

Steiner-Hardness: A Query Hardness Measure for Graph-Based ANN Indexes



Zeyu Wang¹ Qitong Wang² Xiaoxing Cheng³ Peng Wang^{1*} Themis Palpanas² Wei Wang¹

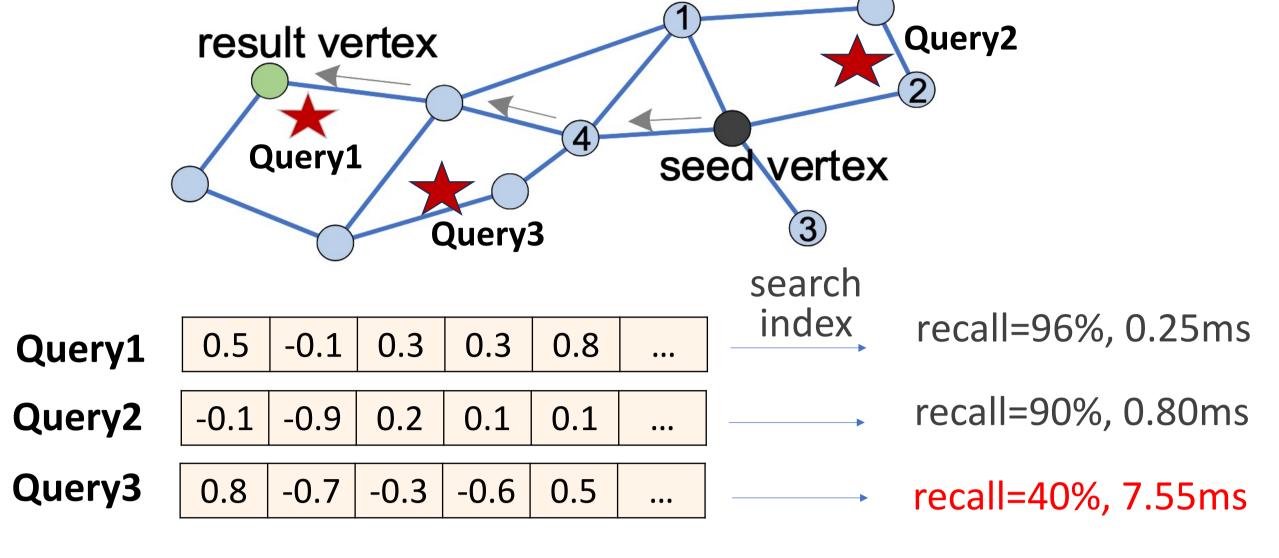
¹Fudan University, Shanghai, China ² Université Paris Cité, Paris, France ³Tongji University, Shanghai, China

<u>zeyuwang21@m.fudan.edu.cn</u> Repo: https://github.com/CaucherWang/Steiner-hardness

1 Background

ANN Query on High-D Vectors

- Find the top-k vectors with smallest distance (L2/IP/Cosine)
- Graph Indexes are the SOTA for ng-approximate result

