

## 向量索引類型比較報告：FLAT、IVF\_FLAT、IVF\_SQ8、IVF\_PQ、HNSW、ANNOY

### 1. 各索引類型的起源與隱喻說明

- **FLAT**

*隱喻：地毯式搜索*

想像你在一個巨大的圖書館裡找書，但沒有任何分類或指引，你只能一本一本地翻閱比對。這就是 FLAT 索引——暴力搜尋，逐一比較每個向量。

- **IVF\_FLAT**

*隱喻：分區尋寶*

圖書館被劃分成許多區域（分桶），你先找到最有可能藏書的幾個區域，再在這些區域裡一本本地找。IVF\_FLAT 就是先將向量分群（倒排檔案），查詢時只在最相關的幾個分群內做 FLAT 搜索。

- **IVF\_SQ8**

*隱喻：壓縮地圖尋寶*

在 IVF\_FLAT 的基礎上，將每本書的資訊壓縮成簡單的代碼（將浮點數壓縮成 8 位整數），大幅降低記憶體用量。你在分區內找書時，看的只是簡化過的書目卡片。

- **IVF\_PQ**

*隱喻：拼圖壓縮尋寶*

把每本書拆成幾個部分（子向量），每部分都用最接近的拼圖塊代表，整體大幅壓縮。你在分區內找書時，對比的是這些拼圖塊的組合，速度快但細節可能遺失。

- **HNSW**

*隱喻：多層地圖快速導航*

圖書館裡有多層樓、每層樓有捷徑，導航系統會先帶你到大致區域，再用捷徑快速接近目標。HNSW 就是這種圖結構索引，能快速縮小搜尋範圍。

- **ANNOY**

*隱喻：多條小徑隨機穿梭*

你在圖書館裡設置了很多條隨機小徑，每次找書時走不同的小徑，最後綜合多條路徑的結果。ANNOY 是多樹結構，查詢時隨機穿梭多棵樹，快速近似搜尋。

### 2. 主要特性總覽

索引類型	特性簡述
FLAT	精確、無壓縮、速度慢、記憶體需求大
IVF_FLAT	分桶加速、近似搜尋、速度快、可調 recall/速度、記憶體需求中等
IVF_SQ8	IVF_FLAT + 壓縮、記憶體需求低、速度快、精度略降
IVF_PQ	IVF_FLAT + 進階壓縮、記憶體需求最低、速度最快、精度下降明顯

索引類型	特性簡述
HNSW	基於圖的近似搜尋、查詢快、記憶體需求較高、建索引慢、支持動態資料
ANNOY	多樹結構、查詢快、記憶體需求低、建索引快、精度可調、適合靜態資料

### 3. 適用應用場景

索引類型	適用場景
FLAT	小型資料集、需極高精度、測試、驗證
IVF_FLAT	中大型資料集、需平衡查詢速度與精度、記憶體資源有限
IVF_SQ8	資源受限（如雲端/嵌入式）、需大幅壓縮、速度優先於精度
IVF_PQ	超大規模資料集、極端資源受限、對精度要求可接受
HNSW	高查詢頻率、需低延遲、需支援動態新增資料、推薦/搜尋等場景
ANNOY	靜態資料集、需多次查詢、快建索引、資源有限、推薦/相似查詢

### 4. 對應搜尋演算法

索引類型	搜尋演算法說明
FLAT	線性暴力搜尋（Brute-force KNN）
IVF_FLAT	倒排檔案分群 + 分群內線性搜尋（Inverted File + Local Linear Scan）
IVF_SQ8	IVF_FLAT + 標量量化（Scalar Quantization）
IVF_PQ	IVF_FLAT + 產品量化（Product Quantization）
HNSW	分層小世界圖導航（Hierarchical Navigable Small World Graph）
ANNOY	多棵隨機二分樹搜尋（Approximate Nearest Neighbor Oh Yeah）

### 5. 優缺點比較

索引類型	優點	缺點
FLAT	精確、無資訊損失	查詢慢、記憶體需求高
IVF_FLAT	查詢快、可調參數、資源需求較低	精度略降、需調參
IVF_SQ8	記憶體壓縮、查詢快	精度進一步下降
IVF_PQ	壓縮比最高、查詢最快	精度損失大、參數複雜
HNSW	查詢極快、支援動態資料、精度高	建索引慢、記憶體需求大
ANNOY	查詢快、建索引快、資源需求低、適合靜態資料	精度略低、不支援動態新增

## 6. Milvus Python 範例

### FLAT

```
index_params = {
    "index_type": "FLAT",
    "metric_type": "L2",
    "params": {}
}
collection.create_index(
    field_name="embedding",
    index_params=index_params
)
```

### IVF\_FLAT

```
index_params = {
    "index_type": "IVF_FLAT",
    "metric_type": "L2",
    "params": {"nlist": 1024}
}
collection.create_index(
    field_name="embedding",
    index_params=index_params
)
```

### IVF\_SQ8

```
index_params = {
    "index_type": "IVF_SQ8",
    "metric_type": "L2",
    "params": {"nlist": 1024}
}
collection.create_index(
    field_name="embedding",
    index_params=index_params
)
```

### IVF\_PQ

```
index_params = {
    "index_type": "IVF_PQ",
    "metric_type": "L2",
    "params": {"nlist": 1024, "m": 16}  # m為PQ子空間數
}
collection.create_index(
    field_name="embedding",
```

```
    index_params=index_params
)
```

## HNSW

```
index_params = {
    "index_type": "HNSW",
    "metric_type": "L2",
    "params": {"M": 16, "efConstruction": 60}
}
collection.create_index(
    field_name="embedding",
    index_params=index_params
)
```

## ANNOY

```
index_params = {
    "index_type": "ANNOY",
    "metric_type": "L2",
    "params": {"n_trees": 50}
}
collection.create_index(
    field_name="embedding",
    index_params=index_params
)
```

## 結論建議

- 小型資料集、需精確查詢：選 FLAT。
- 中大型、需平衡效能與精度：選 IVF\_FLAT。
- 極大資料集、資源有限：選 IVF\_SQ8 或 IVF\_PQ（PQ 壓縮比更高）。
- 查詢極速、支援動態資料：選 HNSW。
- 靜態資料、快建索引、多次查詢：選 ANNOY。

選擇索引時，請依據資料規模、查詢頻率、資源限制與精度需求綜合考量。