



Vector Database Platform Comparison

This report examines Milvus, OpenSearch, Qdrant, txtai, Weaviate, and Typesense across multiple dimensions to guide production usage. For each system we discuss **features**, **data types**, **architecture**, **search modes**, **ML/NLP integration**, **ecosystem**, **languages**, **usability**, **hosting**, **performance**, **community**, **adoption**, and **licensing/cost**. A summary comparison table and recommendations for typical use cases follow.

Milvus

- **Core features:** Milvus is an open-source vector database optimized for high-dimensional similarity search. It supports **ANN** (approximate nearest neighbor) search over dense embeddings and allows **hybrid queries** (combining vector and scalar filters) and **multi-vector fields** (searching multiple vector columns simultaneously) [1](#) [2](#). Milvus supports data sharding and streaming ingestion. Key capabilities include **dynamic schema changes**, **distributed indexes** (Faiss, HNSW, DiskANN, SCANN), **GPU acceleration**, and sparse/dense vector hybrid search [1](#) [2](#).
- **Data types:** Primarily works with vector embeddings (e.g. text or image embeddings), but also stores associated metadata. Milvus collections can include scalar fields (numbers, text) and one or more vector fields. It does not natively index raw images or audio – instead, images/audio must be pre-encoded to vectors (e.g. via TensorFlow/PyTorch models) [2](#). Metadata (payload) fields support filters and can be used in hybrid queries.
- **Architecture & scalability:** Milvus uses a shared-nothing, **disaggregated architecture** with separate *storage* (MinIO/S3) and *compute* nodes. It has four main layers: Access, Coordinator, Proxy, and Storage [3](#). Collections are sharded across nodes for horizontal scale. It offers both a *standalone* single-node mode (for prototyping) and a *clustered* mode on Kubernetes or Docker. Automatic load-balancing and sharding allow Milvus to scale to billions of vectors.
- **Search capabilities:** Supports **semantic (vector) search** and **keyword search** (via text fields, optionally using BM25 on sparse features) and **hybrid search**. Milvus can combine dense and sparse vectors for a query (e.g. BERT plus TF-IDF/BM25) [2](#). It also supports multi-modal queries by searching multiple vector fields (e.g. text+image) simultaneously and merging results [2](#).
- **ML/NLP integration:** Milvus does not include built-in embedding models but integrates well with external ML tools. Users typically feed Milvus with embeddings generated by Hugging Face/SBERT or other frameworks. It has Python/Java/Go/Node clients, making it easy to use with Python ML stacks or frameworks like LangChain (there is a `langchain-milvus` integration) and Hugging Face Transformers [4](#) [2](#). Milvus is also used as a vector store in ML pipelines (e.g. RAG workflows) and has integrations in frameworks like Haystack, LlamaIndex, etc.
- **Ecosystem & plugins:** Milvus is backed by Zilliz, with a growing ecosystem. It provides a built-in **web UI** and there are connectors (e.g. to Dremio, OpenMetaverse, etc.). Milvus Lite (embedded mode) allows simple file-based storage for quick prototypes [5](#). Zilliz Cloud offers a hosted Milvus service. There are no formal “plugins” like in Weaviate, but Milvus continually adds index types and features.
- **Language support:** Official client SDKs exist for Python, Go, Java, Node.js, and C#. APIs are REST/gRPC; command-line tools exist as well. The Python client (`pymilvus`) is mature and widely used. Setup examples show multi-language code (see LangChain docs) [4](#) [6](#). The APIs are fairly

straightforward, and since Milvus has CRUD operations (create collection, insert vectors, search), usage is similar to other databases.

- **Ease of use & setup:** Milvus provides a Docker container and Helm charts for quick start. Milvus Standalone is easy to run locally for development (all components in one container) [7](#). Schema definitions are dynamic (collections can be created on-the-fly). The documentation is comprehensive. Out of the box it requires a bit more setup than a single-node text search engine (like Elasticsearch), but the matrix of options (indexes, parameters) is well documented. Zilliz also provides Milvus Lite for single-file use, making local testing trivial [5](#).
- **Hosting:** Can be self-hosted on any cloud or on-prem via Docker/K8s. Zilliz offers **Zilliz Cloud**, a managed Milvus service with a simple managed API [8](#). On AWS/Azure/GCP one can run Milvus on VMs or Kubernetes. No official hosted SaaS by AWS, but third-party managed Milvus services exist.
- **Performance:** Milvus is highly optimized for indexing speed. In benchmarks, Milvus had extremely fast ingest and indexing times relative to others [9](#). For example, one benchmark found Milvus indexed 10M 96D vectors in ~16 minutes, far faster than Elasticsearch (which took hours) [9](#). Query latency for Milvus can be sub-10ms on large datasets with GPU support, though it varies by index type and data size. Its search accuracy (recall) is high due to configurable indexes. According to Qdrant's benchmark, Milvus had the fastest index build time, though Qdrant outperformed it on query throughput at higher dimensions [9](#).
- **Community & docs:** Milvus is popular and well-documented. It has an active GitHub (tens of thousands of stars) and community forum. Documentation covers all features and use cases. The LangChain and other AI communities often use Milvus, and you'll find tutorials (Milvus + Haystack, etc.). Support is available via Zilliz (commercial) and an open Slack forum.
- **Production adoption:** Milvus is used in production by companies like Zilliz customers, Taobao, Tencent, etc. It's CNCF-graduated (CNCF incubating or graduated) and considered production-ready. Many RAG and recommendation systems use Milvus.
- **Licensing & cost:** Milvus is Apache-2.0 open source (free). There are no license costs for self-host. Zilliz Cloud is paid (pay-as-you-go) but no pricing was included here. Overall, Milvus is free to use with optional commercial support.

