

Graph-Based Vector Search: An Experimental Evaluation of the State-of-the-Art

Ilias Azizi^{1,2}, Karima Echihabi², Themis Palpanas³, Vassilis Christophides¹

1. ETIS, ENSEA, CNRS, France

2. College of Computing, UM6P, Morocco

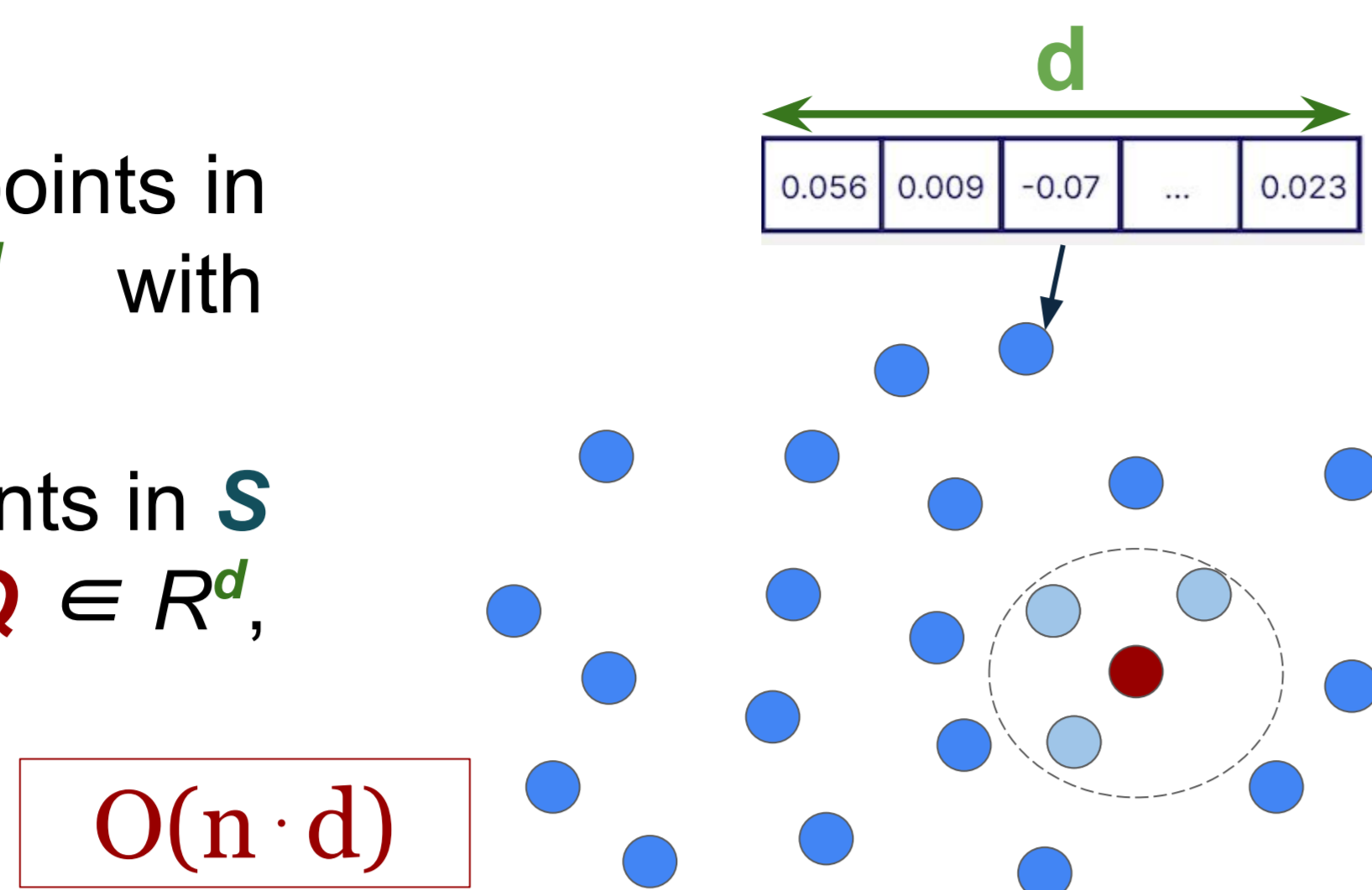
3. LIPADE, University of Paris Cité, France

