



# Lecture Notes: Vectors & Vector Databases

## 1. The Problem: Why do we need Vectors?

To understand why vectors exist, we must look at how recommendation systems (Amazon, Netflix) evolved.

- **Solution 1: Manual Categorization (Arrays)**
  - **Method:** Engineers manually create arrays/lists of similar items (e.g., Gym items in one list, fruits in another).
  - **Failure:** Not scalable. If Amazon has 1 million products, you can't manually map them. It also fails to find "hidden" relationships (e.g., people buying diapers and beer together on Fridays). [\[04:53\]](#)
- **Solution 2: Graph-Based (Co-occurrence)**
  - **Method:** Create nodes for products and increase the "weight" of the edge between two products whenever they are bought together.
  - **Failure:** "Cold Start Problem." If a new product is launched, it has no purchase history, so the system can't recommend it. It also consumes massive memory ( $\$N \times N$  matrix). [\[11:18\]](#)

---

## 2. What are Vector Embeddings?

A **Vector** is essentially a list of numbers that represents the "features" or "meaning" of an item (movie, product, word) in a multi-dimensional space. [\[42:51\]](#)

- **The Dimensional Model:**
  - Imagine a 2D graph where the X-axis is "Action" and the Y-axis is "Comedy."
  - A movie like *Dhamal* would be high on Comedy and low on Action.
  - A movie like *John Wick* would be high on Action and low on Comedy.
- **Scaling to High Dimensions:** In production (like Llama models), we don't just use 2 dimensions. We use **4,096+ dimensions** covering nuances like "Romance," "Emotional Depth," "Realism," etc. [\[46:00\]](#)
- **The Goal:** By converting items into these numbers, similar items end up physically close to each other in the mathematical "Vector Space."

---

## 3. Measuring "Closeness" (Search & Retrieval)

Once everything is a vector, how does the system find similar items? It uses distance math. [\[52:13\]](#)

1. **Euclidean Distance:** Measures the straight-line distance between two points. It cares about the magnitude (how "strong" the features are).

- 
2. **Cosine Similarity** (★★★ Preferred for LLMs): \* Measures the **angle** between two vectors.
- It doesn't care how long the vectors are; it only cares if they are pointing in the same direction (meaning they have the same context/meaning).
  - **Range:** -1 to 1. (1 = identical meaning, -1 = opposite). [\[54:39\]](#)

---

## 4. Vector Databases: The Purpose

Traditional databases (SQL/NoSQL) are built for exact matches (e.g., "Find user with ID 101"). They are terrible at "Similarity Search."

- **Vector Databases** (like Pinecone, Milvus, Qdrant) are specifically designed to:
  1. Store long lists of numbers (Vectors).
  2. Perform **Nearest Neighbor Search** at lightning speed across millions of records. [\[58:40\]](#)
- **Use Case:** When you search "System Design" on YouTube, the system converts your text into a vector and finds videos whose title vectors are mathematically closest to yours—even if they don't use the exact words "System Design." [\[57:45\]](#)

---

## 5. Interview "Deep-Dive" Concept: Vector Arithmetic

A famous example in AI interviews is:

**Vector(King) - Vector(Man) + Vector(Woman) = Vector(Queen)** [\[47:17\]](#)

- **Explanation:** If you take the vector for "King" and subtract the "Male" attributes, you are left with the concept of "Royalty + Power." When you add the "Female" attributes to that concept, the mathematically closest vector in the space is "Queen." This proves that vectors capture **semantic meaning**, not just keywords. [\[51:01\]](#)

---

## Key Interview Terms to Remember:

- **Semantic Meaning:** The actual intent/meaning behind a word rather than its spelling.
- **N-Dimensional Space:** A mathematical space with 'N' number of features.
- **Vectorization:** The process of turning text/images/audio into numbers using a Neural Network. [\[45:06\]](#)

---

This is a **critical** for interview preparation because most candidates know *how to use* a Vector DB (like Pinecone/Chroma), but very few understand *how it works internally*.

---



# Lecture Notes: Vector Database Internals

## 1. The Core Problem: Speed vs. Accuracy

If you have 1 Billion vectors (e.g., all Amazon products), finding the "closest" vector to a user's query is hard.

- **Brute Force (Exact Nearest Neighbor - ENN):** Compare the query vector with **all 1 Billion** vectors.
    - *Pros:* 100% Accurate.
    - *Cons:* Extremely slow. Unusable in production. [\[09:50\]](#)
  - **The Solution (Approximate Nearest Neighbor - ANN):** We accept slightly lower accuracy (e.g., 99%) for massive speed gains. We organize vectors so we don't have to look at all of them.
- 

## 2. Four Key Algorithms for Vector Indexing

These are the internal mechanisms Vector DBs use to search fast. You **must** know HNSW and IVF for interviews.

