

# Faiss: Facebook AI Similarity Search – A Comprehensive Overview

## Introduction and Theoretical Foundations of Similarity Search

In traditional databases, data is stored in structured tables and queries rely on exact matches (e.g. key-value lookups). This approach struggles when we need to find *similar* items rather than exact matches. Modern AI systems represent data (images, text, etc.) as high-dimensional *vectors* (embeddings) that capture semantic content <sup>1</sup> <sup>2</sup>. Similar items produce vectors that are *nearby* in some metric space (e.g. Euclidean distance or cosine similarity), so similarity search boils down to a nearest-neighbor search in a high-dimensional vector space <sup>2</sup>. For example, an image of a city hall can be encoded as a vector such that images of the same building yield vectors close in Euclidean space <sup>3</sup>. Likewise, a text or image classifier might embed inputs as vectors and use a **maximum inner product search** (MIPS) to find items with highest dot-product similarity to a query vector <sup>4</sup>.

Formally, given a set of  $N$  data vectors and a query vector, a similarity search finds the  $k$  nearest neighbors (k-NN) of the query under a chosen distance or similarity metric. The two most common operations are:

- **Nearest neighbor search (L2)**: Given a query vector, retrieve the database items with smallest Euclidean distance to the query <sup>5</sup>.
- **Maximum inner-product search (MIPS)**: Given a query vector (e.g. from a classifier), retrieve items with largest dot-product (cosine similarity if vectors are normalized) <sup>6</sup> <sup>5</sup>.

Performing these operations efficiently at **large scale** (millions or billions of vectors) is challenging. A brute-force scan of all vectors per query is exact but becomes infeasible as  $N$  grows large (the curse of dimensionality) <sup>7</sup> <sup>8</sup>. Traditional SQL databases and search engines are not designed for high-dimensional similarity queries <sup>9</sup> <sup>10</sup>. This gap led to specialized methods for *Approximate Nearest Neighbor* (ANN) search, which trade a small amount of accuracy for huge gains in speed and memory efficiency <sup>11</sup>. In practice, allowing a slight deviation from exact results (e.g. ~90% recall instead of 100%) can yield orders-of-magnitude faster search <sup>11</sup>. The key is to build an **index** on the vectors – a data structure that accelerates search by preprocessing the vector dataset <sup>12</sup>.

**Faiss (Facebook AI Similarity Search)** was developed by Facebook AI Research in 2017 to address these challenges <sup>13</sup>. Faiss is an open-source library that enables efficient similarity search in collections of dense vectors, even at billion-scale. It provides a unified interface to a wide range of ANN algorithms and supports both exact and approximate searches. Critically, Faiss is optimized for high performance: it offers **multiple indexing methods** covering different speed/accuracy trade-offs, is highly optimized in C++ for memory efficiency, and includes **GPU acceleration** for heavy computations <sup>14</sup>. Upon its release, Faiss achieved state-of-the-art results, with implementations 8× faster than previous methods on billion-scale benchmarks and the ability to construct the first k-NN graph on 1B vectors <sup>15</sup>. In summary, Faiss provides the toolkit to make similarity search *fast, scalable, and flexible* – a necessity for modern AI applications dealing with massive embedding datasets.

## Overview of Faiss and Its Key Features

Faiss (short for *Facebook AI Similarity Search*) is a library specifically designed for fast *nearest neighbor search* and clustering on dense vector data <sup>16</sup>. It was open-sourced by Facebook (Meta) AI Research and is widely used in the industry. Faiss can handle *very large* collections – from millions to billions of vectors – while optimizing the three key metrics of similarity search: **speed**, **memory usage**, and **accuracy** <sup>17</sup>. To achieve this, Faiss implements a variety of indexing algorithms (described in the next section) that allow tuning the trade-off between search latency, index memory footprint, and result quality <sup>14</sup> <sup>18</sup>. Unlike naive brute-force search, these indexes perform *compression* or *partial search* to avoid comparing a query to every single vector in the dataset <sup>11</sup> <sup>19</sup>.

**Unique strengths of Faiss** include its breadth of methods, efficiency, and hardware support. Firstly, Faiss provides *several index types* – from exact flat indexes to advanced compressed or graph-based indexes – often composable to meet specific use-cases <sup>14</sup>. This means developers can choose an index that best fits their needs (high recall vs. low latency vs. small memory). Second, Faiss is highly optimized in both memory and compute. Many Faiss indexes use compact vector representations (down to a few bytes per vector) so that even billion-vector datasets can reside in RAM <sup>20</sup>. Its search routines use efficient algorithms (e.g. block mass distance computations) and make maximal use of CPU cache, SIMD instructions, etc., for speed <sup>21</sup> <sup>22</sup>. Third, Faiss offers a state-of-the-art **GPU implementation** for the most important indexes <sup>23</sup>. On compatible hardware, GPU-accelerated Faiss can be an order of magnitude faster than CPU (20× faster on modern GPUs in some cases) <sup>24</sup>, enabling real-time search on very large datasets. Importantly, the Faiss GPU indexes have the *same interface* as CPU indexes, and Faiss provides easy transfer of data between CPU and GPU indices <sup>25</sup>. This flexibility allows seamless scaling from development (CPU) to production (GPU clusters). Finally, Faiss is written in C++ with bindings for Python, making it easy to integrate. The Python API closely mirrors the C++ classes, so prototypes in Python can be translated to optimized C++ if needed <sup>25</sup>.

In summary, Faiss is a powerful, scalable similarity search library that balances **high performance** (through optimized algorithms and GPU support) with **flexibility** (many indexing strategies and combination options). It enables use-cases that were previously impractical, such as interactive search over billions of high-dimensional vectors in applications like visual search, semantic text search, and recommender systems <sup>26</sup> <sup>27</sup>. Next, we delve into the core indexing methods provided by Faiss and explain how each works.

## Core Indexing Methods in Faiss

Faiss implements a variety of indexing methods, each with different strategies to accelerate similarity search. Many of these can be used standalone or combined as *composite indexes*. Below we overview the major index types supported in Faiss, including their theoretical concepts and practical use.

### Flat Indexes (Brute-Force: IndexFlatL2 and IndexFlatIP)

The simplest index in Faiss is the **flat index**, which doesn't do any approximation or compression – it stores all vectors “as is” and computes distances to *every* vector during a search. This is essentially brute-force kNN search, but implemented efficiently in C++/SIMD. Flat indexes guarantee *exact* results at the cost of linear scan time. Faiss provides `IndexFlatL2` for Euclidean distance (L2) and `IndexFlatIP` for inner product (dot-product) similarity. (For cosine similarity, one can use `IndexFlatIP` on normalized vectors <sup>28</sup>.)

Using a flat index is straightforward. We must specify the vector dimensionality  $d$  when constructing the index. No training is needed since no clustering or compression is done (the index is “trained” by default). We can then add vectors to the index and perform searches. Below is a simple example in Python using Faiss: we build a flat L2 index, add a set of vectors `xb` (numpy array of shape  $N \times d$ ), and query it with some vectors `xq` to retrieve the top- $k$  nearest neighbors:

```
import faiss                                     # Faiss module
index = faiss.IndexFlatL2(d)                     # create flat index for L2 distance
index.add(xb)                                    # add database vectors (xb: N x d array)
print(index.ntotal)                              # number of vectors indexed

# Query with a set of vectors xq (nq x d array)
k = 4                                             # number of nearest neighbors to find
D, I = index.search(xq, k)                       # search the index
print(I[:5])                                     # print neighbors of first 5 queries
```

In this snippet, `IndexFlatL2(d)` allocates an index for  $d$ -dimensional vectors. Calling `index.add(xb)` copies the data into the index's storage (flat indexes support only sequential IDs 0.. $N-1$  by default). A flat index does *no* preprocessing of vectors – it just stores them contiguously in memory <sup>29</sup>. The search (`index.search`) computes distances from each query in `xq` to all indexed vectors and returns the top- $k$  results for each query. As expected, flat search is exact but can be slow for large  $N$ , since it's  $O(N)$  per query. Flat indexes are best suited for **small datasets** or when accuracy is paramount and query latency is less a concern <sup>30</sup>. They are also useful to produce a gold-standard result for evaluating the accuracy of other ANN indexes.

**Memory/Speed:** A flat index stores each vector in full (4 bytes \*  $d$  per vector for floats), with minimal overhead <sup>31</sup>. This is quite memory-efficient given no additional structure (aside from storing IDs if needed). However, search speed scales linearly with the number of vectors. For example, scanning 1 million 128-dimensional vectors might take tens of milliseconds on a modern CPU – which is fine for a few queries, but becomes a bottleneck if queries are frequent or  $N$  is much larger. Faiss does optimize the brute-force computations (using batch SIMD distance calculations), and the GPU version can compute distances extremely fast (billions of dot products per second) <sup>32</sup>. But in general, flat indexes trade speed for perfect accuracy. The Faiss wiki notes that flat indexes are the **only choice for exact results**, and they serve as a baseline for approximate methods <sup>33</sup>. If you only have a handful of queries (e.g. offline analysis), a flat index may be acceptable or even optimal (since building a complex index might not pay off) <sup>34</sup>.