Problem

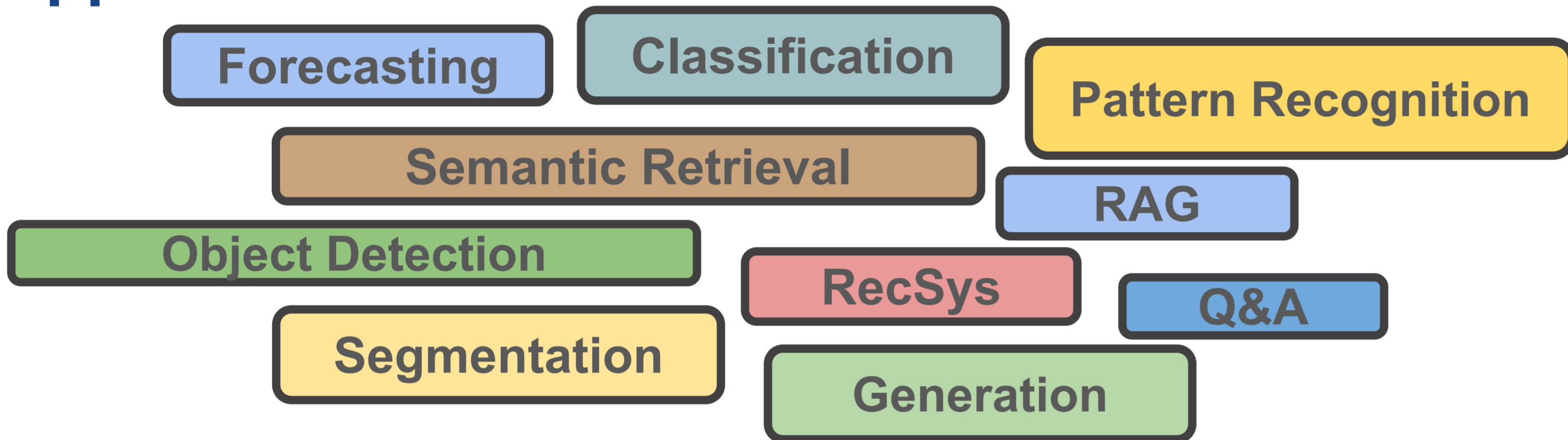
Vector Search

Given: a set S of n distinct points in d -dimensional space R^d with respect to some norm $\| \cdot \|$