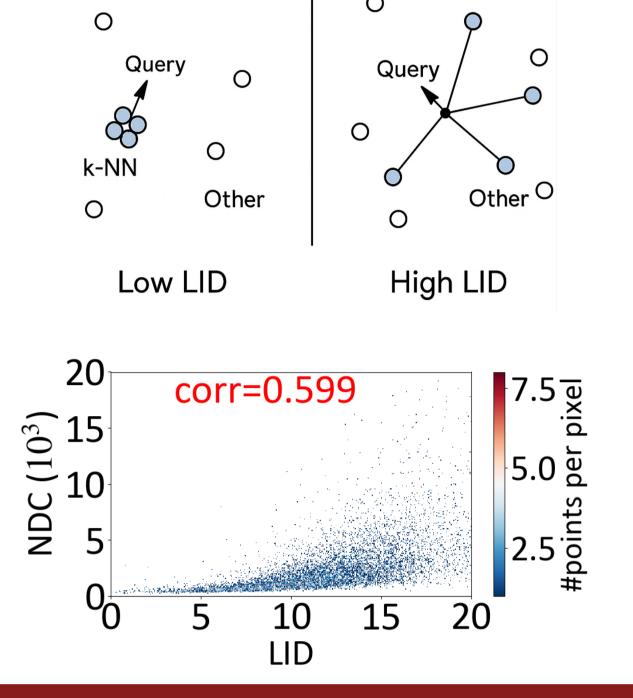


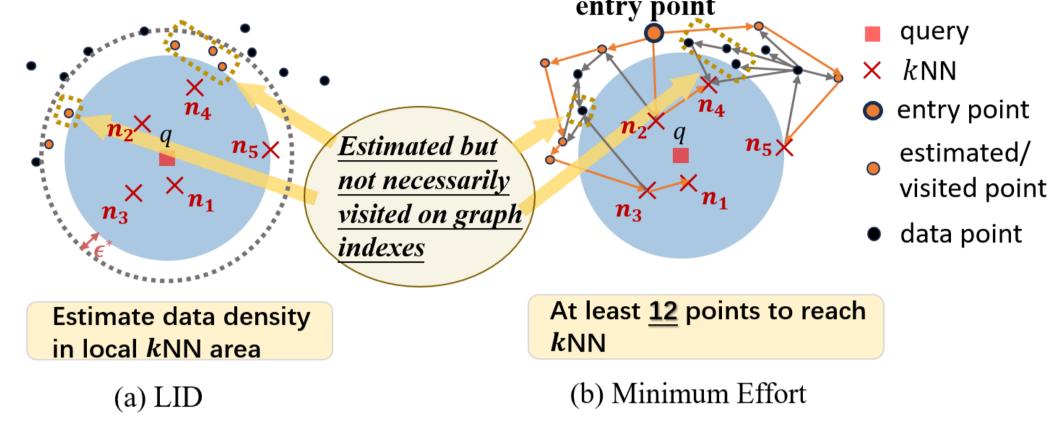
Why does the performance of the queries vary significantly?

How can we define the hardness of vector ANN queries?

2 Limitation of Distance-Based Measures

$$LID(r) = \lim_{\varepsilon \to 0} \frac{\ln(F((1+\varepsilon)r)/F(r))}{\ln((1+\varepsilon)r/r)}$$

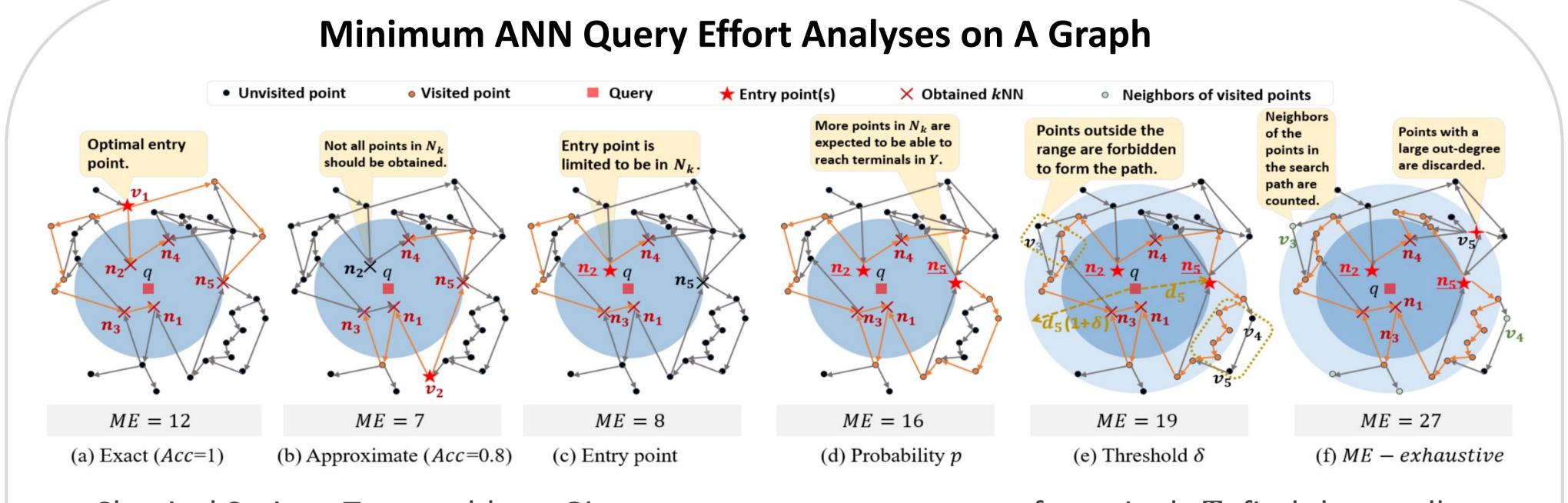




Hardness of queries is decided by graph connections to kNN.

Measures describing local data distribution (e.g. LID) are totally unaware of graph connections!

