

# Advanced Retrieval for Generative Models

## Need for Retrieval

### 1. Overview: Why Do Generative Models Need Retrieval?

Generative models, such as GPT or LLaMA, are powerful because they can **generate coherent text** by learning from massive datasets.

However, they face key **limitations**:

- They have a **fixed knowledge boundary** (whatever was in their training data).
- They **cannot access or recall** real-time or private information.
- They **hallucinate** when asked for specific facts not in their internal memory.

To overcome this, researchers introduced **Retrieval-Augmented Generation (RAG)** —

a method that combines **generation** (neural creativity) with **retrieval** (factual grounding).

### 2. Core Motivation for Retrieval

Problem	Why it Happens	How Retrieval Helps
<b>Stale Knowledge</b>	The model’s training data is frozen in time.	Retrieval fetches up-to-date documents or facts from external sources.
<b>Hallucination</b>	The model “fills in gaps” when it’s uncertain.	Retrieved evidence grounds responses in real data.
<b>Context Limitation</b>	Transformer models have limited context window (e.g.,	Retrieval dynamically selects only the <i>most relevant</i> documents instead of

Problem	Why it Happens	How Retrieval Helps
	8K–200K tokens).	loading everything.
<b>Domain Specificity</b>	Model wasn't trained on niche data (e.g., medical, legal).	Retrieval adds domain-specific info without retraining the model.

### 3. How Retrieval Works (Conceptual Workflow)

Let's break down a typical **Retrieval-Augmented Generation (RAG)** pipeline:

1. **User Query (Prompt)** → e.g., "Explain quantum computing in simple terms."
2. **Retrieval Step:**
  - Convert the query into an **embedding** (vector representation).
  - Search a **vector database** (e.g., FAISS, Chroma, Pinecone) for **similar documents**.
3. **Context Assembly:**
  - Top-k most relevant passages are **retrieved** and **appended to the prompt**.
  - This creates an **extended context** for the model.
4. **Generation Step:**
  - The LLM reads the *retrieved context + user query*.
  - Generates an answer grounded in those sources.
5. **Optional Post-Processing:**
  - Add citations, summarize, or rank multiple responses.

### 4. Mathematical / Intuitive View

If we denote:

- $x$ : user query
- $R(x)$ : retrieval function returning relevant documents
- $p_{\theta}(y|x, R(x))$ : generative model distribution (LLM)

Then, instead of generating purely from  $p_{\theta}(y|x)$ ,

The model conditions on the retrieved context:

$$p_{\theta}(y|x, R(x)) = \text{LLM}(y \mid x + \text{retrieved documents})$$

This makes the model **conditionally generative**, grounded by evidence.

---

## 5. Example: Retrieval in Practice

### Example 1: ChatGPT with Web Access

- When you ask about “latest NVIDIA GPU,” ChatGPT retrieves info from the web.
- It combines the retrieval results with language generation to form an accurate answer.

### Example 2: Enterprise Chatbot

- A company chatbot retrieves internal policy PDFs to accurately answer employee queries.
  - No need to retrain the LLM — retrieval bridges the knowledge gap.
- 

## 6. Benefits of Retrieval-Augmentation

Benefit	Explanation
<b>Dynamic Knowledge Updating</b>	The model can “know” new facts by updating the retriever database, not retraining.
<b>Smaller, More Efficient Models</b>	You can use smaller LLMs + retrieval instead of huge pre-trained ones.
<b>Transparency</b>	Retrieved documents can be shown to the user as evidence.
<b>Reduced Hallucinations</b>	Answers become more factual and reliable.

---

## 7. Limitations / Challenges

1. **Retriever Quality** – Poor embeddings or irrelevant retrievals reduce performance.

2. **Latency** – Fetching documents increases response time.
  3. **Context Length Limit** – You can't add too many retrieved docs; it may exceed the token limit.
  4. **Data Privacy** – Retrieving from sensitive data must be done securely.
- 

## 8. Summary (Quick Revision)

- Generative models **alone** lack factual grounding and up-to-date knowledge.
  - **Retrieval** adds an **external memory** that can be updated dynamically.
  - Together, they form **Retrieval-Augmented Generation (RAG)**.
  - Retrieval ensures the model's outputs are **accurate, context-aware, and trustworthy**.
- 

## Vector Databases

### 1. Overview: What Are Vector Databases?

A **Vector Database** is a special kind of database designed to **store, search, and manage vector embeddings** — numerical representations of data such as text, images, or audio.

In traditional databases, you might query with keywords or filters.