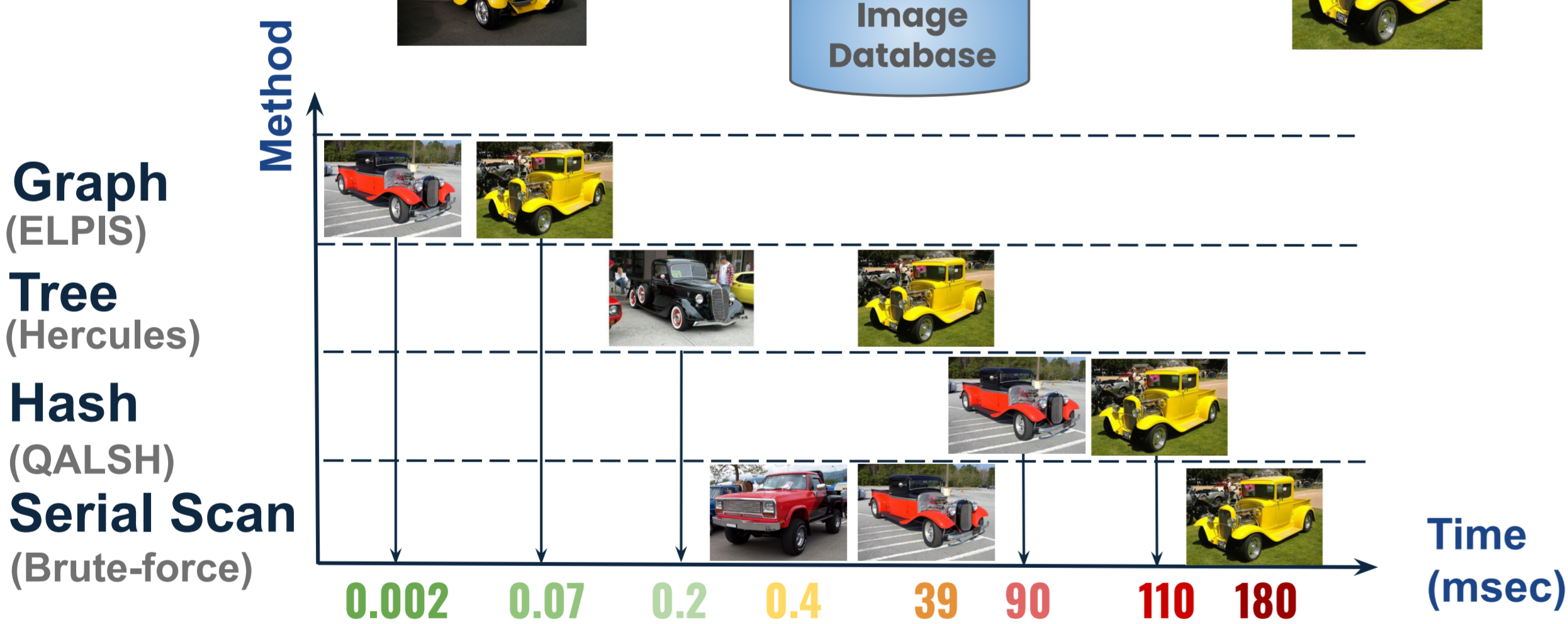
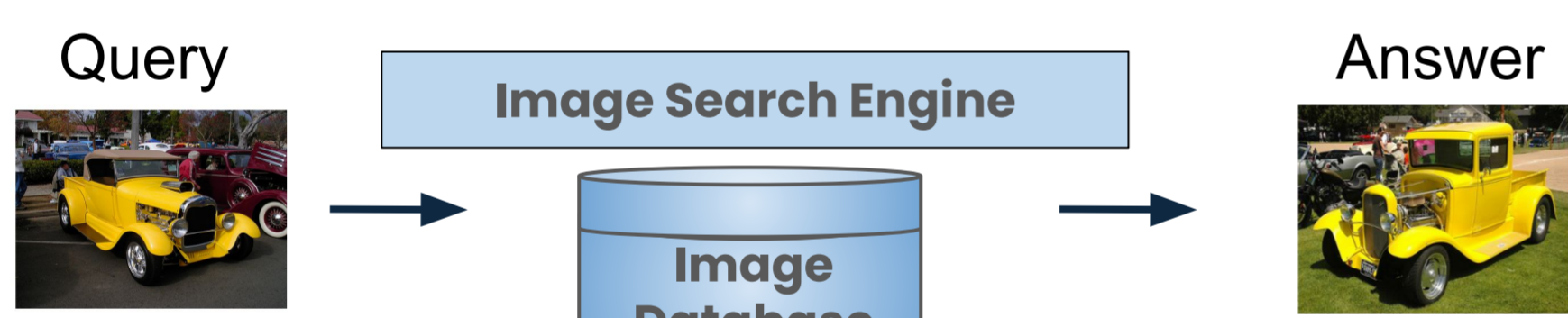
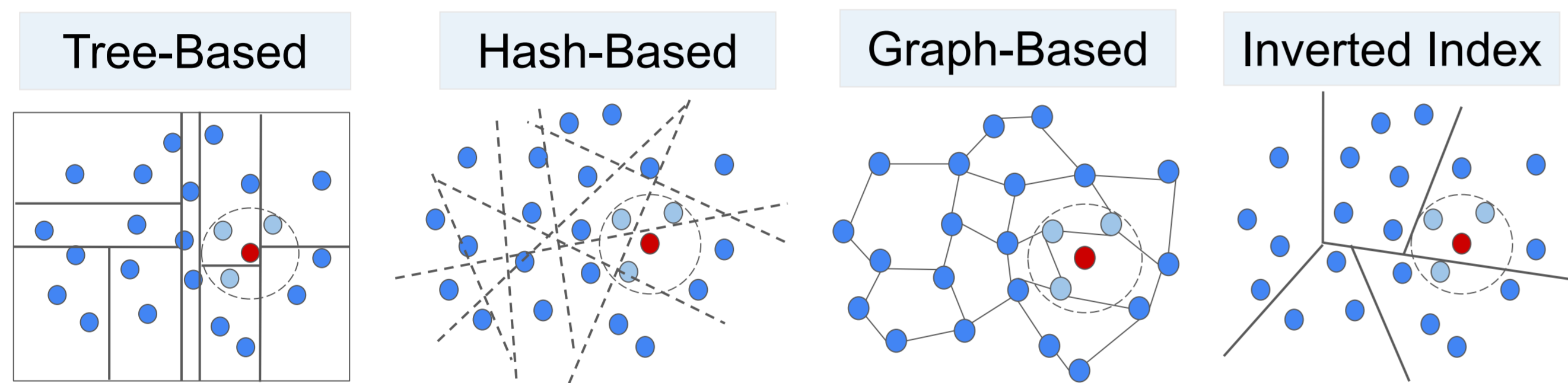
Goal: return the set of K points in S that are **closest** to a query $Q \in R^d$, under $\| \cdot \|$



Applications



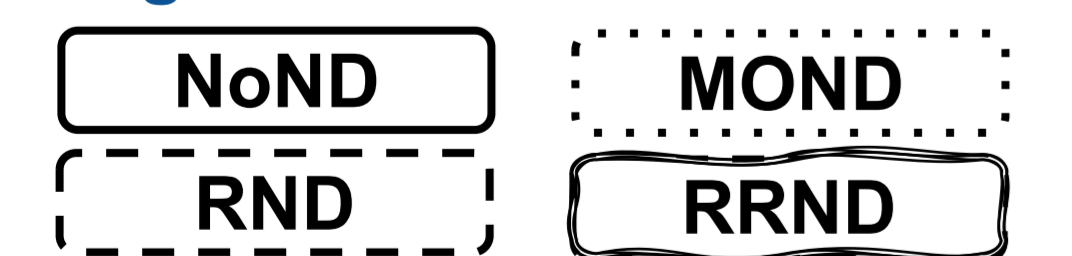
Indexing Approaches



Graph-based Vector Search

- Neighborhood Propagation based
- Incremental Insertion based
- Neighborhood Diversification based
- Divide and Conquer based