### A. Cluster-Based (IVF - Inverted File Index) [\[13:24\]](#)

- **Concept:** Organize vectors into "clusters" (bunches).
- **How it works:**
  1. Use **K-Means Clustering** to group similar vectors (e.g., all "fruit" vectors in one cluster, "electronics" in another).
  2. Calculate a **Centroid** (center point) for each cluster.
  3. **Search:** When a query comes in, compare it *only* to the Centroids.
  4. If the query is close to the "Fruit Centroid," ignore all "Electronics" vectors and only search inside the Fruit cluster.
- **Pros:** Skips 90%+ of the data.
- **Cons:** "Edge Case Problem" - If a vector sits on the boundary of two clusters, you might miss it. (Solved by searching the top-3 closest clusters, not just one). [\[21:25\]](#)

### B. Tree-Based (Binary Space Partitioning) [\[23:13\]](#)

- **Concept:** Recursively cut the high-dimensional space in half, like a kd-tree or Binary Search Tree.
- **How it works:**
  1. Split the space into Left/Right based on a hyperplane.
  2. Keep splitting until you have small buckets of vectors.
- **Why it failed (Spotify used this until ~2023):** It performs poorly in high-dimensional space (The "Curse of Dimensionality"). If you make a wrong turn at the top of the tree, you can never recover the correct neighbor. [\[39:33\]](#)

### C. Graph-Based (HNSW - Hierarchical Navigable Small World) ★★ [\[42:22\]](#)

- **Status:** The Industry Standard (Used by Pinecone, Milvus, Weaviate).
- **Concept:** A multi-layered graph, similar to "Skip Lists" or an express highway system.
- **Structure:**
  1. **Layer 0 (Ground Level):** Contains **all** vectors connected to their nearest neighbors.
  2. **Layer 1 (Expressway):** Contains a few random vectors from Layer 0, linked with longer jumps.
  3. **Layer 2 (Super-Express):** Even fewer vectors for massive jumps.
- **Search Process:**
  1. Start at the top layer. Zoom across the graph to get to the general neighborhood of the query.
  2. Drop down to the lower layer, refine the search.
  3. Drop to Layer 0 for the final precision search.
- **Interview Benefit:** It balances speed and accuracy perfectly. Time complexity is  $O(\log N)$ .

#### D. Compression Method (Product Quantization - PQ) [52:17]

- **Concept:** Reducing the memory footprint of vectors (RAM is expensive!).
- **The Math:**
  1. A standard 1536-dim vector (OpenAI) takes ~6KB.
  2. 1 Billion vectors = **6 Terabytes of RAM**. This is too expensive.
- **How PQ Works:**
  1. Split the long vector into smaller "chunks" (e.g., 4 sub-vectors).
  2. Create a "Codebook" (like a dictionary) of common patterns for each chunk.
  3. Replace the actual decimal numbers with the **ID** of the closest pattern in the Codebook.
- **Result:** Reduces memory usage by **95%+** (e.g., 6TB  $\rightarrow$  12GB), allowing you to fit billion-scale datasets on a single server. [01:12:40]

### 3. Why not just use SQL/NoSQL? [01:23:17]

- **SQL/NoSQL (B-Trees):** Designed for **Exact Match** (e.g., WHERE price = 100).
- **Vector DB:** Designed for **Semantic Similarity** (e.g., WHERE meaning is close to "happy").
- **The "Psycho" Example:** You can't write a SQL query to SELECT \* FROM users WHERE personality = 'psycho'. But with vectors, you can find users whose comment history behaves similarly to known psycho profiles. [01:24:48]

### 4. Real-World Architecture (The "Hybrid" Approach) [01:13:47]

In production, we rarely use just one algorithm.

- **IVF + PQ:** Use IVF to narrow down the search space (speed), and use PQ to store the vectors in compressed format (memory).
- **HNSW + PQ:** The most common modern setup for high-performance retrieval.

## 5. Metadata Storage [01:20:10]

- **Crucial Concept:** You cannot "reverse" a vector to get the original text/image.
- Therefore, a Vector DB stores three things:
  1. **Vector:** The numbers for searching.
  2. **ID:** A unique key.
  3. **Metadata:** The actual content (e.g., text: "The quick brown fox", url: "youtube.com/...").
- When you search, the DB finds the Vector, but returns the Metadata.