In contrast, **vector databases** allow you to query by *semantic meaning*.

👉 Example:

When you search “*doctor*”, the system can also find “*physician*” or “*medical professional*” — because their vector representations are **close in embedding space**.

Thus, vector databases are the **backbone of retrieval systems** used in **Retrieval-Augmented Generation (RAG)** and **semantic search**.

---

### 2. Motivation: Why Do We Need Vector Databases?

Problem	Why It Happens	Vector Database Solution
<b>Keyword Search is Limited</b>	Keyword-based matching doesn't capture meaning (e.g., "AI" ≠ "Artificial Intelligence").	Vector databases use embeddings that represent <b>semantic similarity</b> , not just exact text match.
<b>High-Dimensional Data</b>	Embeddings from models like BERT or CLIP are large (hundreds to thousands of dimensions).	Vector DBs are optimized for <b>high-dimensional indexing</b> and <b>efficient similarity search</b> .
<b>Scalability</b>	Millions of embeddings can't fit in normal RAM efficiently.	Vector DBs use approximate search and memory-efficient storage structures.
<b>Fast Retrieval for RAG</b>	LLMs need quick access to relevant documents before generation.	Vector DBs support <b>real-time top-k similarity search</b> for queries.

### 3. Core Concepts and Workflow

Let's break down how a **Vector Database** works step-by-step:

#### Step 1: Embedding Generation

- Input data (text, image, audio) is converted into **vector embeddings** using a model like:
  - Text: Sentence-BERT, OpenAI Embeddings, etc.
  - Image: CLIP, ResNet-based models.
- Each embedding is a point in a **high-dimensional space** (e.g., 768-dim vector).

#### Step 2: Indexing

- The database builds an **index** to make similarity search efficient.
- Common indexing techniques:
  - **FAISS (Facebook AI Similarity Search)** – product quantization + IVF.
  - **HNSW (Hierarchical Navigable Small World Graph)** – graph-based search.

- **Annoy (Approximate Nearest Neighbors Oh Yeah)** – tree-based structure.

### Step 3: Querying

- A user query is also converted into an embedding vector.
- The DB finds **nearest neighbors** based on **similarity metrics**:
  - **Cosine similarity**
  - **Euclidean distance**
  - **Dot product**

### Step 4: Retrieval

- The top-k most similar embeddings are fetched.
- The corresponding **documents or metadata** are returned to the application (e.g., RAG system).



## 4. Mathematical Intuition

If we denote:

- $q$  = query embedding
- $v_i$  = stored embedding for document  $i$

Then the **similarity** between  $q$  and  $v_i$  is measured as:

$$\text{sim}(q, v_i) = \frac{q \cdot v_i}{|q||v_i|}$$

The database returns top-k items where  $\text{sim}(q, v_i)$  is highest.

This is the essence of **nearest neighbor search** in vector space.



## 5. Example: How It Fits into a RAG Pipeline

### Workflow Example:

1. A user asks: *"Explain the Tesla autopilot system."*

2. The query → embedding vector  $q$ .
3. The vector DB searches for similar documents in the company's knowledge base.
4. It returns top 3–5 relevant passages.
5. These are appended to the prompt and sent to the **LLM**.
6. The model generates an accurate, context-aware answer.

## 6. Popular Vector Databases and Tools

Tool	Description	Strength
<b>FAISS</b> (Meta)	Library for efficient similarity search	Great for offline or research use
<b>Pinecone</b>	Managed vector DB service	Scalable, easy integration with APIs
<b>Weaviate</b>	Open-source with hybrid search (text + vector)	Semantic + keyword search
<b>Chroma</b>	Lightweight local DB, great for prototypes	Easy integration with LangChain
<b>Milvus</b>	Distributed, high-performance vector DB	Enterprise-grade scalability

## 7. Benefits

- **Semantic search** — retrieves meaning, not just words.
- **Real-time retrieval** — critical for live RAG applications.
- **Scalability** — supports millions/billions of embeddings efficiently.
- **Flexibility** — works with text, images, audio, multimodal data.
- **Integration** — easy to connect with LLM frameworks like LangChain, LlamaIndex.