OpenSearch (Amazon OpenSearch Service)

- **Core features:** OpenSearch is a community-driven fork of Elasticsearch, a general-purpose search engine with added **vector search capabilities**. It provides full-text (keyword) search, analytics, alerting, and now a **knn_vector field type** for embedding search [10](#). It supports approximate k-NN algorithms (HNSW via Lucene, Faiss, NMSLIB) for ANN search [11](#). Hybrid queries are possible by combining vector filtering with text queries and other filters. OpenSearch also includes machine learning plugins to *generate* embeddings (using Painless scripts or ingestion pipelines) [12](#). In summary, it offers a full-text search engine with integrated vector search.
- **Data types:** Stores arbitrary JSON documents. Supports scalar fields (text, number, date, geo, etc.) and a special **knn_vector** field (dense vectors). One can also use Lucene's **sparse vector fields** (for neural ranking/token weights). OpenSearch can ingest embeddings as a field or compute them via an *ingest pipeline* (e.g. using Hugging Face or OpenAI models) [12](#). Unlike Milvus or Qdrant, OpenSearch does not natively index images; images must be converted to text vectors beforehand.
- **Architecture & scalability:** OpenSearch inherits Elasticsearch's architecture: a cluster of nodes with data and master roles. It uses Lucene segments for storage and supports sharding and replication for scaling and redundancy. It has been optimized in versions 3.x to improve concurrency (vector and text queries) [13](#). OpenSearch 3.0 introduced "concurrent segment search" by default, giving up

to 2.5x faster vector query performance ¹³. CPU and memory usage are managed per shard. For large datasets, it can run on many nodes (e.g. AWS Service scales elastically).

- **Search capabilities:** OpenSearch provides **semantic vector search** (via ANN) and **standard keyword search** (inverted index with BM25). Hybrid search is supported by combining full-text queries with vector queries, including boosting and filtering. In 2.x it added **neural sparse search** (sparse vector tokens) for semantic ranking over text. The vector and text queries can be combined using Boolean queries.
- **ML/NLP integration:** OpenSearch has a built-in ML Commons plugin allowing on-node model inference (BERT, etc.) during ingest or query. You can register Hugging Face/ONNX models and use a **text_embedding** processor in ingest pipelines ¹². This means you can call OpenAI or local Torch models to fill a vector field upon indexing ¹². Also, since it uses Elasticsearch APIs, it integrates with Elastic-trained models and frameworks. LangChain and others support OpenSearch (via its JSON API). ML-powered features include anomaly detection and k-NN queries. The AWS OpenSearch Service also provides these features in managed form.
- **Ecosystem & plugins:** OpenSearch is a broad platform. Beyond vector, it supports plugins for alerting, SQL, traces, etc. It has a community of dashboards (Kibana fork), tools (Beats, Logstash). For vector specifically, OpenSearch's v3 is refining disk-based vector storage (binary quantization) for cost efficiency ¹⁴ ¹³. It is more general than most vector-only DBs. Unlike others, there is no "module" system – instead, features are built into the core distribution or plugins. Many AWS services (like OpenSearch Service) provide easy deployment.
- **Language support:** OpenSearch uses the Elasticsearch REST API. Official clients exist for dozens of languages (Python `opensearch-py`, Java, .NET, etc.). Any HTTP client can talk to it. This makes integration easy if you've used Elasticsearch. Query syntax is JSON DSL.
- **Ease of use & setup:** It can be as simple as running an OpenSearch Docker container or using AWS OpenSearch. Schemas ("mappings") define vector and text fields. The learning curve is moderate due to Elasticsearch heritage. Many engineers are already familiar with the ecosystem. The vector search features require setting `index.knn: true` and defining `knn_vector` fields (as shown in docs ¹⁵). AWS offers a serverless OpenSearch option with vector collections (abstracting shards) ¹⁶.
- **Hosting:** OpenSearch can be self-hosted (Linux, Docker, Kubernetes). The most common hosting is **Amazon OpenSearch Service (Managed)**, which offers both provisioned and serverless clusters. Other cloud providers do not officially offer it (though Azure supports Elasticsearch). Since it's open source, there are no lock-ins. AWS Service also provides a vector database mode (OpenSearch Serverless collections with k-NN) ¹⁶.
- **Performance:** With recent updates, OpenSearch 3.x is much faster than earlier versions (text query latency ~10x faster than 1.x) ¹⁷. For vector search, concurrent segment execution yields up to 2.5x speedups with no accuracy loss ¹³. Ingest throughput is high (one test sustained ~20k docs/s for 70M docs on a single r5 node ¹⁸). Benchmarks comparing vector DBs show OpenSearch (Elasticsearch) lagging behind specialized DBs: one study found Elasticsearch took ~80 minutes to index 1M 96-D vectors vs. Milvus in under 1 min ⁹. However, Elasticsearch/OpenSearch scaled query performance has improved. In practice, query latencies (with optimized HNSW settings) can be on the order of tens of milliseconds at the million-vector scale.
- **Community & docs:** OpenSearch inherited the large Elasticsearch community. It has extensive docs and community forums. Vector search is new, but documentation (at docs.opensearch.org) is thorough ¹⁰ ¹³. AWS provides blog posts and tutorials (e.g. vector search with RAG). The project is backed by AWS and a community of contributors.
- **Production adoption:** OpenSearch is battle-tested for logging, search, and is production-grade. Vector search was added in 2021 (Elasticsearch k-NN plugin open sourced as OpenSearch), so

adoption is ramping. AWS clients use it for RAG and recommendations. Its robustness and HA features (snapshot, node recovery) make it suitable for mission-critical workloads.

- **Licensing & cost:** OpenSearch is Apache-2.0 (open source). AWS OpenSearch Service has pricing (instance-hours, etc.), but the software itself is free. No license costs. AWS even offers a free trial tier. As a fork of Elasticsearch 7.x, it avoided Elastic's SSPL.

