

# Double Filtering Using Short and Long Quantized Projections

## — Memory-Efficient Online k-NN Search for Large-Scale Datasets —

(beyond PUBMED23 in the Indexing Challenge)

N. Higuchi<sup>1</sup>, Y. Imamura<sup>2</sup>, T. Shinohara<sup>3</sup> K. Hirata<sup>3</sup> T. Kuboyama<sup>4</sup>  
1. Sojo Univ. 2. Third Inc. 3. Kyushu Inst. Tech. 4. Gakushuin Univ.

### Key Features of Our Method

- Double Filtering
  - Two types of index: short and long projections
  - 1st stage by fast filtering by enumeration
- Asymmetric distance between query and projection
  - Better than Hamming distance
- Memory-efficient for large datasets

### Results for Recall@30 = 70%

				Filtering Cost (ms/q)		Latency (ms/q)			
w	dim	qbit	k <sub>1st</sub>	k <sub>2nd</sub>	1 <sup>st</sup>	2 <sup>nd</sup>	avg	std	avg <sup>i</sup>
18	192	1	180,000	600	0.03	0.49	0.64	0.09	1.05
20	192	1	150,000	400	0.06	0.48	0.66	0.10	1.12

w: sketch width  
dim, qbit: QSMAP dimension and quantization  
k<sub>1st</sub>, k<sub>2nd</sub>: #candidates of 1<sup>st</sup> and 2<sup>nd</sup> filtering  
1<sup>st</sup>, 2<sup>nd</sup>: filtering cost (ms/q)  
avg, std: latency (ms/q), avg<sup>i</sup>: reported by SISAP

Experimental Environment: Docker container with 16GB RAM and 8 CPU,  
running on Windows 11 WSL, Ryzen 9 3950X CPU, 64 GB RAM

### Filtering by Sketch Enumeration

```
// q: query, k': number of candidates
// PriorityOrder(q): sketch enumerator
// σ: sketch transformation
// σ-1(ς): { x ∈ DS | σ(x) = ς }, the set of points whose sketch is ς
FilteringBySketchEnumeration(q, k')
  C := ∅;
  for ς in PriorityOrder(q)
    C := C ∪ σ-1(ς);
    if |C| ≥ k' break;
  return C;
```

### k-NN Search by Double Filtering

Two Types of Index: Short and Long  
Short = Binary Sketches  
Long = Quantized Projections of S-Map (QSMAP)