## 8. Challenges

1. **High storage cost** – large embeddings require significant memory.
  2. **Index tuning** – choosing right index type (HNSW, IVF, etc.) affects accuracy/speed.
  3. **Data freshness** – frequent updates can slow down large indexes.
  4. **Privacy concerns** – embedding sensitive data must be done securely.
- 

## 9. Summary (Quick Revision)

- Vector databases store **embeddings** instead of raw text.
  - They allow **semantic similarity search** — finding meaningfully related data points.
  - Core operations: *embedding* → *indexing* → *similarity search* → *retrieval*.
  - Essential for **RAG**, **semantic search**, **recommendation systems**, and **AI assistants**.
  - Popular tools: FAISS, Pinecone, Weaviate, Milvus, Chroma.
- 

## Nearest Neighbor Search (NNS)

---

### 1. Overview: What Is Nearest Neighbor Search?

**Nearest Neighbor Search (NNS)** is the process of finding the data points in a dataset that are **most similar** (or closest) to a given query point.

In the context of **machine learning and vector databases**, each data item (like a sentence, image, or audio clip) is represented as a **vector embedding** in high-dimensional space.

When you issue a query:

- It is also converted into a vector.
- The system searches for **neighboring vectors** that are most similar to it.

 **Analogy:**



Imagine you're standing on a map (your query point), and you want to find the **closest cities** (neighbors) around you — that's what NNS does, but in 100s or 1000s of dimensions.

## 2. Motivation: Why We Need NNS

Problem	Explanation	How NNS Helps
<b>Semantic Search</b>	Keyword search can't find meaning-based matches.	NNS retrieves items <i>semantically close</i> to your query.
<b>High-Dimensional Data</b>	ML embeddings are vectors with hundreds of dimensions.	NNS efficiently finds the most similar vectors.
<b>RAG Systems</b>	Generative models need relevant data before answering.	NNS retrieves the most relevant context for the prompt.
<b>Recommendation Systems</b>	"Find similar items" to what a user liked.	NNS finds items closest to the user's preferences.

## 3. How Nearest Neighbor Search Works (Step-by-Step)

Let's say we have:

- A dataset of embeddings:  $D = v_1, v_2, \dots, v_n$
- A query embedding:  $q$

The goal is to find the **top-k nearest neighbors** of  $q$  in  $D$ , i.e., the most similar vectors.

### Step 1. Represent Data as Vectors

Each object (text, image, etc.) is represented by a dense vector embedding:

$$v_i = f(x_i)$$

where  $f$  is an embedding model.

### Step 2. Define a Similarity Metric

To measure closeness, we use metrics such as:

- **Cosine Similarity:**

$$\text{sim}(q, v_i) = \frac{q \cdot v_i}{|q||v_i|}$$

- **Euclidean Distance:**

$$d(q, v_i) = \sqrt{\sum_j (q_j - v_{ij})^2}$$

- **Dot Product:** For normalized vectors, it approximates similarity.

### Step 3. Compute Similarity

Compare the query vector with every stored vector to measure similarity or distance.

### Step 4. Rank and Retrieve

Sort all results by similarity (or distance) and return **top-k** neighbors — the most relevant or similar items.



## 4. Exact vs. Approximate Nearest Neighbor Search

Type	Description	Pros	Cons
<b>Exact NNS</b>	Compares the query to <i>every vector</i> in the dataset.	100% accurate	Very slow for large datasets (O(n))
<b>Approximate NNS (ANN)</b>	Uses clever data structures to <i>approximate</i> nearest neighbors quickly.	Fast and scalable	Small loss in accuracy

Most **vector databases** (like FAISS, Milvus, Pinecone) use **ANN algorithms** because they handle **millions or billions of vectors** efficiently.



## 5. Key Algorithms for Approximate NNS

## 1. HNSW (Hierarchical Navigable Small World Graph)

- Builds a graph where each node is connected to its neighbors.
- During search, the algorithm “jumps” across the graph to find the nearest nodes quickly.
- **Used in:** Milvus, Weaviate, Pinecone.

## 2. IVF (Inverted File Index)

- Clusters vectors into multiple groups (using k-means).
- Only searches in the *closest clusters* instead of the whole dataset.
- **Used in:** FAISS.

## 3. Product Quantization (PQ)

- Compresses vectors into smaller representations for memory efficiency.
- Balances speed and accuracy.

## 4. Annoy (Approximate Nearest Neighbors Oh Yeah)

- Uses multiple random projection trees.
- Lightweight and easy to use (used in Spotify recommendations).

---

## 6. Example: NNS in a Vector Database

Let's say we have these text embeddings:

Text	Vector (simplified)
"I love AI."	[0.8, 0.6]
"Machine learning is amazing."	[0.75, 0.65]
"The sky is blue."	[0.1, 0.2]

If the query is “Artificial Intelligence is great” → embedding [0.78, 0.64]

Then:

- Cosine similarity with “I love AI”  $\approx 0.999$

- Cosine similarity with "Machine learning..."  $\approx 0.997$
- Cosine similarity with "The sky is blue"  $\approx 0.3$

So, the **nearest neighbors** are the first two — semantically related to AI.

