

Vector-Database

Overview

- Database of high-dimensional vectors
- Vectors generated by transformation or embedding function
- Useful for natural language processing, computer vision, recommendation system and more

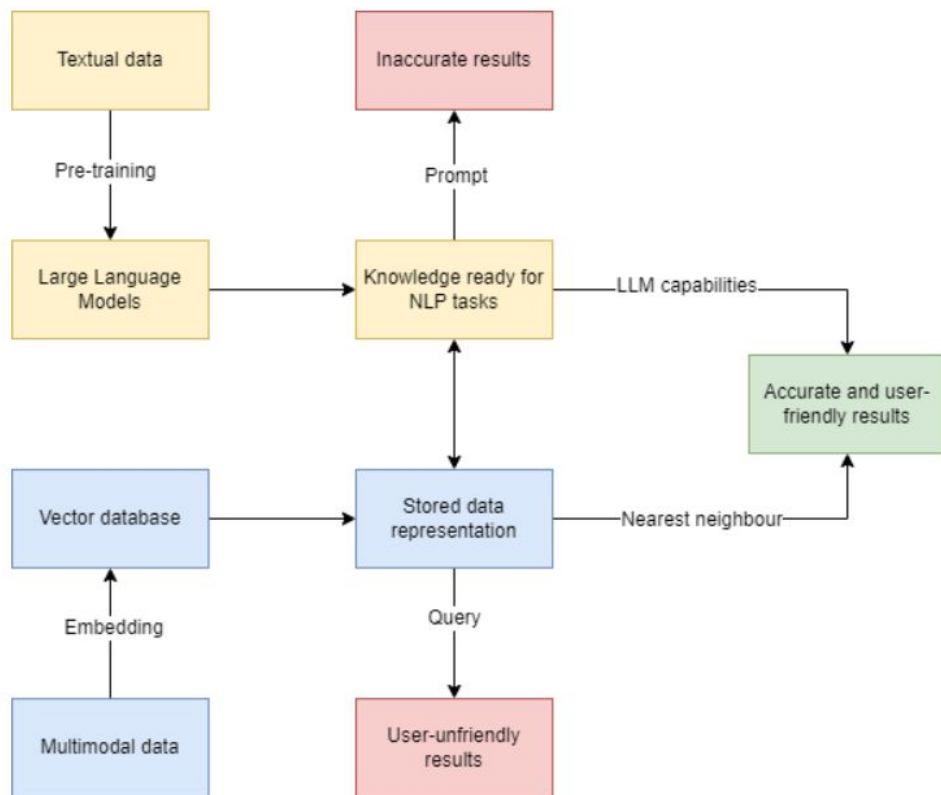
Advantages

- Fast and accurate similarity search and retrieval
- Support for complex and unstructured data
- Scalable and performant
 - Sharding, partitioning, caching, replication

Challenges

- Complex searching
- Different Vector Data Types
 - Sparse, different dimensions
- Easy API and connectors for frameworks to use
 - TensorFlow, PyTorch, Scikit-learn

LLMs



LLMs \longleftrightarrow Vector Databases

- Reduce hallucination
- Real-time knowledge
- Model compression
- Text generation
- Text augmentation
- Easy to update
- Can provide sources
- Reduce time and cost

Vector-Search: Brute-Force

- Calculate distance from query vector to all vectors in the database
- Advantages
 - Guarantees the exact nearest neighbors.
 - Simple to implement.
- Disadvantages
 - Computationally expensive for large datasets
 - Does not scale well with the number of vectors or the dimensionality of the vectors
- Examples: Euclidean/cosine similarity

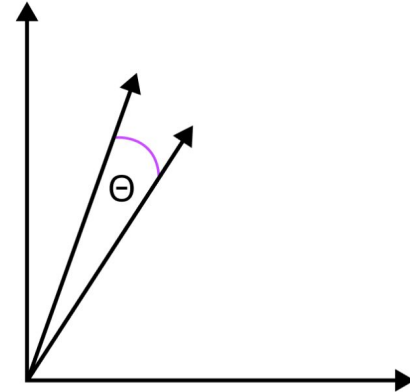
Vector-Search: Brute-Force



```
def cosine_similarity(vec1, vec2):  
    vec1 = np.array(vec1)  
    vec2 = np.array(vec2)  
    dot_product = np.dot(vec1, vec2)  
    norm_vec1 = np.linalg.norm(vec1)  
    norm_vec2 = np.linalg.norm(vec2)  
    return dot_product / (norm_vec1 * norm_vec2)
```