Neighborhood Diversification



Randomized Seed Selection

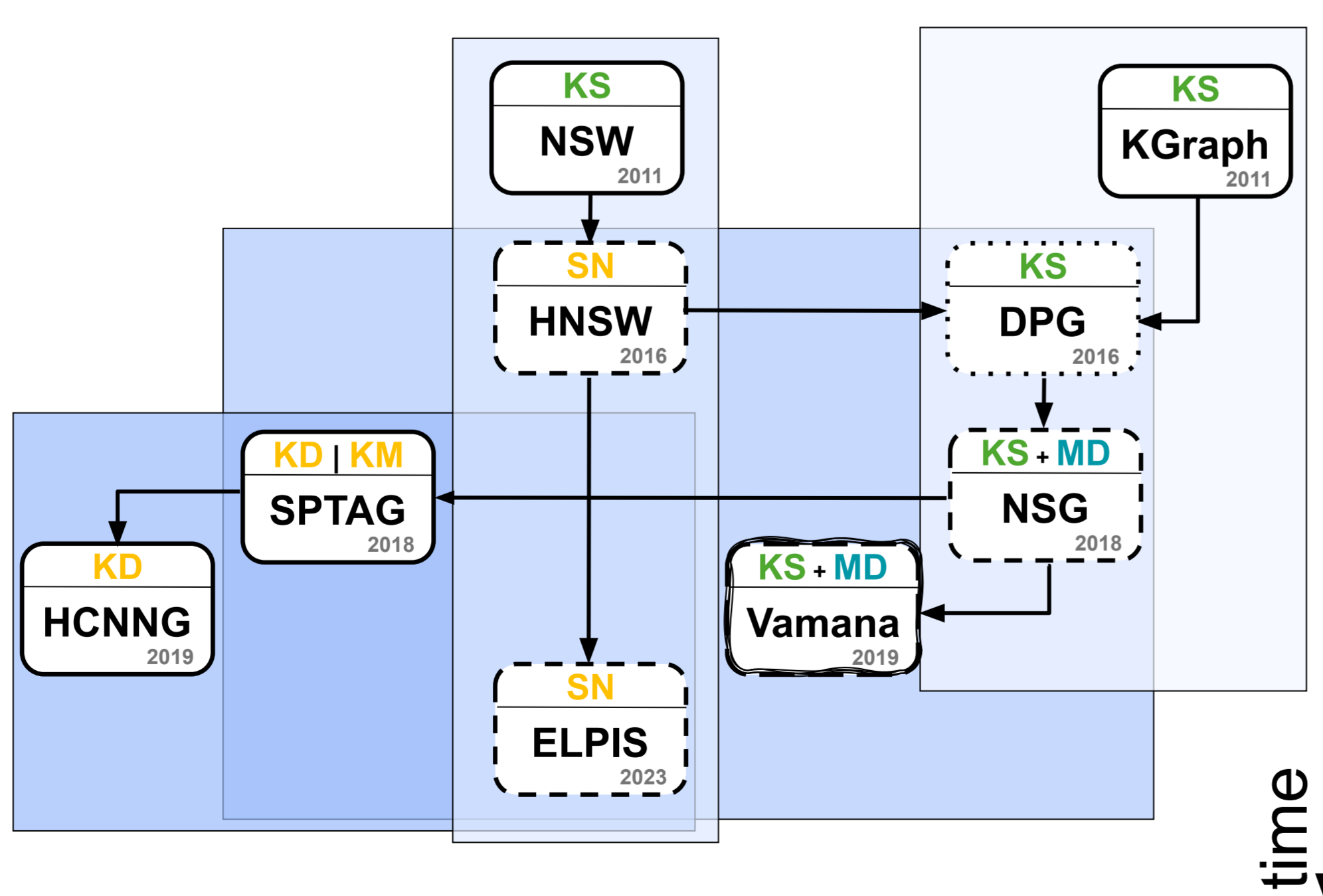
- K Random Sampling (KS)

Index-based Seed Selection

- Stacked NSW (SN)
- KD, KMean Balanced Trees

Predefined Seed Selection

- Medoid (MD)