Qdrant

- **Core features:** Qdrant is an open-source vector search engine (Rust-based) focused on high performance and scalability. It supports **dense vector search**, filtering by payload, and **sparse vector support** for keyword-like relevance (sparse vectors generalize BM25) ¹⁹. It provides **scalar filtering** (range, geo, full-text on payload) with Boolean logic, and **hybrid queries** by combining dense and sparse vectors. It offers **asynchronous I/O, write-ahead logging (WAL)** for durability, and **SIMD acceleration** for faster compute ²⁰. Qdrant's design emphasizes low memory usage (96-bit quantization reducing RAM by up to 97%) and performance on multi-dimensional data ²¹.
- **Data types:** Qdrant stores high-dimensional numeric vectors along with a JSON "payload" per vector. Payloads can be text, numbers, geo, arrays, etc., and are indexed for filtering ²². It can attach arbitrary JSON data to points. Qdrant does not natively process raw images or documents; these must be converted to vector embeddings externally. Thus, supported data are numeric vectors plus user-defined metadata. It recently added sparse (TF-IDF-like) vectors for semantic token weights ²³.
- **Architecture & scalability:** Qdrant is distributed: it supports **horizontal scaling** via sharding collections across nodes and replication for high availability ²⁴. It allows zero-downtime scaling (new shards can be added without downtime) ²⁴. Clusters can have a mix of CPUs/GPUs (GPU indexing supported). Qdrant Cloud (managed service) and Hybrid Cloud (self-managed) are offered. The Rust core is optimized for performance; indexes are in-memory by default, with options to offload to disk. Qdrant supports dynamic rebalancing and node coordination, making it suitable for large multi-tenant deployments.
- **Search capabilities:** Qdrant performs **nearest-neighbor vector search** using HNSW (via its NMSLIB engine) by default. It also supports **keyword search** through payload filtering (exact match, range, geo, etc.) and the new **sparse vector search** which applies neural token weights (like semantic scoring). Hybrid search is done by chaining: e.g. filter by keywords (payload), then do vector search on the result set. It recently added query planning to optimize mixed searches ²⁰.
- **ML/NLP integration:** Qdrant has many integrations: official guides exist for using Cohere embeddings, Haystack, LangChain, LlamaIndex, ChatGPT Retrieval Plugin, and Microsoft Semantic Kernel with Qdrant ²⁵. In practice, developers use any embedding model (OpenAI, HuggingFace, Cohere, etc.) to generate vectors, then store in Qdrant. While Qdrant doesn't include embedding models, its **FastEmbed** library provides a Scala/BERT-based embedding framework that can be used with Qdrant. There's also a Qdrant Model Context Protocol (MCP) server to serve models. LangChain and similar RAG frameworks support Qdrant natively (see LangChain docs for a Qdrant connector).
- **Ecosystem & plugins:** Qdrant's ecosystem includes the **FastEmbed** library (for reranking), **Qdrant Cloud/Hybrid**, and community projects. It has a built-in **web UI** for cluster management. The API is REST/gRPC; one can use OpenAPI to generate clients. Qdrant Cloud is a fully-managed SaaS. There's also Qdrant MCP Server for model hosting. Its open GitHub has many third-party projects (docs, clients, example apps). No plugin system per se, but a rich API and event triggers.
- **Language support:** Official client libraries exist for **Python** (qdrant-client), **JavaScript/TypeScript**, **Rust**, **Go**, **.NET**, and **Java** ²⁶. The REST API makes any language possible. Clients are well-

maintained. The Python client even includes convenience features like integrated filtering. Because Qdrant uses simple HTTP/gRPC, integration is considered easy.

- **Ease of use & setup:** Qdrant is easy to run via Docker or binary. A single-node launch (e.g. `docker run qdrant/qdrant`) gets you started. The HTTP API is straightforward (create collection, insert points, search). Schema definitions include vector size and optional payload index rules. The learning curve is low for basic use. Advanced features (e.g. distributed mode) require more setup but are documented. The Qdrant documentation is clear, with quickstarts and examples (e.g. for Haystack/LLM). Overall, ease-of-use is rated high among vector DBs due to simplicity of API.
- **Hosting:** Self-hosted Qdrant is free (Apache-2.0) and can run on-prem or in cloud (container/K8s). Qdrant also offers **Qdrant Cloud** (managed vector DB) with a free tier (first 1GB cluster free) and pay-as-you-go pricing ²⁷. They also have a Hybrid Cloud edition (self-hosted but billed hourly). There is no AWS or GCP managed Qdrant service beyond these.
- **Performance:** In published benchmarks (Azure D8s VM, 1M–10M vectors), Qdrant achieved the highest throughput (RPS) and lowest latencies in most scenarios ²⁸. It was notably 4x faster than before on one dataset. A key point: “*Qdrant achieves highest RPS and lowest latencies in almost all scenarios*” ²⁸. In the same test, Milvus had the fastest indexing but slower query throughput at large scales. Qdrant’s built-in 96-bit quantization dramatically reduces memory for large vectors. Its search latencies (single-threaded p95) often measure in a few ms for ~1M vectors, and it scales well with threads.
- **Community & docs:** Qdrant has a strong community (e.g. 24k GitHub stars) and active Discord. Documentation is detailed (architecture, tutorials). The Qdrant team publishes benchmarks and blogposts. Community contributions (plugins, benchmarks) exist. Support tiers (Enterprise features, priority support) are available from Qdrant.
- **Production adoption:** Qdrant boasts enterprise users (TripAdvisor, Cognizant, Hubspot, Bosch, etc. as seen on their site ²⁹ ³⁰) and open-source adoption is growing. The testimonials highlight ease and performance (e.g. “*Qdrant powers our demanding recommendation and RAG applications*”) ³¹. Many modern stacks (RAG pipelines, recommendation engines) use Qdrant in production at scale.
- **Licensing & cost:** Qdrant core is Apache-2.0 (open source) ³². There is no license cost for self-host. Qdrant Cloud and Hybrid have pricing (\$ per hour or vector). The site indicates a free 1GB managed cluster tier ²⁷. Overall, Qdrant is free to use at any scale if self-hosted, with optional paid managed/cloud offerings.