3 Steiner-Hardness

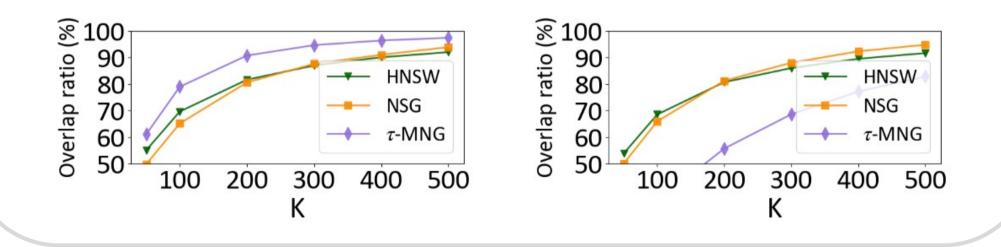


- Classical Steiner-Tree problem: Given a root vertex r, a group of terminals T, find the smallest subgraph G', s.t. for any $t \in T$, there is a path from r to t on G'.
- <u>Strict Lower Bound of ANN query</u>: The size of the optimal Steiner-tree for any possible root vertex (seed vertex), where the terminals are kNN of the query.
- Constraint on the lower bound to make it closer to the real effort: (1) the root vertex is limited to one of kNN (i.e., skip Phase-1 search), (2) not all kNNs are required to be accessed (recall), (3) only the candidates close to the query could form G', (4) the decision cost on each step.

Steiner-Hardness on All Graph Indexes

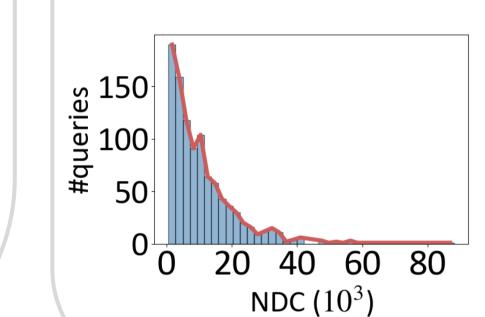
We use ME on MRNG as *Steiner*-hardness

- MRNG is pruned from KGraph
- Strict RNG pruning rule (no specific tricks)