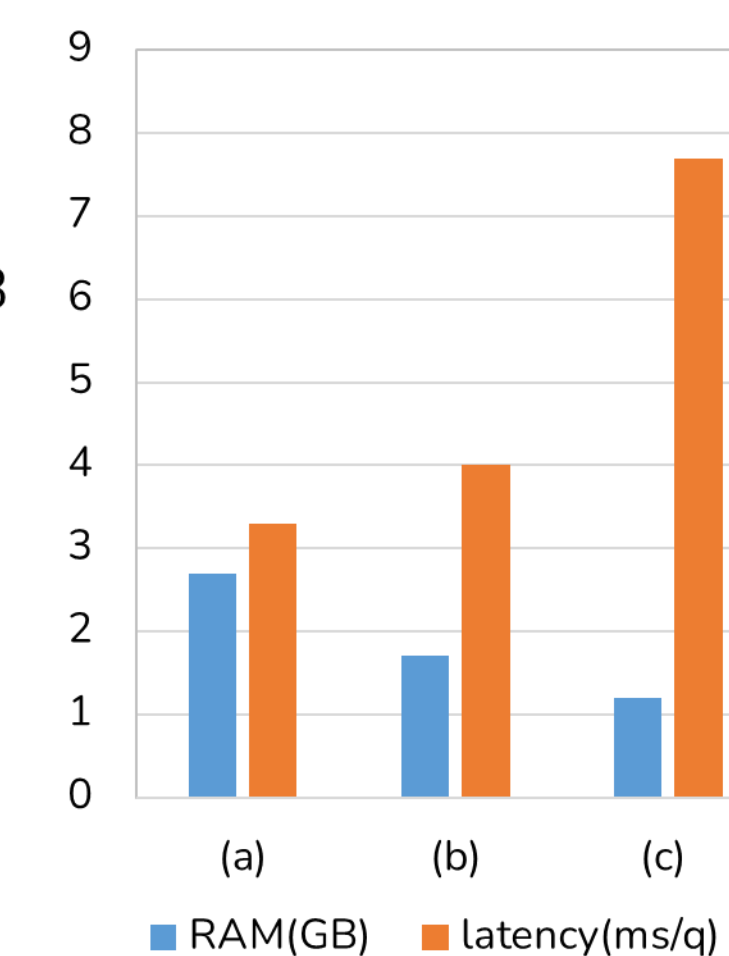
Double  
Filtering

Full Dataset ( $n > k_{1st} > k_{2nd} > k$ )  
[Coarse Filtering by Sketch Enumeration]  
> k<sub>1st</sub> Candidates  
[Fine Filtering by QSMAP]  
> k<sub>2nd</sub> Candidates  
[Reranking by Exact Distance]  
> kNN Answer

### Results on Smaller RAM

Index (Bucket + QSMAP) is very small  
Dataset on SSD  
→ Double-Filtering can run on smaller RAM < 3 GB  
Sketch width = 20-bit, Recall@30 = 70%

	dim	qbit	k <sub>1st</sub>	k <sub>2nd</sub>	latency	reranking	RAM(GB)
	384	2	300,000	30	4.0	1.5	2.7
(a)	384	2	160,000	40	3.3	1.8	2.7
(b)	384	1	260,000	60	4.0	2.7	1.7
(c)	192	1	500,000	150	7.7	6.1	1.2



### Why Fast ?

- Average bucket size  $|\sigma^{-1}(\varsigma)|$  is greater than 1
  - Typically  $\geq 10$  (e.g.,  $\sim 23$  for PUBMED23 with 20-bit sketches)
- To collect top-k' candidates:
  - Expected sketches to enumerate  $\approx k' / (\text{average bucket size})$
  - Much smaller than the total number of sketches
- Result:
  - Only a small fraction of sketches are needed
  - Explains why the 1st-stage filtering cost is negligible

### Index Construction & Total Memory Layout

#### Index Construction

##### Sketch

Width(w): 18-bit or 20-bit  
(A) Bucket Table (keyed by sketches)

##### QSMAP

192 Dimensions  
1-bit Quantization  
(B) Array Sorted in Sketch Order

##### Memory Requirement

(A):  $2^w \times \text{int} = 1 \text{ MB or } 4 \text{ MB}$   
(B):  $192 \times 1 \times 23 \text{ M (bits)} = 552 \text{ MB}$

#### Total Memory Layout

##### Index Part

= (A) + (B) < 1 GB

##### Dataset

Float (32-bit) Vectors = 30 GB  
→ Char (8-bit) Vectors = 8.8 GB  
(with negligible effect on accuracy)

##### Our Implementation

Index + Dataset < 9.8 GB < 16 GB

Fast Reranking on RAM!

### Feature Directions

- Optimal pivot selection for quantized projections
- Application to other large-scale datasets (e.g., DEEP1B)
- Similarity search in non-Euclidean or coordinate-free spaces

### Filtering by Narrow Sketch Enumeration

- Problem: Full data scan is costly ( $n = 23\text{M}$ )
- Enumeration method (example for recall@30 = 70%):
  - Required candidates  $k' = 150,000$   
→ Expected sketches to enumerate  $\approx 150,000 / 23 \approx 6,522$
  - Candidates  $\approx 0.65\%$  of full dataset → much smaller than  $n$
- Even with narrow sketches:
  - Enumeration efficiently retrieves a small set of candidates