## Inverted File Index (IVF) – Partitioning into Voronoi Cells

One of the most popular ANN strategies in Faiss is the **inverted file index (IVF)**, also known as a *cell-probe* or *coarse quantization* method. The core idea is to partition the vector space into *clusters* (cells) using a clustering algorithm like *k-means*, and then restrict search to only a few relevant clusters instead of the entire dataset <sup>35</sup> <sup>36</sup>. This significantly reduces the number of distance computations per query.

In Faiss's IVF implementation (`IndexIVF*` classes), one first chooses a number of clusters (often denoted *nlist* = number of lists). During *training*, Faiss performs *k-means* on a sample of the dataset to learn *nlist* centroids (each centroid is a  $d$ -dimensional vector). These centroids define the partition of the space. The

IVF index uses an underlying *quantizer* to assign any vector to its nearest centroid (by L2 distance) <sup>37</sup>. When we add vectors to an IVF index, each vector is stored in the list (bucket) of its nearest centroid – this is the inverted file, essentially a collection of *nlist* buckets of vector IDs <sup>38</sup> <sup>39</sup>.

At query time, searching an IVF index involves two stages <sup>38</sup> <sup>40</sup>:

1. **Coarse search (quantization):** The query vector is first projected to the nearest centroid(s). Instead of searching all lists, we choose a subset of the most promising clusters. Faiss uses a parameter *nprobe* = number of clusters to search. The query is compared to all centroids, and the top *nprobe* closest clusters are selected <sup>38</sup>. These clusters likely contain the nearest neighbors of the query.
2. **Refinement search within clusters:** The search then visits only those *nprobe* lists and computes exact distances between the query and the vectors stored in those lists <sup>38</sup>. From these, the overall top-*k* nearest neighbors are produced.

By adjusting *nprobe*, one can trade search accuracy vs. speed: searching more clusters (*nprobe* closer to *nlist*) yields higher recall (approaches exhaustive search) at the cost of more comparisons <sup>41</sup>. In practice, IVF can drastically accelerate search because *nprobe*  $\ll$  *nlist* and typically each list holds  $\sim N/nlist$  vectors (so each query only compares to  $\sim (nprobe/N)$  vectors on average) <sup>42</sup>. The main failure mode is if a query's true nearest neighbor falls in a cluster that wasn't among the *nprobe* searched – hence larger *nprobe* improves accuracy <sup>43</sup>.

**Using IVF in Faiss:** An IVF index is often created by specifying a *coarse quantizer* (often a flat index for the centroids) and the number of lists. For example, an index with 100 clusters can be built as follows:

```
# Define coarse quantizer and IVF index
quantizer = faiss.IndexFlatL2(d)                # underlying quantizer
(centroid finder)
index_ivf = faiss.IndexIVFFlat(quantizer, d, nlist) # IVF index (lists with
full vectors)
index_ivf.train(xb)                             # train the quantizer on
data
index_ivf.add(xb)                               # add data points to the
index
# Once built, we can search the IVF index:
index_ivf.nprobe = 5                            # e.g. search 5 nearest
clusters
D, I = index_ivf.search(xq, k)
```

In this snippet, `IndexIVFFlat` means the index uses IVF with *flat storage in each list* (i.e. storing full vectors in the inverted lists). We pass a pre-defined `IndexFlatL2` as the *quantizer* and specify *nlist*. We then call `train()` with sample data, which runs *k*-means to compute centroids (this step is required for IVF indexes) <sup>44</sup>. After adding data, we set `nprobe` at query time to decide how many clusters to probe for each query <sup>45</sup>.

**Memory/Speed:** IVF indexes are an excellent **scalability** option. Storing vectors in clusters has minimal memory overhead – just an extra cluster ID or list pointer per vector (Faiss stores an 8-byte id by default)

<sup>46</sup> <sup>47</sup> . For `IndexIVFFlat`, each vector is stored in full plus 8 bytes for its ID in the list <sup>46</sup> . The big win is that search complexity becomes  $O(nprobe * (N/nlist))$  on average instead of  $O(N)$ . For instance, if  $nlist = 10000$  and  $nprobe = 10$  for a million-vector dataset, each query compares to roughly 10% of the vectors instead of 100%. Faiss further optimizes by computing the query-centroid distances quickly (often negligible cost compared to the second stage) and by supporting GPU acceleration for IVF search (centroid assignment and list scanning can be done on GPU) <sup>23</sup> <sup>48</sup> . Empirically, IVF is a good *general-purpose ANN index* with balanced trade-offs: it offers high recall with appropriate parameters and scales to large  $N$  with reasonable memory. Pinecone's benchmarks note that IVF provides “good search-quality, good search-speed, and reasonable memory usage” for large-scale search <sup>49</sup> <sup>50</sup> .

A critical parameter is  $nlist$  (number of clusters). A common heuristic is to set  $nlist$  on the order of  $\sqrt{N}$  (e.g. for 1M vectors, a few thousand clusters) <sup>51</sup> , or more generally to balance the cost of scanning clusters vs. scanning vectors in lists <sup>52</sup> . The Faiss wiki suggests `nlist = C * sqrt(n)` for some constant  $C$  (like 4 to 16) as a rule of thumb <sup>52</sup> . Too few clusters and each list is large (less pruning benefit); too many clusters and you must raise  $nprobe$  or risk missing neighbors. Faiss also allows using alternative quantizers – e.g. a hierarchical quantizer or even an HNSW index to assign vectors to IVF lists (improving assignment speed/accuracy) <sup>53</sup> , but the basic principle remains clustering.

## Product Quantization (PQ) – Compact Vector Codes

**Product Quantization (PQ)** is a core technique in Faiss for compressing vectors into short codes, markedly reducing memory and speeding up distance computations <sup>21</sup> <sup>54</sup> . PQ is a form of vector quantization that partitions the feature space into a Cartesian product of lower-dimensional subspaces and quantizes each subspace independently <sup>55</sup> <sup>56</sup> . In simpler terms, PQ breaks each  $d$ -dimensional vector into  $m$  smaller subvectors, then quantizes each subvector to the nearest prototype in a small codebook. Each subvector can be encoded by a few bits (e.g. an 8-bit code indicating one of 256 possible centroids for that subvector). The concatenation of these subvector codes forms a compact code for the entire vector.

For example, suppose  $d = 128$  and we choose  $m = 16$  subvectors, each of length 8. We train 16 separate quantizers (via k-means) – one for each subspace – each with say 256 centroids (8 bits). Then any 128-D vector is split into 16 chunks of 8-D; each chunk is replaced by the index of the nearest centroid in the corresponding codebook. The vector is thus represented by 16 bytes (128 bits) instead of  $128 \times 4 = 512$  bytes (if using 32-bit floats). Typical Faiss usage is  $m$  between 8 and 64. The number of bits per subquantizer can also be varied (commonly 8 bits, but Faiss supports 4 or even 6, etc.) <sup>57</sup> . An  $m$ -byte PQ code is effectively an integer in a  $256^m$  discrete space.

Faiss uses PQ in several ways: **IndexPQ** applies PQ to create a compressed flat index (exhaustive search in the compressed domain), and **IndexIVFPQ** combines IVF partitioning with PQ compression of vectors in each list (often called “IVFADC” in literature) <sup>58</sup> <sup>59</sup> . There are also variants like `IndexIVFPQR` which add a second-level refinement PQ for re-ranking <sup>60</sup> <sup>61</sup> . PQ indexes require a training phase to learn codebooks (quantizers) for each subspace, usually done on a sample of vectors.

**Intuition:** PQ dramatically compresses vectors with minimal loss of accuracy when done properly. The reason it works is that high-dimensional data often has correlated components or lies on a manifold, so we can get away with coarse-grained representation in each subspace. Distance computations can also be accelerated: Faiss precomputes distance lookup tables for each codebook with respect to a query, so it can compute approximate distances by summing subcomponent distances (this is faster than computing full

Euclidean distance in high-dim) <sup>62</sup> <sup>63</sup>. In an IndexPQ (flat PQ), search is done purely in the compressed domain: the index stores only PQ codes, and during search Faiss computes approximate distances by decoding partially or using precomputed tables <sup>64</sup> <sup>62</sup>. This yields an *approximate* result (since we're not using exact original vectors), but the error can be small if PQ is fine-grained enough. Often a *re-ranking* step is applied: find top-k by compressed distance, then re-order those by the true distance using the original vectors (if available) <sup>65</sup> <sup>66</sup>. IndexIVFPQR in Faiss automates this by storing additional PQ codes to refine the distance for top candidates <sup>60</sup>.