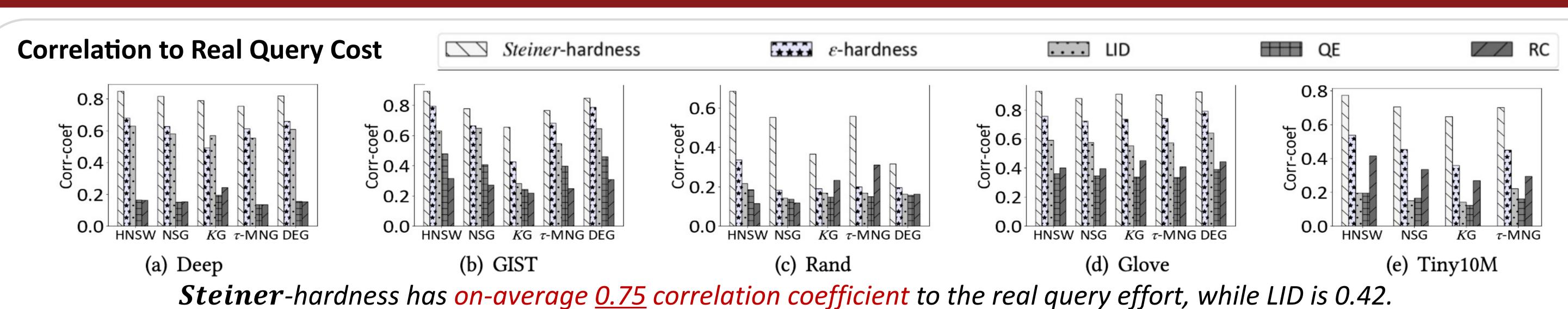
Unbiased Steiner-Workload Generation

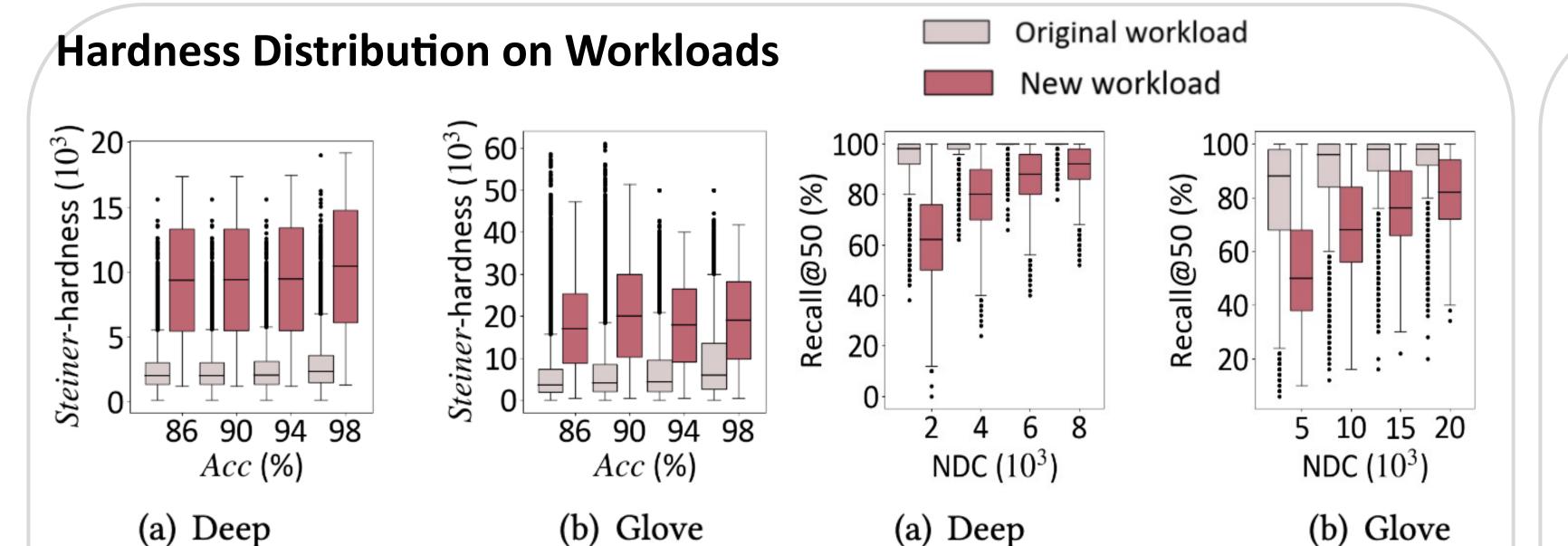
Current workloads are dominated by easy queries, leading to over-optimistic results.



- Over-sample queries from the same data distribution
 Calculate *Steiner*-hardness for new samples
- 3. Pick queries from different hardness ranges

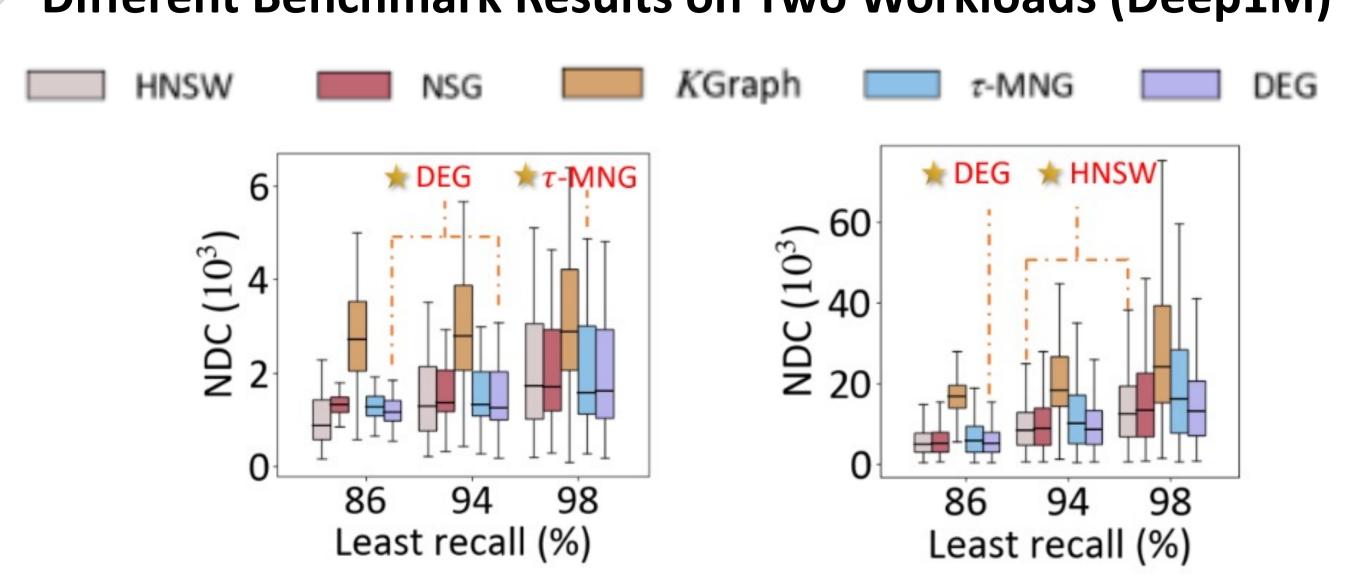
4 Experiments





More hard queries are included in the **Steiner**-workload, making an even distribution of query hardness.

Different Benchmark Results on Two Workloads (Deep1M)



On new workload: 10x worse avg. performance, larger performance variance, and $HNSW_0$ performs best.