



Karima Echihabi¹ (3rd year Ph.D Student, Supervised by Themis Palpanas² and Houda Benbrahim¹) ² Université de Paris & French University Institute (IUF), France ¹ Mohammed V University, Morocco

1. Problem and Motivation > Similarity search is core operation for many critical data processing tasks data cleaning, data integration and data analytics Similarity search finds objects in a collection close to a given query according to some definition of sameness • Objects can be data series, images, graphs, deep network embeddings, etc. Similarity search algorithms can either return exact or approximate answers • Exact methods return correct results but are expensive Approximate methods sacrifice accuracy to achieve better efficiency > 1) no exact method is overall winner > 2) no approximate method offers query

- accuracy guarantees, or scales to TB of disk-based data (with low footprint)
- > We conducted two most extensive experimental evaluations in this area Lernaean Hydra papers: exact (PVLDB' and approximate techqniques (PVLDB'20 considered techniques from both data se Ο Proposed extensions to existing indexe Outperform LSH, kNNGs and inverted in **Developed Hercules:** novel versatile sin Builds upon insights gained from Lernae designs of existing approaches, their strengths and weaknesses Supports all flavors of exact and approximate search Ο Supports progressive query answering Ο

3. Approach and Uniqueness

High-Dimensional Vector Similarity Search From Time Series to Deep Network Embeddings

ßs	HERCULES: vers	atil	es
	Supports exact / a	appi	(0)
	> best in accuracy i		ne .
19)	outperforms data series ind inverted indexes, kNNGs or		
20) eries	and high-d vector communities		
es to support δ-ε-approximate search ndexes under many scenarios			
nilarity search method ean Hydra papers about intricate			

2. Background and Related Work

Similarity search has been extensively studied in past 25 years



A research direction that is more promising than the two current trends

- LSH-based algorithms, which provide guarantees, but are slow
- kNNGs and inverted indexes, which are fast but provide no guarantees

Hercules: novel similarity search method for exact/approximate search **Current Prototype**

- Exact search: Best overall. At least 2x faster than best index, 20x-100x faster than Optimized Sequential Scan (OSS) on balanced workloads, only index that outperforms OSS on hard workloads
- Approximate search: Best overall on disk (time/accuracy) with low footprint

> A public archive of code/datasets/results, a stepping stone for a benchmark

4. Results and Contributions



Université de Paris