Experimental Evaluation

Neighborhood Diversification

Smaller index size and memory footprint

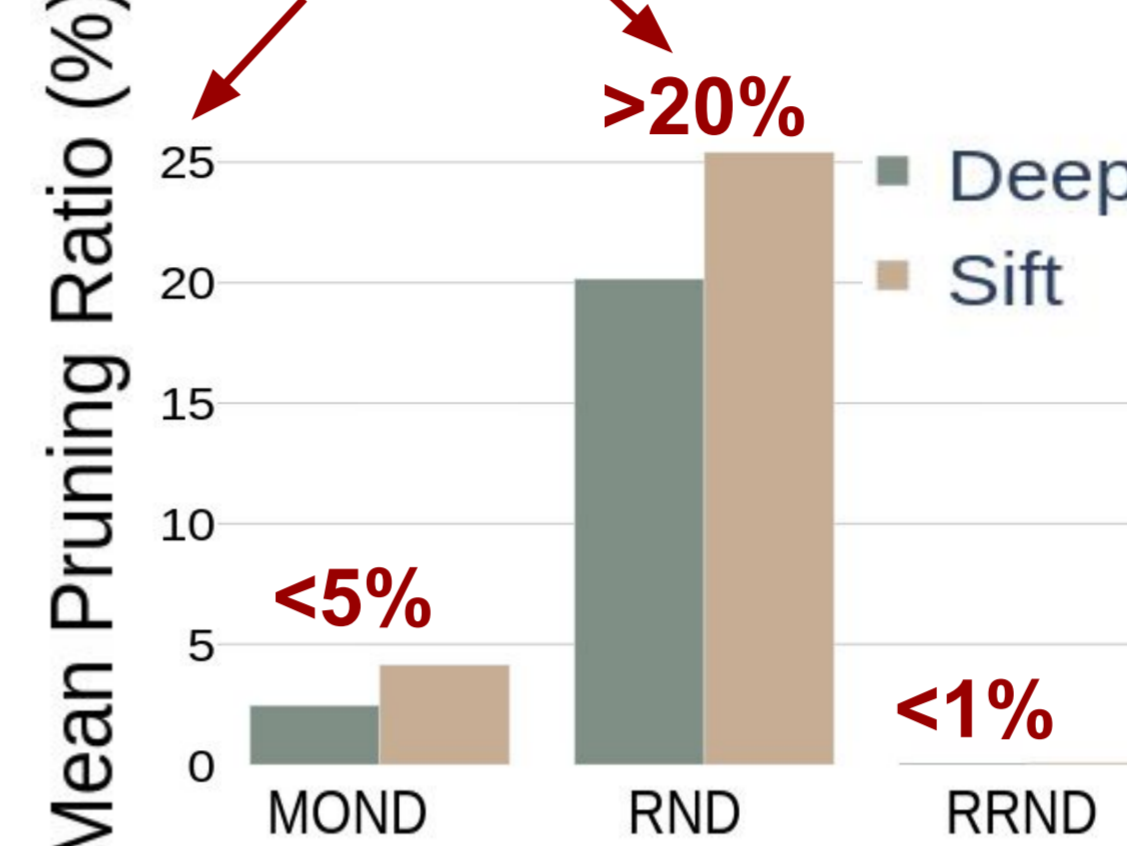


Figure 1. Average edge pruning ratio across nodes

Seed Selection

Optimizing data structures for Seed Selection enhances search efficiency on large scale datasets