txtai

- **Core features:** txtai is a Python-based “**AI framework**” featuring an **embeddings database**. It unifies **dense and sparse vector indexes, knowledge graphs, and a relational store** under one hood ³³. Its vector search supports both **sparse (TF-IDF/BM25-like)** and **dense (neural)** retrieval, combined via SQL-like queries. In addition to search, txtai provides full-text processing pipelines: it can apply Transformers for question-answering, summarization, transcription, translation, etc., and orchestrate them in pipelines or agents. It also supports topic modeling and graph analysis. In short, txtai is more than a database; it’s a toolkit for semantic search and language workflows.
- **Data types:** txtai indexes text documents and their embeddings. It can generate and store embeddings for **text, documents, audio, images, and video** ³⁴. It also has an SQL-accessible store for metadata. Internally, it uses SQLite/PostgreSQL for relational data and FAISS or Annoy for vector data. Multimodal support is achieved by embedding images/audio into vectors as well. txtai does not natively store raw images/files – everything is ultimately converted to embeddings and text.

- **Architecture & scalability:** txtai is a library (not a standalone server). It can run locally or distributed as microservices (via FastAPI). It supports **SQLite for small scale or PostgreSQL for larger scale** (for metadata) and FAISS/Annoy for vectors. It can be scaled with containers: e.g. multiple txtai instances behind a FastAPI front-end with shared DB. Because it's Python-based, it may not match the raw performance of a Rust/Go engine, but it offers flexibility. There's also **txtai.cloud**, a hosted vector DB service (preview) for managed usage. For enterprise, it can be containerized and deployed across clusters.
- **Search capabilities:** txtai provides **semantic (vector) search** as its core. It also allows **SQL queries** over metadata and joins between text and graph data. It uses Faiss for ANN. Keyword search is possible via SQLite full-text indexes. Hybrid search (combining dense and sparse) is possible because txtai supports sparse vector indexes (BM25-like) alongside dense embeddings ³³. Essentially, one can query both text and embeddings via the same API. However, txtai is primarily oriented to neural search, so traditional inverted-index features (like fuzzy search) are not its focus.
- **ML/NLP integration:** txtai is built on Transformers (Hugging Face) and has first-class support for Hugging Face models and others. It comes with pre-configured pipelines: e.g. one-line config to use a Hugging Face sentence transformer as the embedding model ³⁵. It supports any Hugging Face or SpaCy model and can use OpenAI APIs as part of pipelines. It has built-in LLM orchestration (MCP – Model Context Protocol) and can call LLMs for tasks (RAG, summarization). Integration with LangChain/Haystack is possible via its database interface, but txtai itself fills a similar niche as those libraries. Bindings exist for **JavaScript, Java, Rust, and Go** (via its MCP and API) ³⁶, so it can be used outside Python too.
- **Ecosystem & plugins:** As a relatively new framework, txtai has fewer “plugins.” However, it has a rich set of *examples/notebooks* (over 60 demos) for various use cases (QA, video search, chatbot) ³⁷. The core is maintained by Neurolab (John Snow Labs), and there's a Slack channel and GitHub issues for support. The “txtai cloud” suggests an aim to be a service. There is no marketplace of extensions, but many components (FAISS, FastAPI, sqlalchemy) are modular.
- **Language support:** Primarily Python 3.10+. The main interface is Python (`import txtai`), but it exposes a REST/MCP API with clients in JavaScript, Java, Rust, and Go ³⁶. This means you can embed txtai in multi-language stacks. However, since the database part is SQLite/FAISS, concurrency is limited to what the chosen DB allows unless using txtai's single-process model or the cloud service. The development experience is smooth in Python, but if you want a pure Python solution that includes the search engine, you get it; otherwise, it's an embedding service with API.
- **Ease of use & setup:** txtai is easy to install (`pip install txtai`) and comes with sensible defaults (“batteries included”) ³⁶. A simple example shows indexing two sentences with one call ³⁵. It includes a built-in FastAPI server (`uvicorn "txtai.api:app"`) for REST calls. Users can configure pipelines via YAML (specifying models). It is generally easier to set up than a full cluster, but performance tuning may require knowledge of FAISS or SQL DB. Because it's a Python library, one needs to manage the environment and dependencies.
- **Hosting:** txtai can run anywhere Python can. Options: run on a VM, container (Docker is provided), or as a serverless function (with limitations). The project mentions **txtai.cloud** (a hosted offering currently in preview). It also suggests an Enterprise service. There is no big cloud provider offering, but one could deploy it on AWS/GCP/Azure. For small-to-medium apps, the local/lib approach is fine.
- **Performance:** As a Python-based system, txtai's raw throughput is lower than native engines. However, it leverages FAISS/Annoy (C++) for vectors, so search speed can be good. Benchmarks are scarce, but expect single-machine QPS in the low hundreds and latency ~10ms for a few million vectors (depending on FAISS index settings). Its strength is feature-richness, not extreme performance. Since it can use PostgreSQL or SQLite, SQL queries scale as those do (SQLite is single-

threaded; Postgres better). Users typically optimize by reducing index size or pre-embedding. Its embedding pipelines can be slow if the model is large, though caching helps.

