# Tutorial 7: Vector Databases

## Question Set:

**Q1.** What is a vector database?

**ANSWER:**

A vector database is a specialized type of database designed to efficiently store, manage, and query high-dimensional vector data, which is commonly used in machine learning, natural language processing, and computer vision tasks.

Unlike traditional databases that are utilized for data query, vector databases are optimized for operations like similarity search, where the goal is to find the nearest vectors to a given query vector based on metrics such as Euclidean distance or cosine similarity. These databases leverage advanced indexing techniques, such as product quantization, inverted files, and hierarchical graph structures, to handle large-scale datasets and ensure fast retrieval times. As a crucial component in AI-driven applications, vector databases enable the scalable and efficient handling of complex data, powering everything from recommendation systems to semantic search engines.

Foundations of vector databases include:

- Vector Embeddings:
  - Dense representation of data items (like text, images, or sounds) in a high-dimensional space.
  - Each dimension represents a feature of the data, and the distance between vectors indicates their similarity.
  - Also known as dense representations, feature vectors, latent vectors, embedding vectors, etc.
- Similarity Search:
  - Search based on the similarity of data points, rather than exact matches
  - Common metrics used to determine similarity include Euclidean distance, cosine similarity, and inner product.
- Indexing:
  - Critical for vector databases to speed up similarity searches among high-dimensional data
  - Techniques such as k-d trees, hierarchical navigable small world graphs (HNSW), and inverted files are used to enable efficient nearest neighbor searches.

**Q2.** Why do we need vector databases?

**ANSWER:**

Traditional databases are not optimized for operations that involve complex mathematical computations, such as similarity search or nearest neighbor search. Vector databases, on the other hand, are built to perform these tasks efficiently, allowing for fast retrieval of data based on vector similarity rather than exact matches.

Vector databases are specifically designed to handle and efficiently manage high-dimensional vector data, which is essential for many modern AI and machine learning applications.

In applications like recommendation systems, image and video search, natural language processing, and anomaly detection, data is often represented as vectors in high-dimensional spaces. These vectors encapsulate complex relationships and patterns within the data, making it possible to compare and find similar items even in large datasets. Vector databases use advanced indexing methods and storage optimizations to ensure that these operations can be performed at scale, with low latency, and with high accuracy, making them indispensable for powering AI-driven systems in real-time environments.

**Q3.** Discuss Indexes in Faiss.

- Flat Index
- Hierarchical Navigable Small Worlds Index
- Inverted File Product Quantization Index
- Locality Sensitive Hashing Index

**ANSWER:**

Faiss is a powerful library for similarity search and clustering of dense vectors, particularly designed for efficient nearest neighbor search on high-dimensional vectors. Different indexing methods are available in Faiss, each tailored to various performance and memory needs.

A **vector Index** is a data structure used in computer science and information retrieval to efficiently store and retrieve high-dimensional vector data, enabling fast similarity searches and nearest neighbor queries.

**1. Flat Index:**

The Flat Index is the simplest Index in Faiss, where all vectors are stored in memory and compared in a **brute-force** manner. During a query, the Index computes the distance from the query vector to all vectors stored in the dataset and returns the closest matches.

- Usage:
  - Best for small datasets where computational resources and time are not major concerns, as the method involves a complete search through all vectors.
  - Ensures **exact nearest neighbor search**, making it highly accurate but potentially slow for large datasets.
  - Deliver the best accuracy of all Indexing methods.

- Advantages:
  - High accuracy, as every vector in the dataset is compared.
  - Simple and straightforward to implement.

- Disadvantages:
  - Computationally expensive for large datasets, as it requires computing the distance for all vectors.
  - Requires significant memory to store all vectors in RAM.

## 2. Hierarchical Navigable Small Worlds Index (HNSW)

HNSW is a graph-based Index that constructs a graph of vectors, where each vector is connected to its nearest neighbors. Nodes and edges in HNSW represent the vector embeddings and relationships among these embeddings, respectively. The search process starts from an entry point in the graph and navigates through the graph to find the closest neighbors of the query vector.

- Usage:
  - Designed for **approximate nearest neighbor search** and is suitable for very large datasets where exact search may not be feasible.
  - Frequently used in scenarios where low-latency and real-time performance are critical, such as in recommendation systems.

- Advantages:
  - Significantly faster than brute-force search, especially for large datasets.
  - Provides a good trade-off between search accuracy and performance.
  - Scalable to large datasets, as it navigates only through relevant parts of the graph.

- Disadvantages:
  - Slightly more complex to implement compared to a Flat Index.
  - May not guarantee accurate nearest neighbor search results due to the approximate nature of the algorithm.

## 3. Inverted File Product Quantization Index (IVFPQ)

The Inverted File with Product Quantization Index is a two-stage Index that first partitions the dataset into multiple clusters (inverted file system) and then applies product quantization within each cluster to reduce memory usage and search time. The vectors are assigned to a cluster, and only a subset of clusters is searched based on the query vector, reducing the number of distance computations. Product quantization compresses the vectors within each cluster, allowing for efficient use of memory while maintaining reasonable search accuracy.

- Usage:
  - Ideal for **large-scale approximate nearest neighbor search** where memory efficiency and search speed are critical.
  - Commonly used in scenarios like large image retrieval systems, where exact results are not always necessary.

- Advantages:
  - Scalable to very large datasets, since the search space is reduced by only considering relevant clusters.
  - Product quantization reduces memory usage while retaining much of the search accuracy.
  - Efficient and fast, especially for very large datasets.

- Disadvantages:
  - Search is approximate, so it may not always return the exact nearest neighbors.
  - More complex to implement and tune compared to simpler index types like Flat Index.