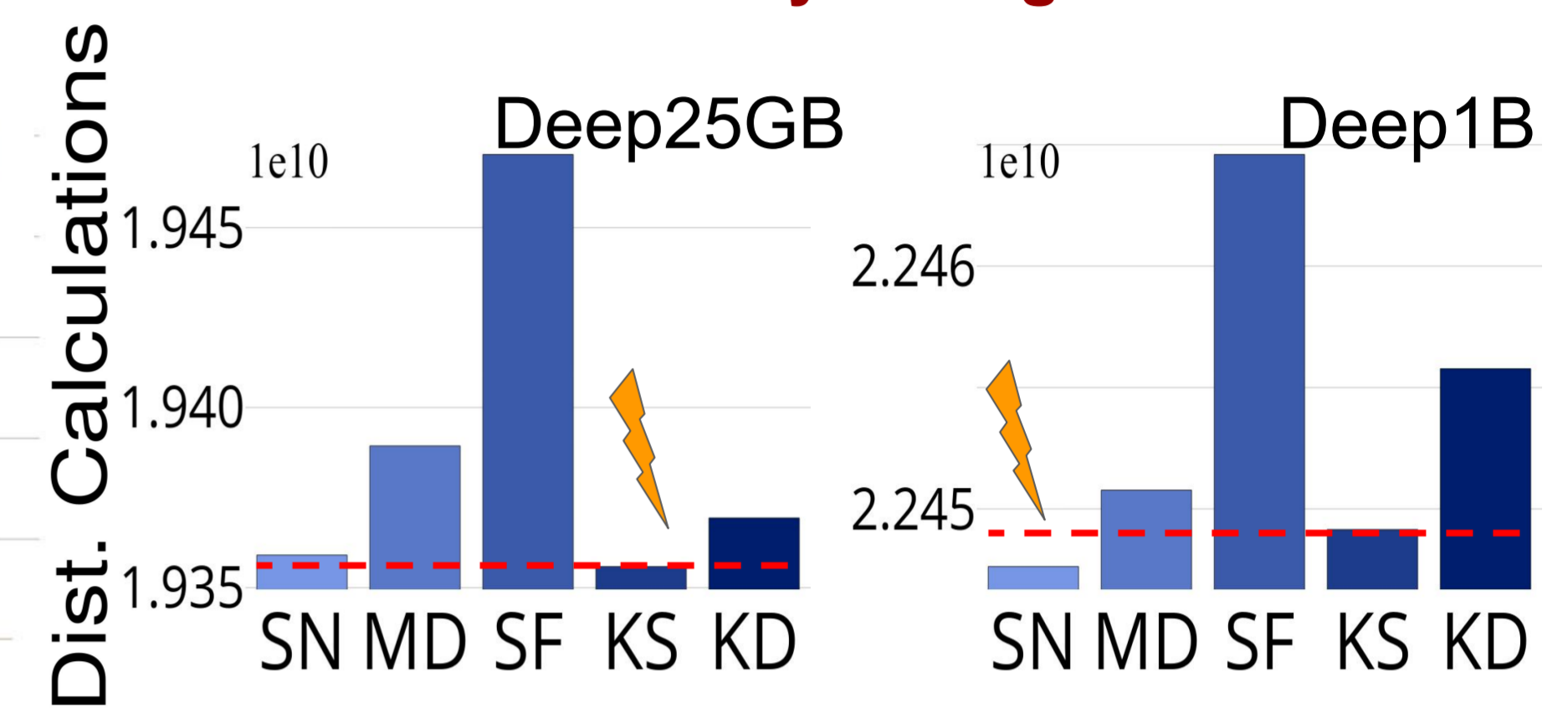


Figure 2. Impact of Seed Selection choice on Search Performance

SOTA Methods Evaluation



Incremental Insertion based graphs are the most **scalable**

Neighborhood Diversification based graphs are the most **efficient on easy workloads**

