

I VecFlow-Chamfer: A GPU-based Data Management System for High-Performance Multi-Vector Search on Superchips



Chenghao Mo¹, Ben Karsin², Philip Adams³, Minjia Zhang¹

¹: SSAIL Lab, UIUC ²: Nvidia ³: Microsoft

What is Multi-Vector Search?

Doc 1: The **cat** sits on the mat
do d1 d2 d3 d4 d5

Query: Where does the **cat** sit
q0 q1 q2 **q3** q4

Chamfer Score(Query, Doc i) = $\sum_i \max_j (q_i, d_j)$

Query token	Best match	Sim score
Where	mat	0.61
does	The	0.52
the	the	0.95
cat	cat	0.98
sit	sits	0.89

Motivation

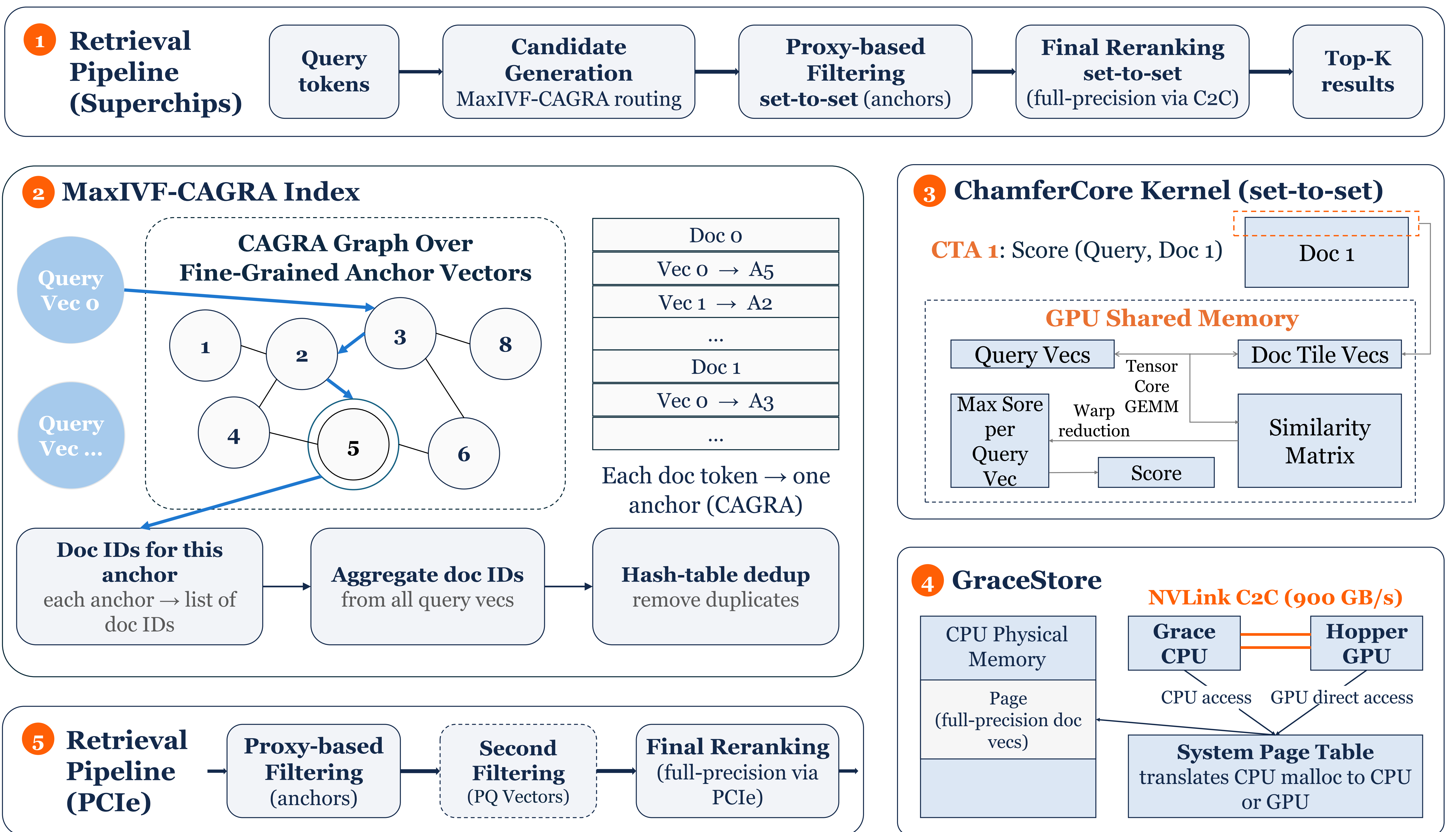
- Multi-vector gives **fine-grained** representation
- Heavy Compute and memory loads suit **GPUs**
- Storage needs a GPU + CPU **hierarchy**

Challenges

Multi-vector search **costs** far more

Single-vector dim 768, multi-vector dim 128, ~68 token/doc (MS MARCO).

VecFlow-Chamfer Architecture Overview



Evaluation Highlights

- ▲ PLAID ■ VecFlow-Chamfer
- Up to **16.9x** faster, and **+7.5 pts** recall with 98.45% R@100 in 0.97 ms vs PLAID's 91% in 16.35 ms
- MaxIVF-CAGRA: **~25%** fewer candidates, higher coverage
- Sub-ms Chamfer scoring **31K docs in 0.12 ms** (MS MARCO)

Scan to Share

Code:
github.com/Supercomputing-System-AI-Lab/VecFlow

Paper:
dl.acm.org/doi/abs/10.1145/3786706