Cosine Similarity

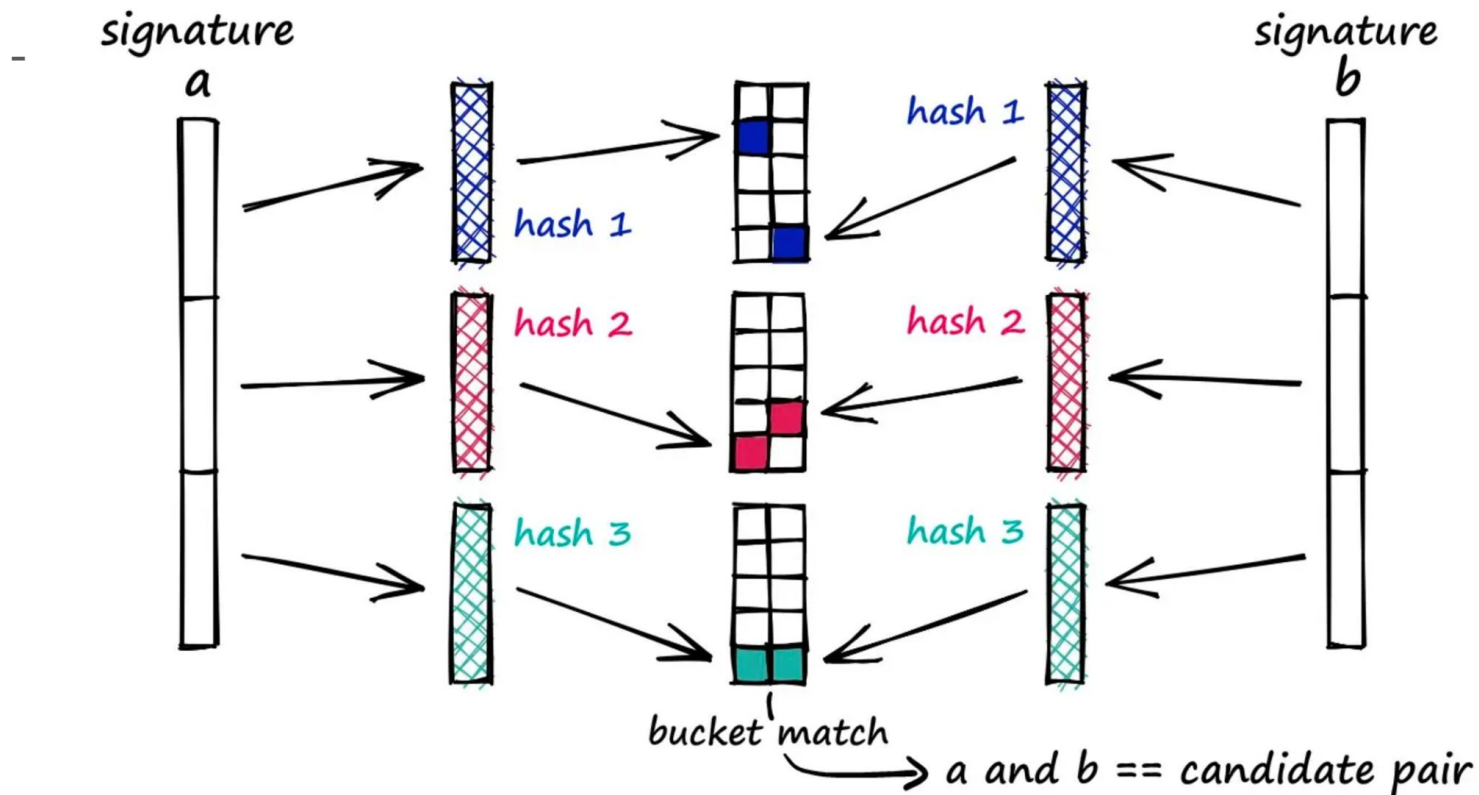
$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



Vector-Search: ANN Search

- Approximate Nearest Neighbor (ANN) Search with Locality-Sensitive Hashing (LSH)
 - Hashing function that maximizes collision (for similar inputs)
- Advantages
 - Faster search for high-dimensional data
 - Tunable trade-offs between speed and accuracy
- Disadvantages
 - Provides approximate, not exact, results
 - Hash function selection and tuning can be complex