**Using PQ in Faiss:** We can create a standalone PQ index as follows:

```
m = 16                                # number of subquantizers
nbits = 8                             # bits per subquantizer code (8 gives 256
centroids each)
pq_index = faiss.IndexPQ(d, m, nbits)
pq_index.train(x_train)                # train PQ codebooks on sample (or data)
pq_index.add(x_base)                  # add base vectors (stores compressed codes)
D, I = pq_index.search(x_query, k)    # search (distance in PQ-compressed space)
```

Here IndexPQ(d, m, nbits) creates a flat PQ index that will compress vectors into  $m \times nbits$  bits. After training, adding vectors will store only their PQ codes (not the full vector). The search returns IDs and approximate distances based on codes <sup>67</sup> <sup>68</sup>. In practice, one would often use IndexIVFPQ for large datasets: first partition into IVF lists, then within each list store PQ-compressed residuals (difference between vector and centroid) <sup>69</sup> <sup>70</sup>. Faiss allows constructing an IndexIVFPQ by providing a coarse quantizer and PQ parameters. For example:

```
coarse_quantizer = faiss.IndexFlatL2(d)
index_ivfpq = faiss.IndexIVFPQ(coarse_quantizer, d, nlist, m, nbits)
index_ivfpq.train(x_train)
index_ivfpq.add(x_base)
index_ivfpq.nprobe = 5
D, I = index_ivfpq.search(x_query, k)
```

This builds an IVF index with *nlist* clusters, and within each cluster stores vectors as PQ codes of size  $m \times nbits$  bits <sup>69</sup> <sup>70</sup>. The train() step here performs both the k-means for coarse quantizer and the training of the PQ codebooks on residuals. At search time, Faiss will assign the query to nearest centroids (like IVF), then for each candidate code compute an approximate distance using PQ (with precomputed tables), and return the top-k. One can optionally then refine these results with the original vectors if higher accuracy is needed (IVFPQ doesn't do that by default unless using IVFPQR).

**Memory/Speed:** PQ offers **huge memory savings**. For instance, using  $m=16$ ,  $nbits=8$  compresses a vector to 16 bytes, regardless of original dimension. Storing 1 billion 128-D vectors with such a PQ would require only ~16 GB (plus overhead) instead of 512 GB if stored in full floats <sup>19</sup>. This compression is what enables Faiss to tackle billion-scale (Faiss focuses on methods that compress vectors since only those scale to extremely large datasets) <sup>20</sup>. Search speed also improves because distance computations are faster on small codes (often using lookup tables). The trade-off is some accuracy loss – the retrieved neighbors might not always

be the true nearest. However, PQ can achieve high recall if enough bits are used. Faiss even implements *optimized product quantization (OPQ)*, which applies a learned rotation to the vectors before quantization to make them more amenable to PQ (improving accuracy for the same code size) <sup>71</sup> <sup>72</sup>. OPQ is often recommended as a pre-processing if using PQ heavily.

In summary, PQ is a powerful method for **memory-efficient** indexing. It is often used in combination with IVF (coarse quantization + fine PQ quantization) to get the benefits of both – coarse clustering for pruning and PQ for compression. Indeed, the *IVF-PQ* combo (a.k.a. ADC: asymmetric distance computation) is referred to as “probably the most useful indexing structure for large-scale search” <sup>70</sup>, as it scales to billions of vectors with manageable memory and high search speeds, while maintaining good accuracy when tuned properly.

## Locality-Sensitive Hashing (LSH) – Hash Bucket Indexing

**Locality-Sensitive Hashing (LSH)** is another technique Faiss provides ( `IndexLSH` ) for ANN search, based on hashing vectors into buckets such that similar vectors collide in the same buckets with high probability. LSH uses *randomized projections* or rotations of the vector space to produce hash codes. The idea is to generate binary codes for vectors and use Hamming distance as a proxy for original distance – vectors with small Euclidean distance are likely to share many bits. A query only needs to search vectors with hash codes near to the query’s hash.

Faiss’s LSH implementation takes an input dimension  $d$  and a number of bits  $nbits$  for the code. Internally, Faiss uses an improved scheme where if  $nbits \leq d$ , a random rotation is applied to the vectors followed by taking the sign of each of  $nbits$  coordinates (if  $nbits > d$ , it uses a “tight frame” projection to get more bits) <sup>73</sup>. This generates an  $nbits$ -bit binary code for each vector, effectively mapping vectors to corners of an  $nbits$ -dimensional hypercube. These codes are stored in a flat binary index and compared via Hamming distance <sup>74</sup> <sup>75</sup>. During search, the query is similarly hashed to a binary code, and the index finds items with the most similar codes (smallest Hamming distance). One can control the recall by adjusting  $nbits$ : more bits means a finer-grained hashing (fewer collisions, so higher precision) at the cost of larger codes and slower search.

**Using LSH in Faiss:** Construction is simple – choose  $nbits$  (often on the order of a multiple of dimension, e.g.  $2 \times d$  as a starting point <sup>76</sup>) and initialize the index:

```
n_bits = 2 * d
lsh = faiss.IndexLSH(d, n_bits)
lsh.train(x_train)          # train may be no-op for LSH (random rotations)
lsh.add(x_base)             # add base vectors (now stored as binary codes)
D, I = lsh.search(x_query, k)
```

LSH does not require data-dependent training in the classic sense (the “train” call in Faiss just sets up random projections if needed). After adding vectors, the index holds their binary codes <sup>74</sup> <sup>77</sup>. The search returns the IDs of the nearest neighbors by Hamming distance. Note that those are approximate – Hamming distance in hash space correlates with original distance but is not perfect.

*Characteristics:* LSH has the attractive property of theoretical guarantees on approximation quality. However, in practice it often needs fairly long codes to reach high recall in high dimensions. Pinecone's analysis highlights that LSH performance can vary widely and is **highly sensitive to the curse of dimensionality** – as dimension grows, you need disproportionately more bits to maintain recall <sup>78</sup> <sup>79</sup>. For example, to achieve ~90% recall on 128-D data, one might need a code of 8192 bits (64× the dimension) <sup>76</sup>, which is a 1024-byte code – *larger* than the original 128-D float vector (512 bytes) in that case! <sup>76</sup> This undermines the benefit of hashing. With shorter codes, LSH can be extremely fast but will miss many true neighbors (low recall) <sup>80</sup>. Thus, Faiss's LSH (`IndexLSH`) is best suited for **low-dimensional vectors or small datasets** where moderate-length codes suffice <sup>79</sup> <sup>81</sup>. It's generally less effective on very high-dimensional embeddings (e.g. 100+ dims) compared to quantization or graph methods.

Memory-wise, an LSH index requires  $nbits/8$  bytes per vector <sup>82</sup> (plus some overhead for storing the buckets). For moderate nbits, this can be quite compact. And search time is sub-linear, as the hash lookup narrows candidates (Faiss's `IndexLSH` essentially does a table scan on binary codes using bit operations, which is fast in C++). But as noted, to get good accuracy one might increase nbits and effectively end up storing a large binary code, negating memory savings <sup>78</sup> <sup>83</sup>. In summary, LSH in Faiss is an option for approximate search, but for high-dimensional data one often finds that other approaches (IVF, PQ, HNSW) achieve a better speed-accuracy trade-off. The Faiss wiki even notes two fundamental drawbacks of traditional LSH: (1) it often requires many hash functions (many bits) to reach good recall, inflating memory usage, and (2) random hashes are not tuned to the data distribution, making them suboptimal in practice <sup>84</sup>. Faiss mitigates (2) by using random rotation (which is data-independent but better than axis-aligned random projections) <sup>74</sup>, but the issue of needing large nbits remains for difficult data. So, use `IndexLSH` for small or low-d scenarios; otherwise, other Faiss indexes may be preferable <sup>85</sup> <sup>81</sup>.

## HNSW (Hierarchical Navigable Small-World Graph) – Graph-Based Search

**HNSW** is a graph-based indexing method that has become one of the top-performing ANN algorithms. Faiss includes HNSW indexes (e.g. `IndexHNSWFlat`, `IndexHNSWPQ`, etc.), which implement the **Hierarchical Navigable Small-World** graph approach. An HNSW index organizes the data points in a layered navigable graph: each vector is a node in a graph and edges connect each node to its nearest neighbors (forming a small-world graph where short paths exist between any two nodes) <sup>86</sup> <sup>87</sup>. The “hierarchical” part refers to additional layers of shortcut connections: HNSW builds multiple graph layers where the top layer is a very coarse graph (few nodes, long links) and lower layers have more nodes and local links <sup>88</sup>. A search in HNSW works by first traversing the top coarse layer (quickly finding a region near the query), then descending layer by layer, each time using the graph connections to move closer to the query in the vector space. By the bottom layer (which contains all points), the algorithm has a good starting point and only explores a *limited neighborhood* of the query, rather than all points <sup>88</sup> <sup>89</sup>. This yields very high efficiency while maintaining high recall.

**Key parameters:** HNSW in Faiss has a few parameters that control its performance <sup>90</sup>:

- *M*: the number of neighbor connections per node in the graph (each point will be connected to M nearest neighbors in the final layer graph). Higher M means each node has more links, which increases recall (and also index size) <sup>91</sup>. Typical M might be 16, 32, 64.
- *efConstruction*: the size of the dynamic list for neighbors during index build (higher means more candidate neighbors are considered when adding a new node, leading to better graph quality at the cost of build time) <sup>92</sup>.



- *efSearch*: the size of the priority queue for neighbors during search (this basically controls search breadth – larger efSearch explores more nodes and improves accuracy, at the cost of query speed)

<sup>92</sup> .

In Faiss, you construct an HNSW index usually by wrapping another index that handles storage/compression. For example, `IndexHNSWFlat(d, M)` gives an HNSW index over an underlying flat vector store (no compression) <sup>93</sup> <sup>94</sup> . If memory is a concern, one could use `IndexHNSWPQ` to store PQ-compressed vectors but still use the HNSW graph for search <sup>95</sup> . Basic usage:

```
M = 32
index_hnsw = faiss.IndexHNSWFlat(d, M)           # HNSW index with M neighbors,
flat storage
index_hnsw.hnsw.efConstruction = 64
index_hnsw.hnsw.efSearch = 32
index_hnsw.add(xb)                               # add vectors (builds the graph)
D, I = index_hnsw.search(xq, k)
```

Here we choose M=32 neighbors. We set `efConstruction` and `efSearch` as desired (they default to M in Faiss if not set). After adding data, Faiss automatically builds the HNSW graph. A search will then explore the graph with the given efSearch and return the k nearest neighbors.