## 7. Applications of Nearest Neighbor Search

Domain	Use Case
<b>Information Retrieval</b>	Find documents similar to a query (used in RAG).
<b>Recommendation Systems</b>	Suggest products or movies similar to what a user liked.
<b>Computer Vision</b>	Retrieve similar images or detect near-duplicates.
<b>Anomaly Detection</b>	Outliers are data points with no close neighbors.
<b>Clustering</b>	k-NN algorithms use NNS as the foundation for grouping.

## 8. Challenges

### 1. **Curse of Dimensionality:**

As dimensions increase, distances become less meaningful and harder to compute efficiently.

### 2. **Trade-off between Speed and Accuracy:**

Faster approximate methods might skip some true neighbors.

### 3. **Index Maintenance:**

Frequent updates (inserts/deletes) can slow down large vector indexes.

### 4. **Hardware Demands:**

High-performance vector search often requires GPUs or optimized CPUs.

## 9. Summary (Quick Revision)

- **Nearest Neighbor Search** = finding the vectors most similar to a query vector.
- **Similarity metrics** like cosine or Euclidean distance define "closeness."

- **Exact search** is accurate but slow; **Approximate search (ANN)** is fast and scalable.
- Core algorithms: **HNSW, IVF, PQ, Annoy**.
- Backbone of **RAG, semantic search, recommendation**, and **vision retrieval** systems.

## Locality Sensitive Hashing (LSH)

### 1. Overview: What Is Locality Sensitive Hashing?

**Locality Sensitive Hashing (LSH)** is an algorithmic technique used to **speed up nearest neighbor search** in high-dimensional data by **hashing similar items into the same “buckets”** with high probability.

Instead of comparing every vector with every other vector (which is very slow for millions of items), LSH **narrows down the search space** —

You only compare your query with items in *the same or nearby buckets*.

 **Core idea:**

| Similar data points should have similar hash codes.

### 2. Motivation: Why We Need LSH

Challenge	Explanation	How LSH Helps
<b>High-dimensional data</b>	Exact distance calculations are expensive ( $O(n)$ ).	LSH reduces comparisons to only a subset of data.
<b>Scalability</b>	Datasets can have millions of embeddings.	Efficient indexing using hash buckets.
<b>Fast Approximation</b>	Perfect accuracy isn't always needed.	LSH gives <i>approximate nearest neighbors</i> very quickly.
<b>RAG &amp; Search Systems</b>	Retrieval must be real-time.	LSH enables sub-second retrieval in high-dimensional spaces.

### 3. Intuitive Explanation

Imagine you have a giant room full of books 📖.

If you're looking for books about "Deep Learning," you could:

- ❌ **Exact Search:** Check *every book* — accurate but slow.
- ✅ **LSH:** Put similar books on the *same shelf* based on their topic.

Then, when you search for "Deep Learning," you only look at that shelf.

That's what LSH does — it organizes similar data into nearby *hash buckets*, so you don't have to look everywhere.

---

### 4. Core Mathematical Idea

Traditional hash functions (like those used in hash tables) try to **disperse data uniformly**, minimizing collisions.

But LSH is *the opposite* — it's designed so that **similar vectors collide (hash to the same bucket) with high probability**.

Formally, a hash function family  $H$  is **locality sensitive** if for vectors  $x$  and  $y$ :

If  $\text{sim}(x, y)$  is high  $\Rightarrow P[h(x) = h(y)]$  is high

If  $\text{sim}(x, y)$  is low  $\Rightarrow P[h(x) = h(y)]$  is low

where  $h(\cdot)$  is a hash function from the family  $H$ .

So LSH **preserves locality** — similar items map to the same or close hash codes.

---

### 5. How LSH Works (Step-by-Step)

Let's go step-by-step through the process.

#### Step 1. Choose a Similarity Metric

You first decide what "similar" means —

Common metrics are:

- **Cosine similarity**
- **Euclidean distance**

- **Jaccard similarity** (for sets)

## Step 2. Define Locality-Sensitive Hash Functions

You choose or design a family of hash functions that align with your similarity measure.

For example:

- For **Cosine Similarity**, use **random hyperplanes**:

- Pick a random vector  $r$
- Define:

$$h(x) = \begin{cases} 1 & \text{if } r \cdot x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- Similar vectors likely fall on the same side of the hyperplane → same bit.