- **Community & docs:** txtai's documentation (mkdocs) is thorough, with quickstarts and examples. The GitHub is active and well-maintained (Apache-2.0 license) ³⁸. Community is smaller than the other projects here, but it has backing from John Snow Labs and some enterprise interest. There is a Slack/Discord channel and support forum.
- **Production adoption:** txtai is newer and less battle-tested at massive scale compared to others. It is used in production by some companies building text search and RAG pipelines (e.g. within Rasa). It suits startups and data science teams needing an “all-in-one” solution without hiring DevOps. Major enterprises might opt for dedicated DBs, but txtai is viable for many semantic search applications.
- **Licensing & cost:** txtai is **Apache-2.0 open source** ³⁸. It has no licensing cost. The mention of txtai.cloud implies future paid hosting, but core usage is free. It competes with hosted RAG platforms by being self-hostable at no cost.

Weaviate

- **Core features:** Weaviate is an open-source **AI-native vector database with a knowledge graph** orientation. It stores JSON objects (“objects”) with associated vectors and supports **vector search**, **graph queries**, and **hybrid filters** ³⁹ ⁴⁰. Unique features include *GraphQL API*, a *modular* architecture (many “modules” for different AI models and functions), and *auto-ML* modules for embedding generation (text2vec-transformers, image2vec, QnA, etc.). Weaviate focuses on enabling semantic search with schema-based objects. It also supports **grouped/neural search** (neural sparse search via Lucene pruning) and **aggregation queries**.
- **Data types:** It is multi-modal: supports text, images, and more. Built-in modules allow ingesting images as raw media (auto-vectorized via CLIP), text via BERT or others, audio, etc. The object model (class-based schema) allows rich metadata (string, int, date, geo, etc.), and each object can have *one or more named vector fields*. Weaviate’s “objects” are essentially JSON documents. It also supports storing multiple vectors per object (multi-vector). Sparse indexing is supported through the neural search modules. In short, Weaviate covers text, images, and any data type you can vectorize.
- **Architecture & scalability:** Weaviate is written in Go and designed for cloud-native scale. It uses sharding and replication for distributed deployment (via Kubernetes). It is horizontally scalable: you can run a cluster with many nodes and it will distribute data. By default it uses RocksDB on each node and keeps HNSW graph in memory, but can offload to disk. It also supports fault tolerance (Raft for cluster state). The recommended deployment is Docker/Kubernetes. Weaviate Cloud Service is also available (HNSW graphs auto-partitioned across tenants). Overall, architecture is mature (it is CNCF-incubating or graduated).
- **Search capabilities:** Out-of-the-box **semantic search** on vectors. Weaviate provides standard vector KNN search and **hybrid search** combining vector similarity with filters and keywords ⁴¹. It also supports *graph-based* queries (traversing object relations in the graph) via GraphQL. For text search, it includes a BM25-like mode (neural search). You can do text queries (keyword + synonyms) as well as vector queries in the same API. Weaviate has a unique “nearText” operator (with concepts and certainty) and “nearObject” (vector similarity) for retrieval. Hybrid search (text + vector) is achieved by mixing GraphQL filters with “nearVector” clauses.
- **ML/NLP integration:** Weaviate's modular design includes **text2vec-transformers** (embedding with Transformers) and **clip2vec** (for images) as built-in modules. It also supports **text2vec-contextionary** (fastText-based) and **text2vec-gpt2**. Custom modules can be added (e.g. Hugging Face, GPT, Stable Diffusion for image, etc.) – see Weaviate “Modules” reference. It also integrates

with other tools: There are documented connectors for LangChain, Haystack, etc. Weaviate itself can call OpenAI for embeddings via its OpenAI module. Essentially, Weaviate aims to manage ML models internally, so it has tight integration with popular NLP/ML frameworks under the hood.

- **Ecosystem & plugins:** Weaviate has a rich ecosystem. Notably, it is a CNCF project. It provides a **GraphQL API** for queries, with auto-generated Swagger docs and client libraries. There is **Weaviate Cloud** (serverless and managed options), and integrations like Pinecone-like “vector indexer” for hybrid. Clients and connectors exist (LangChain, LlamaIndex, Superset, etc.). It also supports **Azure Cognitive Search** integration (as an index). Weaviate Agents (a new feature) provide RAG pipeline orchestration within the DB. While it doesn’t have third-party “plugins” directory, its built-in Modules act like plugins.
- **Language support:** Weaviate has official client libraries for **Python, Java, Go, JavaScript** (v4 and v3 of some clients) ⁴², and community ones (C#, etc.). Queries can be done via HTTP or GraphQL in any language. The GraphQL API is comprehensive, so any GraphQL client works too. It integrates with ML stacks via its API (e.g. Python’s `weaviate-client`). Overall, language support is broad and approachable.
- **Ease of use & setup:** Weaviate provides Docker and Kubernetes Helm charts. It is relatively straightforward to start (even a one-node Docker image). Defining a schema (classes) via API or YAML is required. The built-in UI (GraphQL playground) helps explore. Because of its many features, learning Weaviate’s concepts (modules, classes, GraphQL syntax) takes some effort. However, the documentation is thorough, with many examples. Many users find the GraphQL interface intuitive once learned. The cloud version further simplifies usage with web console.
- **Hosting:** Self-host on Kubernetes or Docker. There is a **Weaviate Cloud** offering (Serverless or Enterprise Cloud) with pay-per-use pricing. Weaviate Cloud starts at \$25/mo per 1M dimension-month ⁴³. Additionally, one can deploy Weaviate on any cloud VM or Kubernetes cluster. It is not offered as an AWS/GCP managed service.
- **Performance:** Weaviate advertises very fast KNN: “nearest neighbor searches of millions of objects in <100ms” ⁴⁴. Benchmarks on Medium-size hardware confirm sub-50ms queries at million-scale. Performance depends on vector dimensionality and HNSW parameters. According to Qdrant’s report, Weaviate showed relatively slower performance in their last round of tests ⁴⁵, but the company has optimized since. Index creation in Weaviate can take longer due to heavy integration. In practice, Weaviate scales well, and the asynchronous nature of GraphQL means queries don’t block the DB.
- **Community & docs:** Weaviate has extensive documentation and an active community (Slack, forum, StackOverflow). It’s mature, with many tutorials (official documentation site with how-tos). CNCF affiliation means visibility. The developers publish benchmarks and blogs. Enterprise support is available via SeMI Technologies. Overall, community support is strong.
- **Production adoption:** Weaviate is in use at large organizations (e.g. Zilliz Cloud uses it, plus customers like Epson, “Fits Music” etc.). It’s often chosen for applications needing vector search + knowledge graph (e.g. product catalogs with structure). Its cloud service has notable users, and it’s cited in vector DB market analyses. It is considered production-ready for most use cases.
- **Licensing & cost:** Weaviate’s core is **BSD-3-Clause** open source ⁴⁶. That is very permissive (allows custom modifications). Weaviate Cloud is commercial (\$), and enterprise subscriptions (support, training) exist ⁴⁷. But the software itself is free. The BSD-3 license means no viral restrictions.