*Performance:* HNSW is known for excellent recall-speed balance. It can often achieve ~95–100% recall while being much faster than brute force <sup>96</sup> <sup>97</sup> . In many benchmarks, HNSW (and its cousin NSG) top the leaderboard for CPU ANN search. Faiss’s HNSW should perform similarly to standalone HNSW libraries (like hnswlib), though Faiss allows combining HNSW with other compressions. The **drawback** of HNSW is memory overhead. The graph links (edges) require storing M neighbors per node, and each neighbor is an integer ID plus possibly a distance. Faiss’s wiki formula for memory is roughly  $d$  floats +  $M * 2 * 4$  bytes per vector for the graph (two 4-byte ints per edge) <sup>98</sup> <sup>99</sup> . For instance, on the SIFT1M dataset (1 million 128-D vectors), an HNSW with M=128 used about 1.6 GB of memory, whereas the raw vectors are ~0.5 GB <sup>99</sup> . So HNSW often uses 2–3× the memory of the raw data (or more if M is large). This is the price for having rapid, greedy graph-based search.

HNSW also has some limitations: it does not support real-time dynamic updates very well (it can add vectors one by one, but deletions are not supported because removing a node would break the graph structure) <sup>100</sup> . And like other graphs, constructing it can be somewhat slow (though still much faster than brute-force distance calculations for large N). But once built, queries are extremely fast – often just a few microseconds for high recall on million-scale data <sup>96</sup> .

In Faiss’s guidance, if you have ample RAM and need very high speed and accuracy, HNSW is recommended as a top choice <sup>101</sup> . It’s “*very fast and accurate*” given enough memory, and by tuning M and efSearch you can push accuracy to near-exact <sup>102</sup> . Pinecone’s summary: “*HNSW – great search-quality, good search-speed, but substantial index sizes*” <sup>103</sup> . Thus, HNSW is ideal for use cases where memory is not the bottleneck but query latency is critical (e.g. real-time search on moderate-scale data). If memory is limited, one might combine HNSW with compression (as Faiss allows) or use other methods.

## Composite Indexes and the Faiss Index Factory

One of Faiss's powerful features is the ability to **combine** multiple techniques into a single index, often called a *composite index*. For example, the IVF and PQ methods described above are frequently combined (IVF+PQ) to get the benefits of both: IVF reduces the search scope, and PQ compresses the vectors for lower memory usage <sup>104</sup> <sup>105</sup>. We can even add a further refinement step to re-rank results using the original vectors or a finer quantization. Other combinations include adding a preprocessing step like PCA or OPQ (optimized PQ) to reduce dimensionality before indexing, or using an HNSW graph as a coarse quantizer for IVF, etc. <sup>106</sup> <sup>107</sup>. Conceptually, a composite index is like a pipeline of transformations and indexes applied in sequence <sup>108</sup> <sup>106</sup>. Faiss provides an **index factory** to easily create such combinations via a simple string notation.

For instance, an index described by the factory string `"IVF256,PQ16+RFlat"` might involve: (a) coarse quantization into 256 clusters (IVF256), (b) within each cluster, compress vectors into 16-byte PQ codes (PQ16), and (c) an optional *Refine* step with a flat index (RFlat) which stores the original vectors for re-ranking the final results <sup>109</sup> <sup>110</sup>. The image above illustrates a composite IVF+PQ index: each vector is assigned to one IVF cell (coarse cluster) and also encoded into subvector codes via PQ (fine quantization), rather than storing the full vector. By stacking these methods, composite indexes can achieve excellent speed and memory trade-offs – but they require careful tuning to avoid excessive accuracy loss <sup>111</sup> <sup>112</sup>.

Faiss's `index_factory` function allows users to create composite indexes from a single descriptor string. This string encodes the sequence of components. For example, instead of manually constructing a quantizer and IVF index as we did earlier, one can do:

```
index = faiss.index_factory(d, "IVF256,Flat")
```

This will internally create a `IndexFlatL2` quantizer and an `IndexIVFFlat` with 256 clusters for dimension  $d$  <sup>113</sup>. The factory can express many combinations. For instance, `"OPQ16_64,IVF4096_HNSW32,PQ16"` would mean: apply OPQ rotation (16 units, output dim 64) -> use IVF with 4096 clusters and HNSW-32 as the assigner -> within each cluster, store PQ16 codes with 16-byte codes. The Faiss wiki and tutorials give guidelines on these strings. As shown in the Meta example, a complex index `"OPQ20_80,IMI2x14,PQ20"` was chosen to fit a 1B-vector dataset in 30GB RAM <sup>114</sup> – this string specifies an OPQ rotation to 80D, a multi-index IVF (IMI) with  $2 \times 14$  bits (i.e.,  $2^{14} \times 2^{14} = 2^{28}$  cells), and PQ20 (20-byte codes) <sup>114</sup>.

The index factory greatly simplifies experimentation: you can swap components by changing the string, rather than writing a lot of code. It ensures the pieces are compatible and configured properly. Under the hood, composite indexes in Faiss are often implemented as one index *wrapping* another. For example, `IndexIVFPQ` is an `IndexIVF` that uses a `IndexPQ` for encoding residuals; `IndexHNSWPQ` is an HNSW index that uses a PQ index for storage, etc. <sup>95</sup>. The factory abstracts that away.

**When to use composite indexes?** Almost always when dealing with large datasets, a composite approach is needed to meet memory or latency constraints. A pure flat index may be too slow, and a pure PQ flat index may still be slow (exhaustive search over PQ codes) or not memory-efficient enough without IVF. So combinations like IVF+PQ are extremely common (this is essentially the approach used in the famous *IVFADC* and *IVFADC+* algorithms <sup>58</sup>). Adding a refinement step (RFlat) to recheck top results can improve

accuracy to nearly exact with little overhead <sup>60</sup>. Using HNSW as a coarse quantizer is a newer idea to speed up assignment for very large *nlist* (as seen in Faiss's guidance for 1M–100M+ vectors) <sup>115</sup>. The general principle is: **coarse quantization (partitioning) to limit search scope, and fine quantization (compression) to save memory**, and optionally refine for accuracy <sup>106</sup> <sup>116</sup>.

Faiss's versatility in mixing these components is a major advantage over some other libraries that may only implement one approach. It does mean more parameters to tune (*nlist*, *nprobe*, *m*, *nbits*, *M*, *ef*, etc.), but also more levers to achieve an optimal solution for a given use-case.

To illustrate index factory equivalence: If we manually create an IVF index and an equivalent one via factory, they give the same results. For example, `faiss.IndexIVFFlat(q, 128, 256)` vs `faiss.index_factory(128, "IVF256,Flat")` produce identical indexes, as confirmed by querying them <sup>117</sup> <sup>118</sup>. The factory simply provides a shorthand. It's especially helpful for composite setups, where writing the code by hand would be error-prone. In practice, one can iterate over different factory strings to find a good configuration, then fine-tune parameters like *nprobe* or *efSearch* for performance.

## Memory Efficiency and Performance Trade-offs

Different Faiss indexes have different memory footprints and speed characteristics. It's crucial to understand these trade-offs when choosing an index:

- **Flat (IndexFlat)** – *Memory*: stores full vectors (4 bytes × *d* each), so memory grows linearly with dataset size. No extra overhead for IDs unless an ID map is added <sup>119</sup>. *Speed*: O(*N*) per query – exact but slow for large *N*. Best for small data or when exactness is a must. Flat indexes require data to fit in RAM (Faiss does not do out-of-core search except via user-managed sharding) <sup>120</sup>.
- **IVF** – *Memory*: adds a small overhead for each vector (e.g., 8 bytes for storing the inverted list index/id) <sup>46</sup>. Overall memory can be only slightly more than flat if using `IVFFlat` (since vectors are stored uncompressed in lists). If combined with compression (IVFPQ, IVFSQ), memory per vector can drop significantly at some accuracy cost. *Speed*: sub-linear per query. The cost is proportional to *nprobe* × (average list length). Using a larger *nlist* (more clusters) means shorter lists (faster per probe) but more memory for centroids and possibly more overhead to search many clusters if *nprobe* is large. There is a sweet spot – e.g., *nlist* around  $\sqrt{N}$  often works well <sup>52</sup>. Also, parallelization: Faiss can search different lists in parallel threads or on GPU streaming multiprocessors, improving throughput. **Key tuning**: *nlist* and *nprobe*. High *nlist* gives finer partition (fewer candidates per list) but might need higher *nprobe* to not miss results. Faiss guidelines propose values of *nlist* based on dataset size (e.g. 100K → *nlist*= $\sqrt{N}$  ≈ 316, 1M → a few thousand, 10M → up to 65k etc.) <sup>115</sup>. Memory for centroids themselves is trivial (*nlist* × *d* floats).
- **PQ** – *Memory*: compresses vectors to *M* × *nbits* bits. For example, PQ with *m*=16, *nbits*=8 yields 16 bytes/vector <sup>121</sup> <sup>67</sup>. That's a 32× compression for a 128-D float vector. So memory is drastically reduced (plus possibly a small codebook storage, e.g. 16×256×128 bits for codebooks, which is negligible). *Speed*: computing approximate distances is very fast using precomputed tables (Faiss uses an LUT of size [*m* × (2<sup>*nbits*</sup>)] floats for each query, then a code's distance is *m* table lookups and additions). So distance calculation cost is *m* per vector instead of *d*. E.g., *d*=128, *m*=16: we do 16 table lookups vs 128 multiplications for a full distance. This can significantly accelerate search,

especially with CPU caching (the LUT often fits in L1). However, PQ introduces *quantization error*, so recall is lower than an uncompressed method. We can mitigate by increasing code size (larger  $m$  or  $nbits$ ) at cost of memory and search time, or by adding a re-ranking of top results using original vectors (which needs storing or retrieving originals). **Note:** When combined with IVF (IVFPQ), memory per vector is just the PQ code + id (e.g. 16 bytes + 8 bytes = 24 bytes, vs 512 bytes originally), enabling truly billion-scale indices in tens of GB <sup>20</sup>. On GPU, Faiss has optimized code for PQ search as well, though PQ search on GPU may be limited by memory bandwidth (since many small code distances must be fetched).