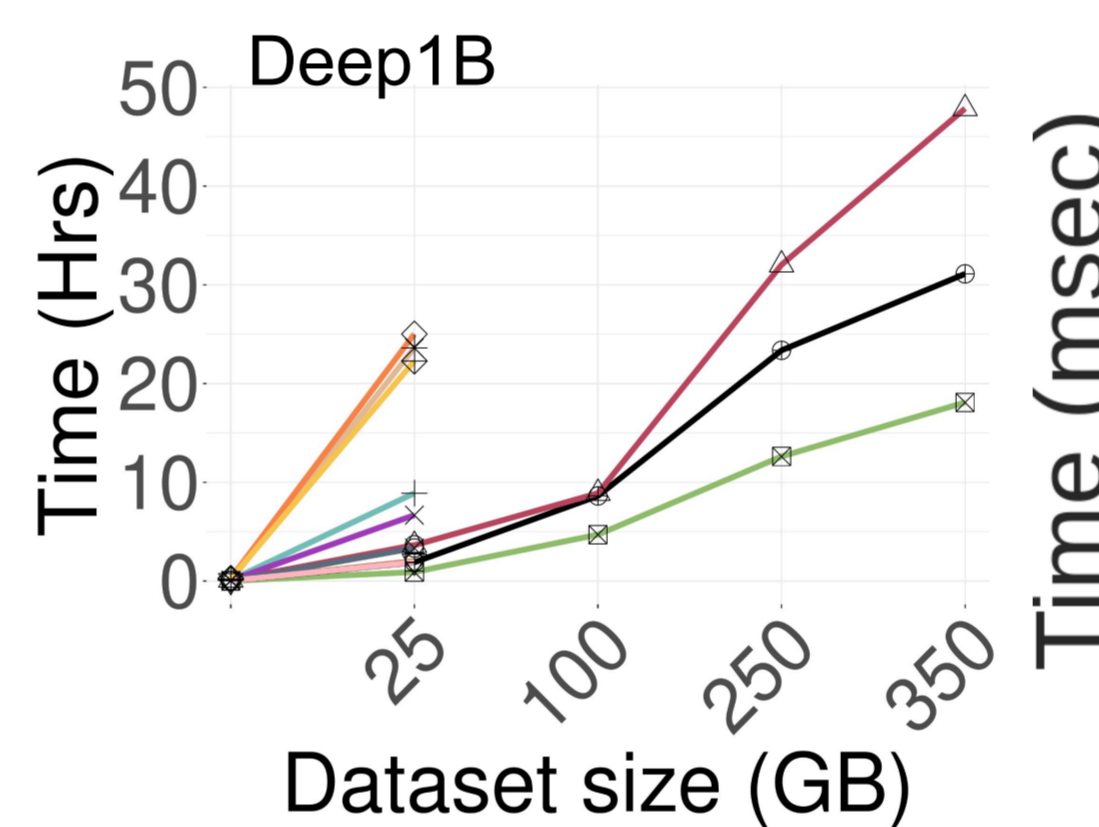


Figure 3. Indexing time

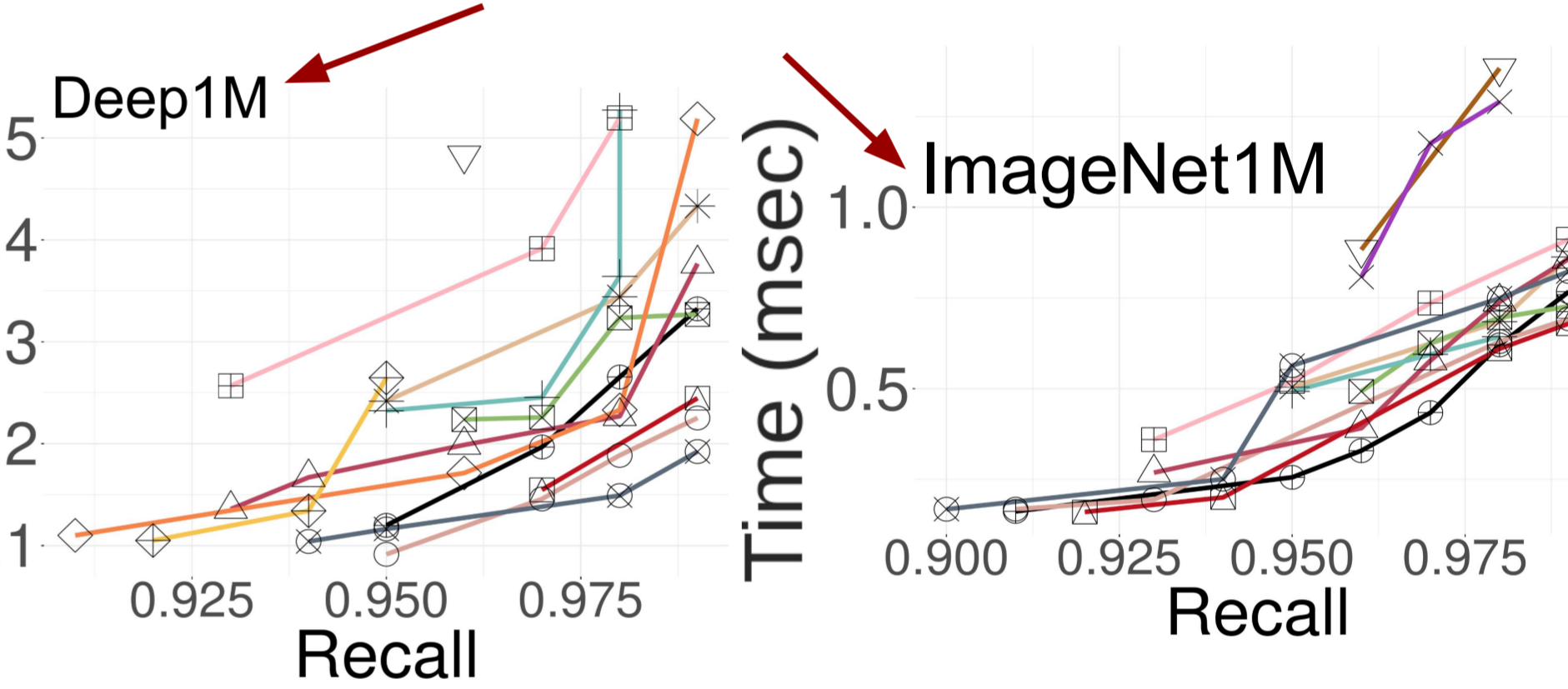


Figure 4. Search Performance on small and easy datasets

Most approaches **fail to scale** efficiently on **large datasets**

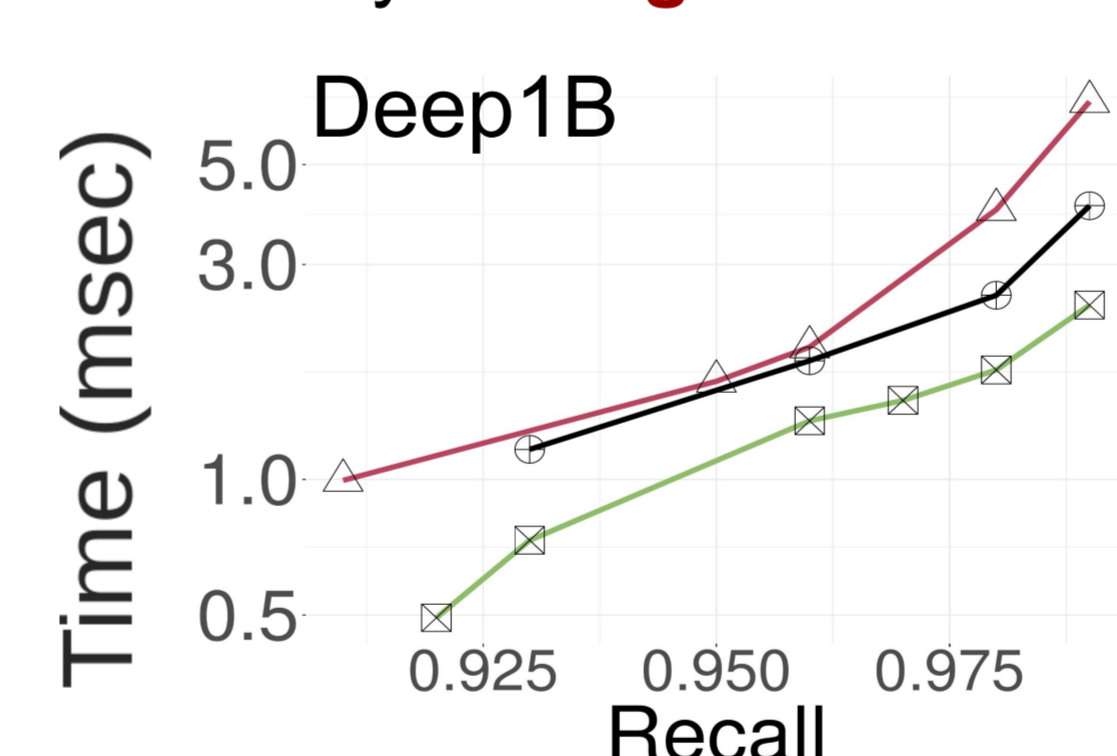


Figure 5. Search Performance on billion scale dataset

Divide-and-Conquer approaches are the most **efficient on hard workloads**