Each hyperplane gives one bit. Multiple hyperplanes = one hash code.

## Step 3. Build Hash Tables

- Create several hash tables, each using a *different* set of random hyperplanes.
- Each table stores points that fall into the same hash bucket.

## Step 4. Query

- Compute the hash code for the query vector.
- Retrieve only vectors from *the same bucket(s)* across tables.
- Compare the query only with those candidates — a small subset of the whole dataset.

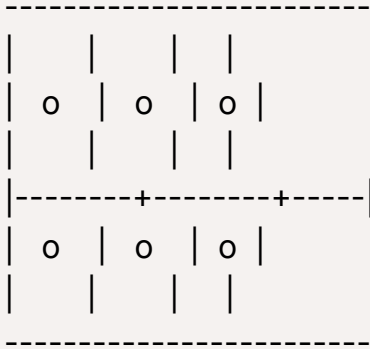
## Step 5. Compute True Similarity (Re-ranking)

- Among the retrieved candidates, compute the *actual* cosine or Euclidean distances.
  - Return the **top-k nearest neighbors**.
-



## 6. Visual Intuition (Text-based)

Imagine a 2D example:



Each region (square) is a **bucket**.

Points in the same region are considered "neighbors."

When a query comes in, we only check within its region instead of scanning the whole grid.

In high dimensions, these "regions" are formed by **random hyperplanes**.



## 7. Example

Suppose we have embeddings for sentences:

- "AI is revolutionizing the world"
- "Machine learning changes industries"
- "The ocean is blue"

When converted into vectors:

- The first two are *semantically close* (cosine  $\approx 0.9$ )
- The third is far away (cosine  $\approx 0.1$ )

Using LSH:

- First two vectors likely hash into **the same bucket** (e.g., 101010 )
- The third hashes into a **different bucket** (e.g., 000111 )



So during a query like "Artificial Intelligence advances,"

The database will look only in the bucket containing "AI-related" embeddings.

## 8. Variants of LSH (Based on Similarity Type)

Similarity Type	Example LSH Scheme	Description
<b>Cosine similarity</b>	Random hyperplane LSH	Divides space with random vectors.
<b>Euclidean distance</b>	p-stable LSH	Uses Gaussian projections.
<b>Jaccard similarity</b>	MinHash	Used for set similarity.
<b>Hamming distance</b>	Bit sampling LSH	Randomly picks subset of bits for hashing.

## 9. Advantages

- ✓ **Speed:** Sub-linear query time for large datasets.
- ✓ **Scalability:** Efficiently handles millions of embeddings.
- ✓ **Simplicity:** Easy to implement and understand conceptually.
- ✓ **No need for training:** Works directly on embeddings.

## 10. Limitations

- ✗ **Approximate results:** You might miss some true nearest neighbors.
- ✗ **High memory usage:** Multiple hash tables are needed for accuracy.
- ✗ **Parameter tuning:** Choosing the number of tables and hash functions affects performance.
- ✗ **Curse of dimensionality:** Efficiency decreases with extremely high dimensions (e.g.,  $> 1000$ ).

## 11. Summary (Quick Revision)

- **Locality Sensitive Hashing (LSH)** groups similar data into the same hash buckets.

- It enables **approximate nearest neighbor search** — much faster than brute-force.
  - Key idea: probability of hash collision increases with similarity.
  - Common techniques: **Random Hyperplane (cosine)**, **p-stable LSH (Euclidean)**, **MinHash (Jaccard)**.
  - LSH trades **accuracy for speed**, making it ideal for **vector search**, **deduplication**, and **RAG pipelines**.
- 

## Hierarchical Search Networks (HSNs)

---

### 1. Motivation — Why Hierarchical Search?

As the number of embeddings (vectors) in a database grows into **millions or billions**, a flat or brute-force search for the *nearest neighbor* becomes computationally expensive.

- **Brute-force (exact search)**:  $O(N)$  comparisons per query
- **Approximate Nearest Neighbor (ANN)**: reduces time by searching *likely regions* rather than the whole space

To achieve this, modern retrieval systems use **hierarchical search structures** like:

- **HNSW (Hierarchical Navigable Small World Graphs)**
- **FAISS IVF+PQ (Inverted File + Product Quantization)**

These structures organize vectors into **layers or levels** that progressively narrow down the search space — similar to how Google Maps zooms in from “World → Country → City → Street”.