- **LSH** – *Memory*: each vector  $\rightarrow$   $nbits$  bits (and typically an ID or pointer in a bucket). For example, 128-D vector to 128-bit code is 16 bytes (much smaller than 512 bytes original). But as mentioned, to get good accuracy often  $nbits$  must be large. Pinecone observed that for  $d=128$ , achieving 90% recall needed 8192 bits = **1024 bytes** per vector <sup>76</sup> – double the memory of the raw vector! That defeats the purpose. In low dimensions or low recall needs, one might keep  $nbits$  small (e.g.  $2 \times d$  bits) to have a compact code <sup>76</sup>. *Speed*: LSH search essentially does a hash table lookup or scans within a few Hamming radius of the query code. Faiss's LSH as implemented performs an efficient binary distance comparison. It can be very fast for moderate  $nbits$ . But if  $nbits$  is huge, the Hamming comparison cost (which is proportional to  $nbits/word\_size$ ) also grows and the advantage wanes <sup>83</sup>. Moreover, if many items land in the same bucket (collision), it has to scan those. For good performance, one might use multiple hash tables or multi-probe LSH, which Faiss's `IndexLSH` does not explicitly implement (it effectively corresponds to one hash function with  $nbits$  bits). In summary, LSH in Faiss is memory-efficient and fast only in certain regimes (small/medium dimension data with moderately large  $nbits$ ). It tends not to be the first choice for high-precision search on complex embeddings, but can still be useful for quick-and-dirty approximate search when memory is ample and suboptimal accuracy is tolerable <sup>79</sup>.
- **HNSW** – *Memory*: as discussed, requires storing neighbor lists. For each vector, approximately  $M$  links are stored, each link an integer (4 bytes) plus possibly a distance. Faiss's HNSWFlat index uses  $\sim(4d + M8)$  bytes per vector (the  $4d$  for the float vector,  $M24$  for neighbor links) <sup>98</sup> <sup>90</sup>. With  $M=32$  and  $d=128$ , that's  $512 + 256 = 768$  bytes per vector ( $1.5 \times$  the original). With  $M=64$ ,  $\sim 1024$  bytes ( $2 \times$  original). So memory overhead is significant, especially since one usually doesn't compress vectors in an HNSW index (though Faiss's HNSWPQ can compress to reduce the  $4d$  part). In addition, HNSW has multiple layers but the top layers contain only a few points, so their overhead is minor. *Speed*: HNSW is extremely fast at query time. It often explores on the order of  $O(M * \log N)$  points or similar – the algorithmic complexity is sub-linear. With typical `efSearch` settings, the actual number of distance evaluations per query might be, say, a few hundred for a million-vector index to achieve  $>90\%$  recall, which is far less than IVF might need if  $nprobe * list\_length$  is larger. This is why HNSW is praised for speed/accuracy. The `efSearch` parameter can be tuned: larger `efSearch`  $\Rightarrow$  more thorough search (higher accuracy, more comparisons, slower). Even at high recall, HNSW can beat or match other methods. The main drawback is memory, as noted, and the slower construction time (though still manageable: constructing an HNSW index is roughly  $O(N \log N)$  with a small constant, possibly a few hours for 1B vectors on a single machine – but usually you'd use IVF+PQ for 1B to save memory). If memory allows, HNSW is a great default ANN method; Faiss's integration means you can also combine it (e.g. use an HNSW coarse quantizer for IVF to speed up clustering for very large  $N$ ) <sup>122</sup>.
- **Composite (IVF+PQ, HNSW+PQ, etc.)** – *Memory*: Combining methods usually aims to reduce memory while keeping speed. IVF+PQ, for example, has memory = PQ code (bytes per vector) + id (8

bytes) + centroid storage ( $nlist * d$ ) + overhead for lists. For large  $N$  with reasonably large  $nlist$ , the centroid storage is negligible. So an IVF+PQ might get you down to 16–32 bytes per vector (for  $m=16$  or  $m=32$  PQ) plus ~8 bytes id = ~24–40 bytes total. This is extremely memory-efficient (we're talking 40 GB for a billion vectors). HNSW+PQ would have vector compressed plus graph links overhead – you save on the vector storage but still have links. For instance, IndexHNSWPQ with  $m=16$  might be ~16 bytes (code) + 8 bytes links \*  $M$ . If  $M=32$ , that's  $16+256=272$  bytes, a big improvement over HNSWFlat's 768 bytes (the PQ compression helps a lot). However, compressing can reduce accuracy if done aggressively. *Speed*: Combining strategies can sometimes slow queries (because you have multiple layers of work: e.g. in IVF+PQ, you first compute centroid distances, then code distances). But these steps are usually efficient and parallelizable. For IVF+PQ, the bottleneck might be scanning the list: even though codes are smaller, you still iterate over roughly  $(N/nlist * nprobe)$  codes. If PQ codes are small, many can fit in cache at once, speeding things up. HNSW+PQ might slow slightly vs HNSWFlat because computing PQ distances is approximate (possibly requiring more exploration for the same recall). In practice, one often chooses composite indexes to meet memory limits, accepting some speed loss or accuracy loss. The Faiss documentation recommends specific composites depending on how tight memory is <sup>123</sup> <sup>124</sup>. For example, if memory is “quite important”, they suggest an OPQ + IVF + PQ with 4-bit codes (very high compression) <sup>123</sup>; if memory is “somewhat important”, maybe just IVF without compression (Flat in lists) <sup>125</sup>; if memory is not a concern, use HNSW or IVF with re-rank which prioritizes speed/accuracy <sup>101</sup> <sup>126</sup>.

In summary, Faiss offers the building blocks to tailor an index balancing memory, speed, and accuracy. A flat index gives 100% accuracy but at full memory and slower speed. At the other extreme, an aggressively compressed index (like IVF+PQ with small codes) uses very little memory and can be very fast, but may sacrifice some accuracy (measured by recall). The optimal choice depends on the application's requirements – some use-cases can tolerate slightly approximate results if it means sub-second query times on huge data, others cannot. Faiss's strength is that you can start with a simple approach and progressively add techniques to meet scaling needs, all within one library.

## Real-World Applications of Faiss

Faiss has become a go-to solution for similarity search in numerous domains, thanks to its efficiency and flexibility. Some prominent use cases include:

- **Recommendation Systems:** Faiss can quickly find similar items or users by comparing embedding vectors. For example, an e-commerce site can embed users' browsing history or item descriptions into vectors, then use Faiss to find related products to recommend. Because Faiss handles huge datasets, it can power real-time recommendations even with millions of products. DataCamp notes that Faiss is a “*game-changer*” for recommendation: given high-dimensional vectors of user behavior or item features, Faiss nearest-neighbor queries can identify similar products or content almost instantly <sup>127</sup>. This personalized search boosts engagement and sales. Companies like Facebook have used Faiss for recommending posts or ads by finding similar user embedding vectors (Faiss was originally built to handle very large recommendation candidate sets).
- **Image and Video Search:** Searching by image similarity is a natural application. Deep learning models can embed images or video frames into vectors that capture visual semantics. Faiss indexes these embeddings to enable queries like “find images similar to this one” or deduplicate near-duplicates in a dataset. Its scalability is crucial for platforms with billions of images (e.g. social media

or stock photo libraries). As an example, one could take all images in a photo app, compute CNN descriptors for each, index them with Faiss, and then given a query image, retrieve all images of the same object or scene by nearest neighbor search <sup>128</sup> <sup>129</sup>. Meta's blog described how Faiss allows quickly searching multimedia documents by similarity – something traditional text-based search can't do <sup>15</sup>. It's also used in video analytics: indexing clip embeddings to find similar scenes. The **visual search** capability is employed in products like content moderation (find near-duplicate images), landmark recognition, or reverse image search engines. According to DataCamp, Faiss *"powers search engines that retrieve visually similar images or videos by indexing high-dimensional vectors from multimedia content"*, enabling use-cases like finding all images of a specific landmark in a large photo collection <sup>130</sup>.