Vector-Search: ANN Search



Vector-Search: ANN Search



```
class LSH:
    def __init__(self, num_vectors, num_hash_tables, dimension):
        self.num_vectors = num_vectors
        self.num_hash_tables = num_hash_tables
        self.dimension = dimension
        self.hash_tables = [defaultdict(list) for _ in range(num_hash_tables)]
        self.random_vectors = [np.random.randn(dimension) for _ in range(num_hash_tables)]

    def _hash(self, vector, random_vector):
        return 1 if np.dot(vector, random_vector) > 0 else 0

    def _get_hash(self, vector):
        return tuple([self._hash(vector, rv) for rv in self.random_vectors])

    def add(self, vector, idx):
        for table_idx, hash_table in enumerate(self.hash_tables):
            hash_key = self._hash(vector, self.random_vectors[table_idx])
            hash_table[hash_key].append(idx)

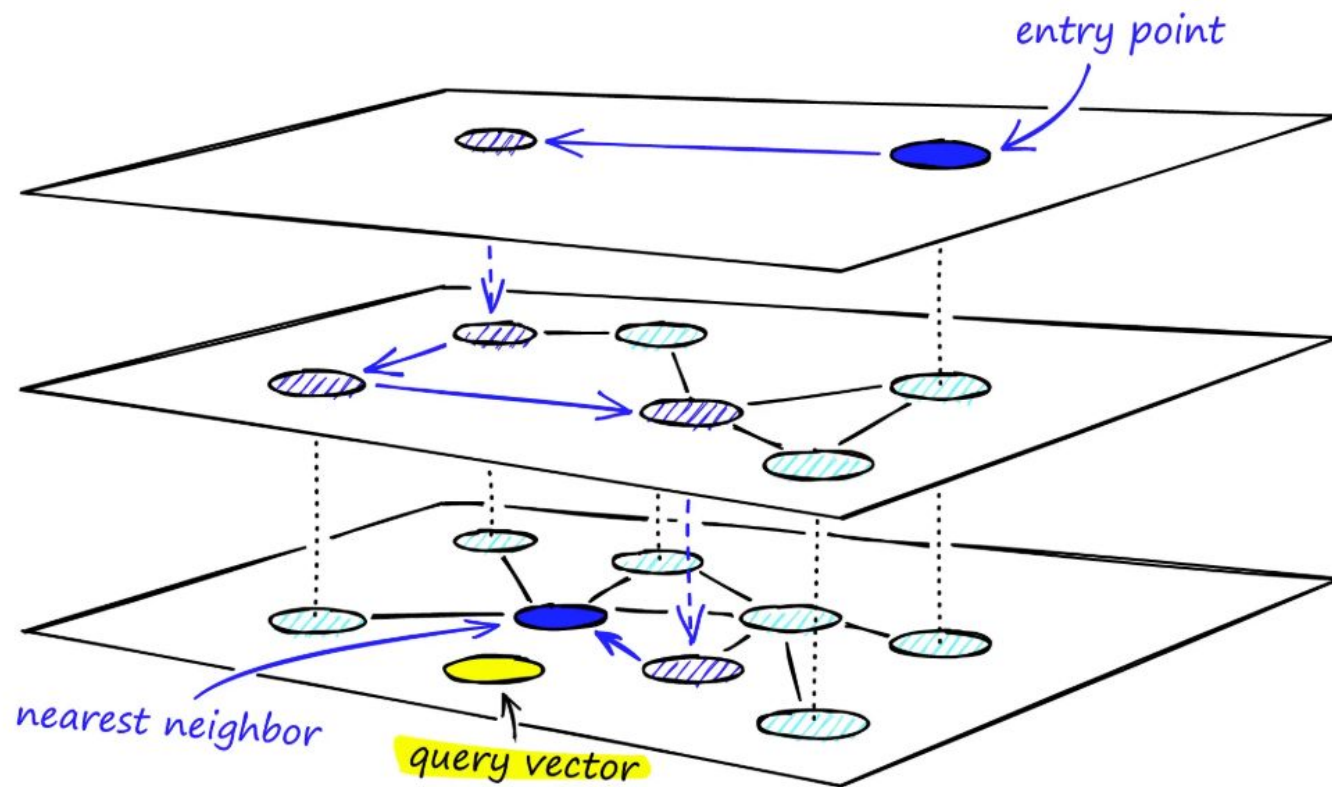
    def query(self, vector, num_candidates=10):
        candidates = set()
        for table_idx, hash_table in enumerate(self.hash_tables):
            hash_key = self._hash(vector, self.random_vectors[table_idx])
            if hash_key in hash_table:
                candidates.update(hash_table[hash_key])

        return list(candidates)[:num_candidates]
```

Vector-Search: HNSW

- Hierarchical Navigable Small World
 - Graph structure where nodes represent vectors, and edges connect vectors that are close to each other
 - During search, it starts from a randomly selected node and navigates towards the query vector by exploring neighboring nodes iteratively
- Advantages
 - Very efficient for high-dimensional vector search
 - Balances speed and accuracy
- Disadvantages
 - High memory overhead due to graph structure
 - Initial graph construction is computationally expensive

Vector-Search: HNSW



Chroma Vector DB

- Create client:

```
1 import chromadb
2 chroma_client = chromadb.Client()
```

- [Alternative] Create persistent client (creates sqlite3-file at location):

```
client = chromadb.PersistentClient(path="/path/to/save/to")
```

- Possible to run in client-server mode

```
1 import chromadb
2 chroma_client = chromadb.HttpClient(host='localhost', port=8000)
```

Chroma Vector DB

- Create collection (where the embeddings, docs, and metadata are stored)

```
1 collection = chroma_client.create_collection(name="my_collection")
```

- Simply adding documents to the collection

```
1 collection.add(  
2     documents=[  
3         "This is a document about pineapple",  
4         "This is a document about oranges"  
5     ],  
6     ids=["id1", "id2"]  
7 )
```

Chroma Vector DB

- Querying

```
1 results = collection.query(  
2     query_texts=["This is a query document about hawaii"],  
3     n_results=2 # how many results to return  
4 )  
5 print(results)
```

- Result:

```
1 {  
2     'documents': [[  
3         'This is a document about pineapple',  
4         'This is a document about oranges'  
5     ]],  
6     'ids': [['id1', 'id2']],  
7     'distances': [[1.0404009819030762, 1.243080496788025]],  
8     'uris': None,  
9     'data': None,  
10    'metadatas': [[None, None]],  
11    'embeddings': None,  
12 }
```