---

### 2. Core Idea of Hierarchical Search Networks

**Goal:** Find nearest neighbors efficiently using a **multi-level graph** structure.

Each level (or hierarchy) represents the data at a different **granularity**:

- Top levels = fewer nodes, global overview

- Lower levels = denser, fine-grained local connections

When a query comes in:

1. The search starts from the **top layer**, where nodes represent rough “clusters” of the space.
2. The algorithm quickly navigates through large regions, moving toward nodes closer to the query vector.
3. Once it reaches a promising region, it **descends** into lower layers that have more detailed neighborhood connections.
4. At the lowest level, it performs a **fine-grained nearest neighbor search** among the most similar candidates.

---

### 3. Hierarchical Navigable Small World (HNSW) — The Core Example

Let’s look at **HNSW**, one of the most popular hierarchical search structures.

#### ◆ Construction:

- Each data point (vector) is assigned a random **maximum level (L)**.
- The network is a set of **layered graphs**, one per level.
- A node at level  $L_i$  is also present in all lower levels  $L_0 \dots L_{i-1}$ .
- Each node connects to a small number of **neighbors** (short-range) and a few **long-range** connections.
- The higher the level  $\rightarrow$  fewer nodes but more long-range links.

#### ◆ Search Algorithm:

1. Start from the **entry point** at the highest level.
2. At each level, use **greedy search** to move to the closest neighbor until no closer node is found.
3. Move down to the next level with that node as the new entry point.
4. Repeat until reaching the lowest level.

5. Perform a **local search** at the lowest level to find the top-k nearest neighbors. This allows logarithmic scaling of search time — roughly  **$O(\log N)$**  instead of  **$O(N)$** .

---

## 4. Analogy — “Zooming In” on Similar Data

Think of searching for a restaurant on Google Maps:

- At the world view (top level): You locate the right city.
- At city view (mid level): You locate the correct neighborhood.
- At street view (lowest level): You find the exact restaurant.

HSNs work similarly — progressively zooming into finer regions of the vector space.

---

## 5. Advantages

- ✓ **Scalability:** Efficient for millions/billions of vectors.
  - ✓ **Speed:** Search complexity close to  $O(\log N)$ .
  - ✓ **Dynamic:** Allows insertion of new vectors without full retraining.
  - ✓ **Accuracy:** High recall even with approximate methods.
  - ✓ **Widely used:** Core of FAISS (Meta), Milvus, Pinecone, Weaviate, and Elasticsearch vector search.
- 

## 6. Limitations

- ⚠ **Memory Overhead:** Stores multiple graph layers.
  - ⚠ **Build Time:** Construction can be slower for large datasets.
  - ⚠ **Parameter Tuning:** Number of connections (M) and search effort (efSearch) affect speed-accuracy trade-off.
- 

## 7. Use in Modern Retrieval-Augmented Generation (RAG)

In **RAG systems**, hierarchical search networks:

- Efficiently retrieve the top-k relevant embeddings for a given query (e.g., question, sentence).

- Work behind vector databases (e.g., FAISS-HNSW index).
  - Enable *fast context retrieval* for generative models like GPT, LLaMA, etc.
- 



## 8. Key Takeaways

- HSNs structure vectors in **multiple levels of graphs** to search efficiently.
  - They balance **speed** and **accuracy** through approximate neighbor navigation.
  - HNSW is the most common implementation, used in major vector databases.
  - It's a key enabler of **real-time semantic search** in RAG and GenAI applications.
- 



# Dense Passage Retrieval (DPR)

---



## 1. Motivation — Why Dense Retrieval?

Before DPR, traditional search systems (such as TF-IDF or BM25) used **sparse retrieval**, treating text as a bag of words.

- They rely on **exact keyword overlap** between the query and the document.
- Problem: If the query uses different wording (e.g., "car" vs. "automobile"), the retriever may **miss** relevant passages.



### Goal of DPR:

Learn a **semantic retrieval** method — where retrieval is based on **meaning**, not just word overlap.

In other words, DPR makes the retriever **understand language**, not just count words.

---



## 2. Core Idea

Dense Passage Retrieval (DPR) represents both **queries** and **passages** as **dense embeddings** (continuous vectors) using deep neural networks.

Then, it measures **semantic similarity** using a **dot product** or **cosine similarity** between embeddings.

So instead of:

Query: "What causes rain?"

Passage: "Rain is caused by condensation of water vapor."

DPR maps both into vectors close to each other in embedding space — even though their exact words differ.