- **Semantic Text Search (NLP):** With the rise of transformer models, textual data (documents, sentences, queries) can be embedded into vector representations such that semantic similarity in meaning corresponds to vector proximity. Faiss is widely used to build semantic search engines that can, for instance, find relevant FAQs given a customer question, or perform neural information retrieval for QA systems. Hugging Face's libraries integrate Faiss to index embeddings of sentences or paragraphs so that given a new query embedding, one can fetch the most semantically similar texts from a corpus <sup>131</sup> <sup>132</sup>. This is key in applications like open-domain question answering (retrieving passages related to a question), document clustering, or improving search in e-commerce by matching queries with product descriptions semantically. Faiss's support for *inner product* search is used for finding maximum cosine similarity matches, a common need in NLP (after normalizing embeddings) <sup>133</sup>. For example, in the HuggingFace LLM course, Faiss is used to build a knowledge base of documentation embeddings and then answer queries by finding the nearest neighbor docs to the query embedding (semantic search) – dramatically improving over keyword search.
- **Anomaly and Outlier Detection:** In anomaly detection tasks, one often embeds data points and then looks for points that have no close neighbors. Faiss can accelerate this by enabling k-NN queries for each point or using range search (find all neighbors within a certain radius). If a point has an unusually large distance to its nearest neighbor, it's a potential anomaly. For example, in fraud detection, financial transactions can be vectorized by their features, and Faiss can quickly find the distance of each transaction to its neighbors; those that are far apart (in embedding space) may indicate rare, potentially fraudulent events <sup>134</sup>. Similarly, for system monitoring or cybersecurity, Faiss could be used to find events/logs that don't cluster with others. The advantage is that Faiss can handle the large volumes of data typically needed to establish "normal" neighborhoods, and still flag outliers in near-real-time. DataCamp mentions how Faiss can identify outliers by finding points that deviate from their nearest neighbors, e.g. flagging an unusual transaction by vector similarity <sup>134</sup>.
- **Clustering and Indexing in ML Pipelines:** Faiss isn't only for search – it also includes routines for clustering (k-means) and can be used as a general vector database inside ML pipelines. For example, one might use Faiss's K-Means implementation to cluster embeddings for unsupervised learning tasks (Faiss provides efficient GPU k-means as well). Some use it for iterative algorithms like *active learning*: Faiss can quickly compute nearest neighbors which may be needed in diversity sampling or core-set selection algorithms. In large-scale ML systems, Faiss often serves as a module to retrieve candidate items (candidate generation stage) before a more expensive ranking stage refines the results <sup>135</sup>. For instance, in a search engine or recommender, Faiss might shortlist a few hundred nearest neighbors out of millions (using an ANN index), and then a second model or re-ranker

examines those in detail. This two-stage approach (ANN retrieval + precise ranking) is common in industry, and Faiss is frequently the ANN engine in the first stage <sup>135</sup>.

Overall, Faiss's impact is seen in any application requiring *fast similarity matching on large datasets*. It underpins features in social media (finding similar posts, images, friends-of-friends suggestions), e-commerce (searching similar products, personalized recommendations), AI research (retrieving examples for few-shot learning or nearest-neighbor classification), and even specialized domains like genomics (where DNA sequences are embedded and nearest neighbor search can find similar gene sequences) or drug discovery (chemical molecule embeddings). Its versatility with distance metrics means it can be used for cosine similarity (common in NLP) or Euclidean (common in vision) or even dot-product for maximal inner product search in certain recommendation algorithms <sup>5</sup>. And importantly, it's all in a single library that can be optimized on CPU or GPU. Many companies and open-source projects have adopted Faiss; for example, Facebook uses it in various services, Spotify has used Faiss (and Annoy) for music recommendation, and the popular **Milvus** vector database can use Faiss indexes under the hood for its search engine. Faiss truly enables large-scale *vector search* – an essential component of modern AI systems.

## Comparison with Other Vector Search Libraries

Faiss is one of several libraries in the ANN search space. Others include **Annoy** (by Spotify), **HNSWlib** (by Y. Malkov), **ScaNN** (by Google), and a number of specialized or newer tools. Each has its strengths, and the choice can depend on the use case. Here we compare Faiss with a few popular alternatives:

- **Faiss vs. Annoy:** **Annoy** (Approximate Nearest Neighbors Oh Yeah) is a C++/Python library developed at Spotify, which builds multiple randomized KD-trees (or more accurately, *random projection trees*) for ANN search. Annoy is designed to be very simple and memory-efficient, with the ability to memory-map indexes from disk. It's particularly useful for **static datasets** – once you build the Annoy index, you cannot easily update it (it's immutable; changes require rebuilding) <sup>136</sup> <sup>137</sup>. The core algorithm in Annoy is to recursively partition the space with random hyperplanes (random projections), building a forest of trees. Querying involves traversing these trees to get candidate neighbors. **Comparative strengths:** Annoy is praised for its ease of use and low memory overhead in certain cases. It allows using disk instead of RAM (good for very large data if you can tolerate slower access). It's also very fast for *building* an index and querying when data is static. However, Annoy's approach trades accuracy for speed and simplicity. It doesn't always reach the high recall that Faiss can, unless many trees are used (which increases memory). Also, as noted, it's read-only once built – not suitable for dynamic environments. In contrast, Faiss supports incremental additions to indices (and in some cases removal, except HNSW) and provides a wider range of algorithms <sup>138</sup> <sup>139</sup>. Faiss can perform both exact and approximate search, and can leverage GPUs, which Annoy cannot. A Zilliz blog puts it succinctly: Annoy's focus on speed and memory makes it great for static, read-heavy workloads, but that comes at cost of flexibility (immutable index), whereas Faiss's broader range of algorithms allows adjusting speed vs accuracy and handling dynamic data better <sup>140</sup> <sup>141</sup>. Also, Faiss tends to outperform Annoy on very large or high-dimensional datasets in accuracy at given speed, especially when GPU is used or when advanced indexes (PQ, HNSW) are needed <sup>142</sup> <sup>143</sup>. Annoy might be a good choice for small-memory environments or when disk persistence of the index is needed, but for most high-scale applications Faiss offers more power (at the cost of a slightly steeper learning curve and setup).

- **Faiss vs. HNSWlib:** **HNSWlib** is a standalone C++ library (with Python bindings) that implements the HNSW algorithm (graph-based search) very efficiently. It's essentially an extraction of the original HNSW author's code. Faiss actually incorporated HNSW based on that work, so in many ways Faiss's HNSW index and HNSWlib will behave similarly given the same parameters. Both will produce high recall results with low latency, and both require significant memory for the graph. One difference: HNSWlib focuses solely on HNSW, whereas Faiss offers HNSW as one option among many and can combine it with others (like PQ). If your use case specifically demands the absolute fastest CPU ANN search and memory is secondary, HNSWlib is an excellent choice, and Faiss's version is only marginally different in performance. HNSWlib might be a bit ahead in that it's highly tuned for that one task and has features like saving/loading index to disk easily. On the other hand, Faiss provides a more uniform interface and the ability to switch to other methods or use GPUs. In terms of the algorithm, both will give you the same complexity characteristics. Pinecone's benchmarking noted "*HNSW consistently tops out as the highest-performing index*" in quality vs speed <sup>96</sup>. Indeed, HNSW (either via Faiss or HNSWlib) often outperforms other approaches like Annoy or LSH for high recall scenarios. So the question is more about integration: if you are already using Faiss, using Faiss's HNSW is convenient. If you only need HNSW and perhaps a simpler API, HNSWlib is fine. Faiss's HNSW is not available on GPU (currently HNSW is CPU-only, though efforts like Nvidia's cuGraph may someday change that), whereas Faiss's other methods (IVF, PQ) can be GPU. So if you need GPU, Faiss would push you toward a GPU-friendly index instead (or use CPU HNSW). Both Faiss and HNSWlib support multi-threading for adds and queries (Faiss's search can be multi-threaded per query or one query per thread easily, similar to HNSWlib's multithreading support). In summary, Faiss vs HNSWlib is less "which is better" than "HNSW algorithm vs others": Faiss gives you HNSW plus alternatives; HNSWlib only gives HNSW but might be slightly more lightweight if that's all you need.