Typesense

- **Core features:** Typesense is an open-source **search engine** initially designed for fast typo-tolerant full-text search (Algolia-like), now also supporting vector and semantic search. It offers **instant search (autocomplete)**, typo correction, faceting, geo-search, etc., plus **vector (KNN) search** on embedding fields ⁴⁸. It provides built-in semantic search by automatically embedding queries/documents via Hugging Face or APIs, and supports **hybrid search** combining keyword and vector with configurable weighting ⁴⁸ ⁴⁹. Key selling points are ease-of-use and speed for developer-facing search.
- **Data types:** Typesense indexes JSON documents with various field types (string, int, array, geo, etc.). A special **vector field** can hold numeric arrays for embeddings ⁵⁰. Users can supply embeddings (Option A) or let Typesense auto-generate them (Option B) using OpenAI, PaLM, or built-in models ⁵¹. It does not directly process images/audio – one would pre-compute vector via CLIP or other model. It includes a “*semantic search*” mode where string queries are internally turned into embeddings (e.g. using a default MiniLM model) ⁵² ⁵³. Traditional keyword search is still available on text fields (with optional filters).
- **Architecture & scalability:** Typesense is written in C++ for speed and runs as a single process (monolithic) per cluster node. It supports clustering (multiple nodes) with automatic data sharding and replication. Each node holds one or more shards, and Typesense handles distribution transparently (similar to Elasticsearch). It uses in-memory indexes with optional disk persistence. It has relatively low resource usage. Clustering is more limited than in heavy-duty systems: it supports primary/replica shards but does not have built-in dynamic rebalancing (shard count is fixed at creation). Overall, it can scale to billions of docs by adding nodes, but is optimized for smaller-to-mid scale (up to hundred-million docs).
- **Search capabilities:** **Keyword search:** full-text with relevance tuning, faceting, geo, synonyms, etc. **Semantic search:** uses embeddings to find conceptually similar results ⁵⁰ ⁵². **Hybrid search:** combined by weighting vector similarity and text relevance; Typesense allows adjusting the mix ⁵⁴ ⁵⁵. The Docs mention “Rank keyword search via vector search” and “Re-ranking hybrid matches” features. It can do classic k-NN (HNSW) queries and optionally re-rank results using vector distance. It also supports *edge cases* like filtering by previous queries.
- **ML/NLP integration:** Typesense has built-in integration with ML models: it can call external APIs (OpenAI, Google PaLM/Vertex) or use built-in transformer models from HuggingFace for embeddings (e.g. S-BERT models via HuggingFace hub) ⁵⁶ ⁵³. This is configurable per collection. It also has a hosted option for GPU acceleration of embeddings ⁵⁷. There is no direct LangChain or Haystack connector, but Typesense can serve as any search backend via API. It’s often used with external embedding pipelines (e.g. embedding documents with HF then indexing into Typesense).
- **Ecosystem & plugins:** Typesense has its own ecosystem centered around its API. It offers official **UI widgets** for building search UIs, and a variety of client libraries (see below). Integrations include e-commerce connectors (Magento, Shopify via plugins) and CMS (WordPress plugin exists). Typesense Cloud (SaaS) adds multi-tenancy. There’s no plugin architecture like Elastic, but the core covers many needs. Its Roadmap highlights vector and RAG features.
- **Language support:** Official client libraries exist for **JavaScript, PHP, Python, and Ruby** ⁵⁸. Community libraries cover Go, .NET, Java, Rust, Dart, etc. since it’s a REST API. The API is JSON/HTTP, so virtually any language can be used. The clients provide easy methods for all endpoints (collections, documents, search). For ML use cases, developers may use Python to generate embeddings but query Typesense via any supported client.