### 3. Architecture

DPR uses a **dual-encoder** architecture with **two separate BERT models**:

#### ◆ a. Question Encoder ( $f(q)$ ):

- Takes the **query** (e.g., "Who wrote Harry Potter?")
- Outputs a **vector embedding**  $q = f_{q_i}(query)$

#### ◆ b. Passage Encoder ( $f(p)$ ):

- Takes a **document/passage** (e.g., "J.K. Rowling is the author of the Harry Potter series.")
- Outputs a **vector embedding**  $p = f_{p_i}(passage)$

Then:

$$\text{Similarity}(q, p) = q \cdot p$$

(Dot product or cosine similarity)

During inference, the system finds the top-k passages with the highest similarity to the query.



### 4. Training Objective — Contrastive Learning

DPR is trained using **supervised contrastive learning**.

Each training sample has:

- A **query (q)**
- A **positive passage (p<sup>+</sup>)** — relevant to the query

- **Negative passages ( $p^-$ )** — irrelevant distractors

The training goal:

maximize  $\text{sim}(q, p^+)$  and minimize  $\text{sim}(q, p^-)$

## Loss Function:

$$L = -\log \frac{\exp(\text{sim}(q, p^+))}{\exp(\text{sim}(q, p^+)) + \sum_i \exp(\text{sim}(q, p_i^-))}$$

This is a **softmax loss over similarities**, encouraging positive pairs to have the highest similarity.

---



## 5. Indexing & Retrieval Workflow

Once trained, DPR works in two stages:

### 1. Embedding & Indexing:

- Precompute embeddings for all passages using the passage encoder.
- Store them in a **vector database** (e.g., FAISS, Pinecone, Weaviate).

### 2. Retrieval:


- Encode the incoming query using the query encoder.
  - Perform **Approximate Nearest Neighbor (ANN)** search to find top-k similar passages (using techniques like HNSW or IVF-PQ).
  - Return top-k passages as **retrieved context** for downstream tasks (e.g., RAG or QA).
- 



## 6. Example

**Query:** "Who discovered gravity?"

**Passage 1:** "Isaac Newton formulated the law of universal gravitation." 

**Passage 2:** "Gravity is a force acting between masses." 

Even though both mention "gravity," DPR will score Passage 1 higher because it *semantically answers* the question.

---

## 7. Comparison — Sparse vs. Dense Retrieval

Feature	Sparse Retrieval (BM25, TF-IDF)	Dense Retrieval (DPR)
Representation	Bag of Words	Embedding vectors
Matching	Exact keyword overlap	Semantic similarity
Model	Statistical	Neural (BERT-based)
Memory	Large inverted index	Compact vector index
Generalization	Weak to synonyms	Strong semantic understanding

---

## 8. Real-World Applications

✓ **Retrieval-Augmented Generation (RAG):** DPR finds relevant passages to feed into an LLM's context window.

✓ **Open-domain QA systems:** Used in models like **REALM**, **FiD**, and **RAG**.

✓ **Chatbots & Search Engines:** Improves retrieval quality beyond keyword search.

✓ **Academic & Legal search:** Finds semantically similar paragraphs or papers.

---

## 9. Advantages

✓ Captures **semantic meaning**, not just lexical overlap.

✓ Scalable when paired with **ANN search**.

✓ Strong **generalization** to unseen queries.

✓ Compatible with modern **vector databases**.

✓ Can be **fine-tuned** for domain-specific retrieval.

---

## 10. Limitations

⚠ Requires **large labeled datasets** for training (query–passage pairs).

⚠ **Retraining is needed** for domain adaptation.

⚠ High **memory usage** for storing dense vectors.

⚠ Sensitive to encoder quality — poor embeddings reduce recall.





## 11. Key Takeaways

- DPR replaces keyword-based retrieval with **deep semantic retrieval**.
- Uses **dual BERT encoders** to map queries and passages into the same embedding space.
- Trained with **contrastive loss** on positive/negative pairs.
- Works hand-in-hand with **vector databases** and **ANN algorithms**.
- Core building block of modern **RAG and QA systems**.



## Multi-Vector Retrieval (MVR)



### 1. Motivation — Why Go Beyond Single-Vector Retrieval?

In **Dense Passage Retrieval (DPR)**, each passage (document) is represented by a **single vector** — a single point in the embedding space.

But this approach has a limitation:



A single vector cannot capture multiple semantic aspects of a long or information-rich passage.

Example:

Passage:

“Isaac Newton discovered gravity. He also developed calculus.”

- One part is about **gravity**.
- Another part is about **calculus**.

If the passage is represented by only one vector, it might lie “between” these two meanings — not close enough to queries about *either* topic.



**Multi-Vector Retrieval (MVR)** solves this problem by representing each passage (and sometimes each query) with **multiple embeddings**, allowing **finer-grained semantic matching**.