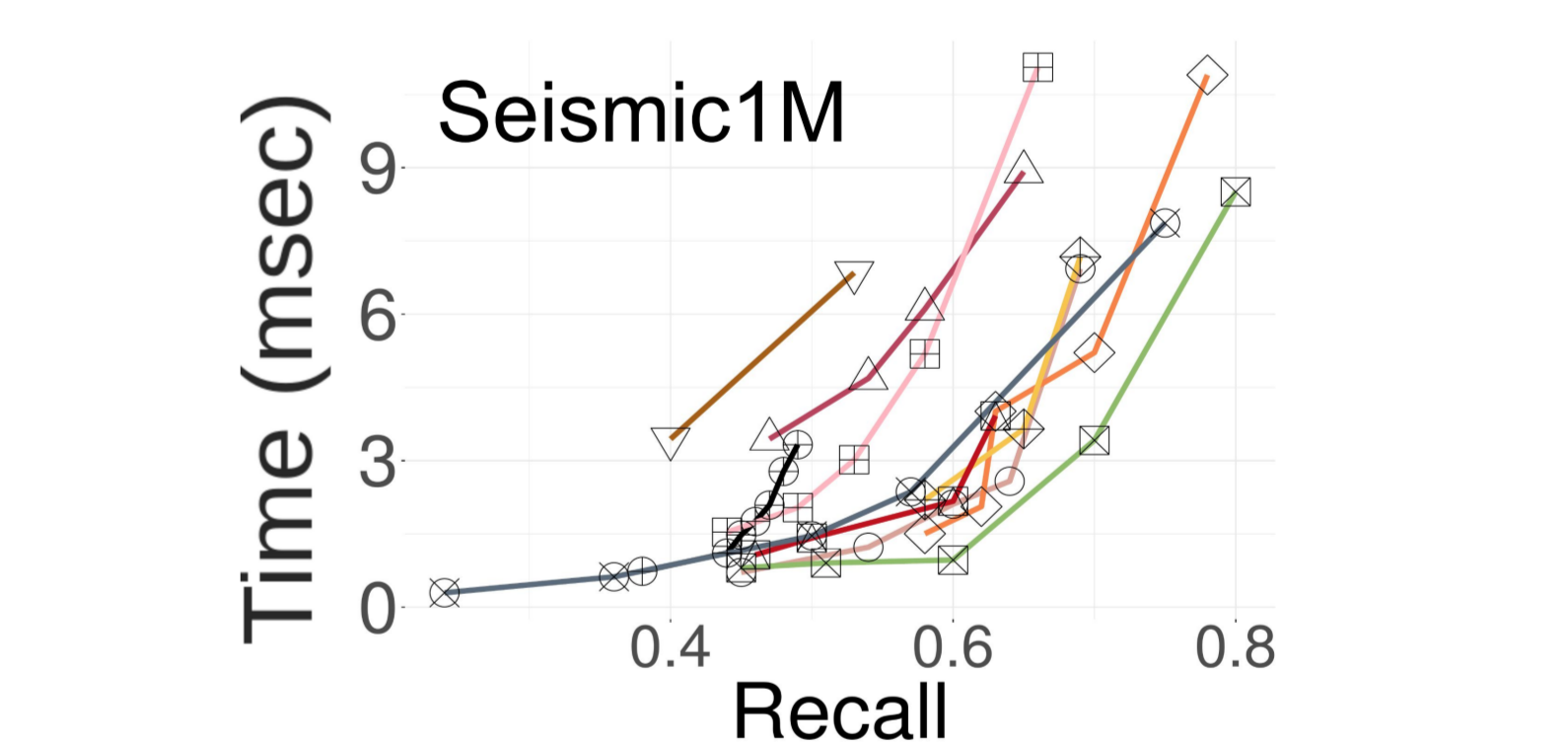


Figure 6. Search Performance on hard dataset

Recommendations

2. **Divide-and-Conquer** and **Relaxed RND** paradigms boost the performance on **hard datasets**

1. **Neighborhood Diversification** graphs achieve the overall best performance on **small and easy datasets**

3. **Incremental Insertion** graph construction offers superior scalability