- **Ease of use & setup:** Typesense is very developer-friendly. Installing (binary or Docker) is straightforward. Defining a schema is simple JSON. They boast “no PhD required”. Typo tolerance and ranking have sensible defaults, and tuning is done via simple parameters. Semantic search requires adding vector fields and specifying embedding models, but even that is guided in docs. The learning curve is low if you’re familiar with typical search (it’s simpler than Elasticsearch). The built-in Cloud service can spin up clusters in minutes.
- **Hosting:** Typesense can be self-hosted on any Linux (Docker or Linux binary). There is a managed **Typesense Cloud** (AWS/EU regions) for those who prefer SaaS. Cloud pricing is tiered; there’s a free tier and paid plans (\$/vector-dim or nodes) ⁵⁹. It is not offered by AWS/Azure as a managed service aside from their own Cloud. The open-source engine has no cost.
- **Performance:** Very fast for text search. Benchmarks (and user reports) indicate sub-millisecond retrieval on typical sized indices for keyword queries. For vector search, it uses HNSW and is efficient at moderate scale. The overhead of auto-embeddings (if using external APIs) can introduce latency, but once vectors are stored, vector search is quick. Typesense claims “blazing-fast” and is optimized for CPU/Memory. It has ~18M Docker pulls, indicating popularity. It’s comparable to other search engines like Meilisearch.
- **Community & docs:** Typesense has good documentation and an active GitHub community (13k stars). The team is responsive and publishes roadmap/features updates. There’s a forum and Slack for support. The site showcases many user logos (Codecademy, Logitech) ⁶⁰, suggesting trust. The project is less known in academic circles than Milvus or Elasticsearch, but growing especially among web developers.
- **Production adoption:** Used widely for website search, e-commerce sites (it competes with Algolia). It’s not primarily known for cutting-edge ML use, but with v0.25+ it added RAG and vector features, making it attractive for startups and SMEs who want semantic search without heavy ops. Its low operational complexity makes it popular for production deployments where ease is valued.
- **Licensing & cost:** Typesense Server is **AGPL-3.0** open source (**note:** Typesense’s engine is under the AGPL license, which is copyleft). This requires that services using it open source their derivative works. (The Cloud service is obviously proprietary but that’s separate.) There is no cost to self-host aside from possible support. The Cloud service has usage pricing (as per [87] and [88]).

Summary Comparison

Feature / Capability	Milvus	OpenSearch	Qdrant	txtai	Weaviate	Typesense
Engine type	Vector DB (ANN)	Search+Analytics Engine	Vector DB (ANN)	Embeddings DB + pipelines	Vector DB + Knowledge Graph	Search Eng (text + vect)
Primary use	High-scale similarity search	Text search + analytics + vector	High-performance similarity search	Semantic search + LLM workflows	Semantic search & graph queries	Typo-tolerant search, now semantic
Vector search	ANN (Faiss/ HNSW, GPU accel.)	ANN (Lucene, Faiss, NMSLIB)	ANN (HNSW), quantized	ANN (FAISS/ Annoy)	ANN (HNSW), multi-vector fields	ANN (HNSW) integrated

Feature / Capability	Milvus	OpenSearch	Qdrant	txtai	Weaviate	TypeSense
Hybrid search (keyword+vector)	Yes (dense+BM25)	Yes (vector+full-text)	Yes (dense+sparse payload)	Yes (dense+sparse via SQL)	Yes (vector + filters/text)	Yes (vector+key weighting)
Data modalities	Vectors + scalar metadata	JSON docs (text, number)	Vectors + JSON payload	Text, audio, images, video (as embeddings)	Objects (text, images, etc.)	JSON docs (numeric, geo)
Embeddings generation	External	Built-in via ML plugins	External or custom	Built-in Transformers + pipelines	Built-in Transformers (modules)	Built-in (HF models, OpenAI/Pa APIs)
API / Query languages	Python/Go/Java/Node/gRPC, SQL-like	Elasticsearch DSL (JSON), SQL plugin	REST API (OpenAPI), gRPC	Python API, REST/MCP (SQL queries)	GraphQL & REST API	REST API (JSON queries), hybrid DSL
Client libraries	Python, Java, Go, Node, C#	Python, Java, .NET, etc.	Python, JS/TS, Rust, Go, .NET, Java	Python (core), others via API	Python, JS, Go, Java, etc.	JavaScript, Python, Ruby (official); Go, Java, .NET, etc. (community)
Ease of use	Moderate (some ops setup)	Moderate-high (known stack)	High (simple API)	High (as a library)	Moderate (schema + GraphQL)	Very high (simple API, great docs)
Scalability	Horiz. (shards, K8s)	Horiz. (shards/replicas)	Horiz. (shards, replica)	Vertical/node (FAISS + SQLite/Postgres)	Horiz. (shards)	Horiz. (shards/replicas)
Deployment	Self-hosted (cloud/k8s) + Zilliz Cloud	Self-hosted or AWS Service	Self-hosted + Qdrant Cloud	Self-hosted (Docker/VM) + txtai.cloud (beta)	Self-hosted + Weaviate Cloud	Self-hosted TypeSense Cloud (SaaS)
Performance (index/query)	Very fast indexing; low latency with tuning ⁹	Good text speed; vector search improved (2.5x faster in v3) ¹³	Excellent query RPS, low latency ²⁸	Good enough for moderate scale; depends on FAISS index	Fast (<100ms for million-scale) ⁴⁴	Extremely fast text search; good vector search

Feature / Capability	Milvus	OpenSearch	Qdrant	txtai	Weaviate	Typesense
Fault tolerance / HA	Yes (cluster mode)	Yes (cluster, shards)	Yes (replication, rolling updates) ²⁴	No built-in (use DB replication)	Yes (replication)	Yes (primary replica sharding)
Community & maturity	High (CNCF, Zilliz)	Very high (fork of Elasticsearch)	Growing (24K stars, active)	Growing (open-source, new)	Growing (CNCF, active)	Growing (1 star, developer-focused)
Production use / adoption	Widely used (AI apps, RAG)	Very widely (enterprise search)	Broad (AI services)	Niche (AI frameworks)	Increasing (enterprise vector)	Widespread (site search)
License	Apache-2.0	Apache-2.0	Apache-2.0 ³²	Apache-2.0 ³⁸	BSD-3-Clause ⁴⁶	AGPL-3.0 (Open source engine)
Hosted/service pricing	Zilliz Cloud (paid)	AWS OpenSearch (per instance)	Qdrant Cloud (free tier + \$) ²⁷	txtai Cloud (preview)	Weaviate Cloud (\$25+/mo) ⁴³	Typesense Cloud (\$)