## 2. Core Idea

Instead of using one embedding per document, MVR breaks it down into **multiple token-level or chunk-level vectors**, each capturing a different semantic unit.

Then, during retrieval:

- Each query vector is compared against **all passage vectors**.
- The **maximum or aggregated similarity score** determines relevance.

Formally:

$$\text{Score}(q, p) = \max_{i,j}; \text{sim}(q_i, p_j)$$

where:

- $q_i$  = the i-th query vector
- $p_j$  = the j-th passage vector
- sim = dot product or cosine similarity

This allows the retriever to match **different concepts** within the same document.

---

## 3. Architecture Overview

### ◆ a. Query Encoder

- Encodes the query into **one or multiple vectors** (depending on design).
- Often uses the [CLS] token for a single vector or the full set of contextual token embeddings for multi-vector queries.

### ◆ b. Document Encoder

- Encodes each passage into a **set of token vectors** (e.g., output of BERT without pooling).
- Each token (or a subset) represents a semantic unit within the document.

### ◆ c. Scoring Function

- For a given query vector (or vectors), compute similarity with all document vectors.

- The final passage score can be:
  - **Max pooling:** take the maximum similarity
  - **Mean pooling:** average all similarities
  - **Attention pooling:** learn weights to combine similarities

This flexibility enables rich, fine-grained matching.

---


## 4. Example

Query:

| "Who developed calculus?"

Passage:

| "Isaac Newton discovered gravity. He also developed calculus."

- DPR: One vector captures both "gravity" and "calculus" mixed — might not be close enough to query.
- MVR:
  - Token 1 vector → "gravity"
  - Token 2 vector → "calculus"
  - Query vector matches strongly with "calculus" token vector. 

Result: The system retrieves the passage correctly because one of its token vectors matches the query meaning closely.

---

## 5. Training — Late Interaction Models (e.g., ColBERT)

The **ColBERT** model (2020, by Stanford) is the classic example of multi-vector retrieval.

It introduces **Late Interaction**, which means:

- The query and passage are encoded **independently**.
- Interaction (similarity computation) happens **after** encoding — at retrieval time.

This makes retrieval scalable because you can **pre-compute passage embeddings** once and store them.

## Loss Function:

Contrastive loss similar to DPR, but using max similarity over token pairs:

$$L = -\log \frac{\exp(\max_{i,j} \text{sim}(q_i, p^+ j))}{\sum p^- \exp(\max_{i,j} \text{sim}(q_i, p^- j))}$$

---

## 6. Retrieval Workflow

### 1. Indexing:

- Each passage is encoded into multiple token embeddings.
- All token vectors are stored in a **vector index** (like FAISS).

### 2. Query Processing:

- Query is encoded into its token embeddings.
- For each query vector, compute nearest neighbors among all passage vectors.

### 3. Scoring:

- Aggregate token-level similarities into passage-level scores.
- Return top-k passages as final retrieved results.

---

## 7. Advantages

- ✓ **Better semantic coverage:** Handles passages with multiple topics.
- ✓ **Higher recall:** Each sub-vector acts as a mini “semantic anchor.”
- ✓ **Interpretable:** You can see which tokens contributed most to the match.
- ✓ **Compatible with ANN search:** Still scalable using indexing structures like IVF-PQ or HNSW.

---

## 8. Limitations

- ⚠️ **Larger storage footprint:** Multiple vectors per document.
- ⚠️ **Higher retrieval cost:** More comparisons during search.
- ⚠️ **Index complexity:** Managing large vector sets is resource-intensive.
- ⚠️ **Latency:** More similarity computations unless optimized with pruning.

## 💡 9. Key Models in Multi-Vector Retrieval

Model	Approach	Key Feature
<b>CoBERT (2020)</b>	Late Interaction	Token-level matching via max similarity
<b>CoBERT-v2 (2022)</b>	Compressed embeddings	Faster, smaller storage
<b>PLAID / SPLADE</b>	Sparse-dense hybrid	Efficiency with interpretability
<b>MEVI (2023)</b>	Multi-embedding video retrieval	Extends MVR to multimodal settings

## 10. Key Takeaways

- **Dense retrieval (DPR):** 1 vector per passage — semantic but coarse.
- **Multi-vector retrieval (MVR):** Multiple vectors per passage — captures fine-grained meaning.
- **CoBERT:** Foundational architecture enabling efficient multi-vector search.
- MVR enables **context-aware, token-level semantic retrieval**, crucial for **RAG**, **QA**, and **multimodal retrieval** systems.