- **Faiss vs. ScaNN:** **ScaNN (Scalable Nearest Neighbors)** is Google's ANN library, released in 2020, focusing on deep-learning embeddings. ScaNN's approach is somewhat like an advanced IVF+PQ with learned quantization – it introduced *anisotropic vector quantization*, which applies a learned transformation weighting dimensions to minimize quantization error <sup>144</sup>. ScaNN also can re-rank with exact inner product of top results. It's particularly tuned for inner product search (e.g. retrieving high-dot-product items). The goal of ScaNN is to optimize the accuracy-speed frontier, often achieving very high accuracy at given recall because of its fine-grained quantization scheme. **Comparisons:** Faiss and ScaNN both target high performance, but Faiss is more general and feature-rich (supporting many modes, GPU, etc.), whereas ScaNN might edge out Faiss in certain benchmark scenarios for recall vs latency on CPU (as per Google's paper). A Zilliz comparison notes: Faiss emphasizes *scalability and compression* (lots of methods including PQ, optimized for very large datasets, with GPU acceleration), while ScaNN "*optimizes accuracy in high-dimensional spaces*" with its anisotropic PQ and is often very effective when accuracy requirements are strict <sup>145</sup> <sup>144</sup>. ScaNN currently is CPU-only (although Google possibly uses internal GPU versions, the open source is CPU with some SIMD). If you need maximum accuracy and are working with, say, TensorFlow or Google's stack, ScaNN is a good choice. But it lacks the breadth: Faiss can do cosine, Euclidean, etc., on GPU, with many indexing options. It's reported that ScaNN can outperform Faiss in certain recall@K vs latency benchmarks on CPU for inner product search (especially when high recall is needed, ScaNN's anisotropic quantization helps preserve accuracy) <sup>146</sup> <sup>147</sup>. However, Faiss's recent developments (like support for re-ranking and its own optimized PQ/OPQ) narrow that gap. Also, Faiss integrated a version of ScaNN's approach for symmetric distances in some CPU code paths (the FAISS authors note that ScaNN's ideas are partly implemented in Faiss for certain cases) <sup>148</sup>. Ultimately, if one were



building a production system, Faiss offers more support (and a bigger community). If one is trying to squeeze the last bit of accuracy on CPU and doesn't mind a TensorFlow dependency, ScaNN might be tested. A Milvus article sums it up: *"FAISS, HNSW, and ScaNN solve ANN with distinct strategies: FAISS emphasizes scalability & compression, HNSW prioritizes fast graph traversal, and ScaNN optimizes accuracy. Choosing depends on dataset size, dimensionality, latency, hardware."* <sup>145</sup>.

In addition to those: **nmslib** (which implements older graph methods and others) can be faster for small data but is less maintained now. **Annoy** we covered. **HNSWlib** we covered. There are also **vector databases** like Milvus, Pinecone, Weaviate, etc., which actually often use Faiss or HNSWlib under the hood while adding distributed, persistent storage and auto-scaling. Those are higher-level systems beyond just libraries. If one's need is an embedded library to add ANN to an application, Faiss is often the top choice given its performance and flexibility. Annoy could be chosen for simplicity or static disk-based scenarios; HNSWlib for when you specifically want the best graph method with minimal fuss; ScaNN for certain high-accuracy retrieval tasks. But Faiss's combination of support for *exact search, various ANN methods, GPU, huge datasets, and an active development (Facebook regularly updates it)* makes it arguably the most **versatile**. It's telling that many vector DBs let you choose Faiss as the indexing backend.

To illustrate, Chloe Williams (Zilliz) writes: *"Annoy uses random projection trees... great for speed on static data, but index is immutable. Faiss uses a broader range of algorithms... adjustable speed-accuracy and supports incremental updates, plus GPU – giving Faiss an edge in large-scale, real-time applications where flexibility is essential."* <sup>140</sup> <sup>142</sup>. And in distributed or very large settings, Faiss can be used across multiple machines (though Faiss itself isn't distributed, one can shard data and build Faiss indexes per shard; some vector DBs orchestrate this). Faiss also integrates well with deep learning pipelines (Facebook has a Faiss/PyTorch integration for similarity tasks, and as mentioned HuggingFace uses Faiss for semantic search).

In summary, Faiss stands out by offering **power and flexibility**, while alternatives might win in specific niches (Annoy for static disk-based simplicity, HNSWlib for pure graph search, ScaNN for slight accuracy gains on CPU inner product search). Often, one might prototype with a simple library like Annoy but then move to Faiss for production to gain efficiency at scale <sup>139</sup>. Or use Faiss for the heavy lifting and an external system for distribution. The good news is these tools are not mutually exclusive – they are all open-source, and one can choose the right tool for each job (some applications even use multiple: e.g. Annoy for one module, Faiss for another, though usually standardizing on one is easier).

## Best Practices for Choosing and Tuning Faiss Indexes

Selecting the right Faiss index and tuning its parameters can greatly impact performance. Here are some **best practices and guidelines** based on Faiss documentation and practical experience:

- **If exact results are required or dataset is small:** Use a **Flat index** (`IndexFlatL2` or `IndexFlatIP`). This gives 100% accuracy. It's suitable when N is small (e.g., < 100k) or when you only perform a few searches (so building a complex index isn't worth it) <sup>34</sup> <sup>149</sup>. Flat search can also serve as a baseline to measure other indexes' accuracy. If you need to assign custom IDs to vectors, wrap it with `IndexIDMap` since flat indexes by themselves only support implicit sequential IDs <sup>119</sup>.
- **If the dataset is moderately large (up to a few million) and memory is not a big concern:** Consider an **HNSW index** for fast, high-accuracy search. Faiss recommends HNSW when you have

“lots of RAM or the dataset is small” because it’s very fast and accurate <sup>101</sup>. Choose  $M$  (connections per node) between 16 and 64 (higher  $M$  = higher recall but more memory). Use `efSearch` to tune recall at query time (start with `efSearch` =  $M$ , and increase if you need better accuracy) <sup>91</sup> <sup>92</sup>. Note HNSW doesn’t support deleting vectors or `add_with_ids` (only sequential add) <sup>150</sup>. If you need to map your own IDs, you can use `IndexIDMap` on top, but removal still isn’t possible. Use HNSW when you can afford  $\sim 1.5\times$  memory of your data for the index and want query speed in the microsecond to sub-millisecond range. It’s often the best choice for real-time queries on million-scale data if you have the RAM.

- **If memory is a concern (large datasets):** Faiss provides compressed indexes. A common strategy: use an **IVF index with product quantization (IVFPQ)**. This drastically lowers memory per vector at some cost to accuracy. When tuning IVF, a key question is how many clusters ( $nlist$ ). Empirical rule: set  $nlist$  proportional to  $\sqrt{N}$  – for example, if  $N=10$  million,  $\sqrt{N} \approx 3162$ , you might use  $nlist = 4096$  or  $8192$  <sup>52</sup>. The Faiss wiki suggests:

- Below 1 million vectors:  $nlist \sim 4 \times \sqrt{N}$  to  $16 \times \sqrt{N}$  <sup>52</sup>.
- 1M – 10M: consider using up to 65k cells (e.g. 65536, and possibly an HNSW quantizer to assign vectors faster) <sup>151</sup>.
- 10M – 100M: up to 262144 cells ( $2^{18}$ ) <sup>152</sup>.
- 100M – 1B: up to 1,048,576 cells ( $2^{20}$ ) <sup>153</sup>.

Essentially, as data grows, increase  $nlist$ . But very large  $nlist$  will increase index build time and memory for centroids, and may require using an efficient assignment (Faiss suggests IVF with HNSW quantizer for 1M+ scale) <sup>122</sup>. Always *train* IVF on a representative sample of your data (e.g., 50k–100k vectors or more, depending on  $nlist$ ). After building an IVFPQ, tune  $nprobe$ : start with a small  $nprobe$  (e.g. 1% of  $nlist$ , or an absolute like 4 or 8) and increase until recall is acceptable.  $nprobe$  directly trades query speed for accuracy – higher  $nprobe$  = more lists searched = slower but more accurate <sup>38</sup> <sup>154</sup>. For example, if  $nlist=1024$ ,  $nprobe=16$  might be a good starting point (search 16 out of 1024 clusters per query, i.e., 1.6% of the data on average).

- **Use OPQ (optimized PQ) if using PQ:** OPQ is a linear transform that makes vectors more amenable to PQ compression <sup>155</sup>. It can significantly boost accuracy for a given code size (for example, OPQ might halve the error rate of PQ in some cases). Faiss’s index factory allows adding `OPQ` before a PQ. If memory is extremely tight, you might even do an OPQ + PQ without IVF (just compress everything), but usually IVF+PQ+OPQ is best. For example, Faiss suggests for strong compression: an index string like `"OPQ16_64, IVF4096, PQ16"` for large data – meaning reduce dimension to 64 with OPQ16 (which also rotates data), then IVF4096, then PQ16 with 4 bits each (i.e.,  $16 * 4 = 64$  bits per vector) <sup>123</sup> <sup>124</sup>. That would be very compact (8 bytes per vector + overhead). If that’s too lossy, use 8 bits per subquantizer (PQx8) for 16 bytes per vector.

- **If only moderate compression is needed and speed is crucial:** One approach Faiss mentions is IVF with *Flat* (no compression) but using fewer lists. For instance, `"IVF16384, Flat"` on a 10M dataset – each vector stored fully, but because we partitioned into 16k lists, we only search some of them, saving time. Memory overhead is small (just 8 bytes per vector extra) <sup>46</sup>. This might be ideal if you have enough RAM for full vectors and want better accuracy than PQ. You can also do a hybrid: IVF with a *Scalar Quantizer* (SQ) which compresses each component to say 8 bits (IndexIVFScalarQuantizer in Faiss, `"IVF..., SQ8"`) <sup>156</sup>. SQ8 would halve memory (each component 1

byte instead of 4 bytes float). It's less sophisticated than PQ but simpler – good if moderate compression is fine.