Recommendations

- **Large-scale high-dimensional search:** If you need **extremely fast indexing** and bulk throughput for billions of vectors, Milvus and Qdrant are top choices. Milvus shines in GPU-accelerated, multi-vector use cases (and Zilliz Cloud managed service), while Qdrant offers exceptional throughput/latency on CPU with integrated filtering ²⁸. Use Milvus for heavy-duty GPU-backed indexing and if hybrid (dense+sparse) search is critical, or Qdrant for cost-efficient CPU and built-in quantization.
- **General-purpose search with semantics:** For a combination of full-text search and vector search in one system, OpenSearch or Typesense work well. OpenSearch is ideal if you already use Elasticsearch (it supports complex aggregations and analytics, plus newly built-in vector search) ¹⁰ ¹³. Typesense is excellent for web/mobile search apps (typo-tolerance, analytics, and ease of setup), now augmented with semantic/vector search ⁴⁸ ⁵². Use OpenSearch if you need enterprise features (security, monitoring) and scale; use Typesense for simpler deployments or if AGPL is acceptable.
- **Knowledge graph and structured filtering:** If your data has a rich schema or graph relationships, Weaviate is recommended. Its GraphQL API and ability to embed knowledge graphs set it apart ⁴⁰ ⁴⁶. It supports advanced use cases like classification, generative Q&A via modules, and is cloud-native. It is a strong choice for ML-first applications needing object semantics (e.g. product search with structured attributes).
- **Embedded/ML workflows:** For use cases involving custom pipelines, RAG, or in-application search (especially in Python), txtai can be a good fit. It's easy to integrate into Python code and provides more than just search (e.g. QA pipelines) ⁶¹. It's not for ultra-large scale, but perfect for teams who want built-in NLP tools plus vector search without managing separate services.
- **Budget constraints:** All except Typesense are open-source free. If license is a concern, Milvus, OpenSearch, Qdrant, txtai (Apache-2) and Weaviate (BSD) have permissive licenses, while Typesense

Core (AGPL) requires caution if deploying as a service. For purely self-hosted projects with no cost, these engines are free; only cloud services (Weaviate Cloud, Qdrant Cloud, Typesense Cloud, Zilliz Cloud) charge fees.

- **Ease of development:** For quick MVPs, Typesense and Weaviate have the lowest barrier (easy APIs, hosted options). Milvus and Qdrant require a bit more setup (docker cluster or instance), but have good SDKs. txtai is extremely easy to get started in Python. OpenSearch requires more DevOps knowledge unless using the AWS managed service.

In summary, choose **Milvus** or **Qdrant** for raw performance and large-scale similarity search; **OpenSearch** or **Typesense** if you need combined text+vector search in a familiar search engine environment; **Weaviate** for AI-native graph/ML features; and **txtai** for integrated ML/AI pipelines in a Pythonic stack. Always validate with benchmarks on your data and consider long-term operational costs (e.g. ease of scaling, available managed services) when making a decision.

1 3 Milvus Architecture Overview | Milvus Documentation

https://milvus.io/docs/architecture_overview.md

2 Multi-Vector Hybrid Search | Milvus Documentation

<https://milvus.io/docs/multi-vector-search.md>

4 5 6 7 8 Milvus | LangChain

<https://python.langchain.com/docs/integrations/vectorstores/milvus/>

9 28 45 Vector Database Benchmarks - Qdrant

<https://qdrant.tech/benchmarks/>

10 12 Vector search - OpenSearch Documentation

<https://docs.opensearch.org/docs/latest/vector-search/>

11 15 Approximate k-NN search - OpenSearch Documentation

<https://docs.opensearch.org/docs/latest/vector-search/vector-search-techniques/approximate-knn/>

13 14 17 18 OpenSearch Project update: Performance progress in OpenSearch 3.0 - OpenSearch

<https://opensearch.org/blog/opensearch-project-update-performance-progress-in-opensearch-3-0/>

16 Working with vector search collections - Amazon OpenSearch Service

<https://docs.aws.amazon.com/opensearch-service/latest/developerguide/serverless-vector-search.html>

19 20 21 22 23 24 25 32 GitHub - qdrant/qdrant: Qdrant - High-performance, massive-scale Vector Database and Vector Search Engine for the next generation of AI. Also available in the cloud <https://cloud.qdrant.io/>

<https://github.com/qdrant/qdrant>

26 API & SDKs - Qdrant

<https://qdrant.tech/documentation/interfaces/>

27 Pricing for Cloud and Vector Database Solutions Qdrant

<https://qdrant.tech/pricing/>

29 30 31 Qdrant - Vector Database - Qdrant

<https://qdrant.tech/>

33 34 35 36 37 38 61 GitHub - neuml/txtai: All-in-one open-source AI framework for semantic search, LLM orchestration and language model workflows

<https://github.com/neuml/txtai>

39 40 41 42 44 Introduction | Weaviate

<https://weaviate.io/developers/weaviate/introduction>

43 47 Vector Database Pricing - Weaviate

<https://weaviate.io/pricing>

46 What Is Weaviate? A Semantic Search Database | Oracle Eesti

<https://www.oracle.com/ee/database/vector-database/weaviate/>

48 49 50 51 54 55 56 57 Vector Search | Typesense

<https://typesense.org/docs/28.0/api/vector-search.html>

52 53 Semantic Search | Typesense

<https://typesense.org/docs/guide/semantic-search.html>

58 API Clients | Typesense

<https://typesense.org/docs/28.0/api/api-clients.html>

59 Pricing | Typesense Cloud

<https://cloud.typesense.org/pricing>

60 Typesense | Open Source Alternative to Algolia + Pinecone

<https://typesense.org/>