Model (LLM)