- **Tuning HNSW parameters:** As mentioned,  $M$  controls memory/accuracy. Common values are 16, 32 for high-dimensional data; increasing  $M$  beyond 64 yields diminishing returns in accuracy vs a lot more memory <sup>102</sup>. `efConstruction` can be set higher (e.g.  $2 \times M$  or more) to improve graph quality; it only affects build time and index quality. `efSearch` can be tuned per query or set globally – for example, you might set `efSearch = 2 \times k` or  $4 \times k$  (if you want very high recall). Higher `efSearch` increases query time roughly linearly. Faiss guideline: if you have the RAM, HNSW is usually the best option for pure speed <sup>101</sup>. So use it if you can accept the memory overhead. If not, IVF+PQ might be next.
- **Few queries scenario:** If you will perform very few searches (like  $< 1000$  queries in total, ad hoc analysis), the Faiss wiki suggests it might not be worth building an elaborate index because the indexing time won't be amortized <sup>34</sup>. In such cases, just use a flat index (brute force) or at most a quick clustering. Building an IVF or HNSW can take time (clustering or graph construction), which is wasted if you hardly query the index. So consider the break-even: for instance, building a PQ or HNSW might take minutes to hours for huge data, which is only worth it if you will run many queries or need fast repetitive querying.
- **Use training sets wisely:** For indexes that require `train()` (IVF, PQ, OPQ, etc.), always use a large, representative sample of your data to train. If your data has outliers or multiple distributions, ensure the training sample covers that. Under-training (like using too few vectors for k-means) can degrade performance significantly. Faiss typically requires at least 30–100 times `nlist` vectors for training k-means <sup>151</sup> (the wiki suggests 30K to 256K vectors for training for typical `nlist` values) <sup>157</sup>. For PQ, you want enough samples to fill the subquantizer codebooks reasonably (e.g. at least a few hundred times the codebook size if possible).
- **Monitoring accuracy:** After building an ANN index, measure its accuracy (recall@1 or @10) against brute force for a test subset <sup>17</sup> <sup>158</sup>. This will tell you if your parameters are tuned well. If recall is too low, increase `nprobe` (for IVF) or `efSearch` (for HNSW) or code size (for PQ) or `nbits` (for LSH) until you hit an acceptable point. Often a small increase in `nprobe` or `ef` can boost recall a lot with moderate cost.
- **Parallelizing and GPU:** If using CPU, Faiss allows parallel searching by multiple threads. You can call `index.search` in multiple threads (it's thread-safe for searching) or use Faiss's built-in `omp_set_num_threads` to control internal parallelism. For adding/training, some indexes can also be parallelized (IVF training uses multiple threads in kmeans; HNSW add is single-threaded by default in Faiss, though one could add in parallel partitions). If you have GPUs, use the GPU index for heavy searches – e.g., `faiss.index_cpu_to_gpu()` to transfer a CPU index to GPU <sup>159</sup>. GPU indices shine especially for brute-force (Flat) and large PQ computations, where massive parallelism helps <sup>24</sup>. There is overhead transferring queries to GPU, so batch queries if possible to amortize that.
- **Index maintenance:** If you need to update the index frequently (add or remove vectors), some indexes are better: Flat indexes and IVFs can add (and IVF lists can remove via `remove_ids`, though this can fragment lists), HNSW can add but not delete, Annoy cannot update at all. If your

data is dynamic, you may consider rebuilding periodically or using an ID map to mask “deleted” entries if the index type doesn’t support removal.

- **Metric considerations:** If you want cosine similarity, remember to normalize vectors and use `IndexFlatIP` or corresponding IP index (since cosine similarity is monotonic to dot product for normalized vectors) <sup>28</sup>. Faiss’s factory defaults to L2 for “Flat” unless you specify otherwise. You can specify `metric=METRIC_INNER_PRODUCT` for IVF too if doing MIPS. The best ANN method might differ slightly: for dot-product (like recommendation scores), sometimes one uses *metric transformations* (like convert to Euclidean in one higher dimension, or use special tricks), but Faiss can handle inner product natively.

In conclusion, choosing a Faiss index is about balancing requirements: - Use **Flat** for absolute accuracy or very small data. - Use **HNSW** for fastest queries when memory is sufficient (tune M, ef). - Use **IVF** (with or without PQ) for large data to reduce comparisons; tune nlist, nprobe. - Add **PQ/SQ/OPQ** for memory saving if needed, accepting some accuracy loss; tune code size. - Use **LSH** only for specific low-dim or fast-but-rough needs. - Consider **composite** indexes (PCA/OPQ + IVF + PQ + recheck) for extreme cases (like billion-scale or very high compression needs).

Faiss’s own wiki provides a decision flowchart: start by asking “do I need exact? is RAM limiting? how big is N?” and so on, which essentially echoes the points above <sup>160</sup> <sup>161</sup>. By following those guidelines and iteratively tuning parameters while measuring recall and speed, one can arrive at an index that meets the application’s latency and accuracy targets. The flexibility of Faiss means there’s often a solution for most operating points in the speed-accuracy-memory space, and with careful tuning, you can deploy similarity search that is both **efficient** and *effective* for your specific problem.

## References:

1. Johnson, Douze, Jégou. “Billion-scale similarity search with GPUs.” arXiv preprint arXiv:1702.08734 (2017). [Faiss announcement and technical details] <sup>162</sup> <sup>163</sup>
2. Faiss Wiki – **Overview of indexes and parameters** (Facebook AI Research) <sup>133</sup> <sup>101</sup>
3. Pinecone Engineering – “Introduction to Faiss” and “Nearest Neighbor Indexes: Flat, LSH, HNSW, IVF” (2023) <sup>164</sup> <sup>97</sup>
4. DataCamp Blog – “What Is Faiss?” (2024) <sup>16</sup> <sup>165</sup>
5. Zilliz (Milvus) Blog – “Annoy vs Faiss”, “Faiss vs ScaNN vs HNSWlib” comparisons (2024) <sup>166</sup> <sup>26</sup>
6. Meta AI – “Faiss: A library for efficient similarity search” (Engineering at Meta blog, 2017) <sup>167</sup> <sup>168</sup>

---

<sup>1</sup> <sup>2</sup> <sup>3</sup> <sup>4</sup> <sup>5</sup> <sup>6</sup> <sup>8</sup> <sup>9</sup> <sup>10</sup> <sup>11</sup> <sup>12</sup> <sup>13</sup> <sup>14</sup> <sup>15</sup> <sup>17</sup> <sup>18</sup> <sup>19</sup> <sup>20</sup> <sup>23</sup> <sup>25</sup> <sup>29</sup> <sup>32</sup> <sup>48</sup> <sup>114</sup> <sup>120</sup> <sup>128</sup> <sup>129</sup> <sup>158</sup> <sup>159</sup>  
<sup>162</sup> <sup>163</sup> <sup>167</sup> <sup>168</sup> **Faiss: A library for efficient similarity search - Engineering at Meta**  
<https://engineering.fb.com/2017/03/29/data-infrastructure/faiss-a-library-for-efficient-similarity-search/>

<sup>7</sup> <sup>26</sup> <sup>135</sup> <sup>144</sup> <sup>145</sup> <sup>146</sup> <sup>147</sup> **What is the role of FAISS, HNSW, and ScaNN in AI databases?**  
<https://milvus.io/ai-quick-reference/what-is-the-role-of-faiss-hnsw-and-scann-in-ai-databases>

<sup>16</sup> <sup>21</sup> <sup>22</sup> <sup>24</sup> <sup>27</sup> <sup>54</sup> <sup>71</sup> <sup>72</sup> <sup>127</sup> <sup>130</sup> <sup>134</sup> <sup>155</sup> <sup>165</sup> **What Is Faiss (Facebook AI Similarity Search)? | DataCamp**  
<https://www.datacamp.com/blog/faiss-facebook-ai-similarity-search>

28 31 35 36 37 38 39 40 41 42 43 45 46 47 51 52 53 55 56 57 58 59 60 61 62 63 64 65 66 67

68 69 70 73 74 75 77 82 84 90 91 92 93 94 95 98 100 121 133 154 156 157 **Faiss indexes** ·

**facebookresearch/faiss Wiki · GitHub**

<https://github.com/facebookresearch/faiss/wiki/Faiss-indexes>

30 44 49 50 76 78 79 80 81 83 85 86 87 88 89 96 97 99 103 164 **Nearest Neighbor Indexes for**

**Similarity Search | Pinecone**

<https://www.pinecone.io/learn/series/faiss/vector-indexes/>

33 34 101 102 115 119 122 123 124 125 126 149 150 151 152 153 160 161 **Guidelines to choose an index** ·

**facebookresearch/faiss Wiki · GitHub**

<https://github.com/facebookresearch/faiss/wiki/Guidelines-to-choose-an-index>

104 105 106 107 108 109 110 111 112 113 116 117 118 **Facebook AI and the Index Factory | Pinecone**

<https://www.pinecone.io/learn/series/faiss/composite-indexes/>

131 **Understanding Semantic Search and FAISS for Efficient Similarity ...**

<https://python.plainenglish.io/understanding-semantic-search-and-faiss-for-efficient-similarity-search-with-python-374dc7f2be84>

132 **Semantic search with FAISS - Hugging Face LLM Course**

<https://huggingface.co/learn/llm-course/en/chapter5/6>

136 137 138 139 140 141 142 143 166 **Annoy vs Faiss on Vector Search - Zilliz blog**

<https://zilliz.com/blog/annoy-vs-faiss-choosing-the-right-tool-for-vector-search>

148 **Billion-Scale Approximate Nearest Neighbor Search [pdf]**

<https://news.ycombinator.com/item?id=35814381>