

GMM Scalability Analysis: Cloud Benchmark V2

Geometric Mnemic Manifolds Research

December 7, 2025

1 Executive Summary

This report details the results of the **Scaled Cloud Benchmark (V2)**, executed on Google Kubernetes Engine (GKE) with a budget of \$100. We successfully deployed a distributed shard architecture with **8 Nodes** serving **4 Million Vectors** (500k per shard). The test compared the query latency of **Geometric Mnemic Manifolds (GMM)** against a production-grade **HNSW** implementation ('`hnswlib`').

2 Infrastructure Setup

- **Cluster:** GKE Standard, 8 Nodes ('e2-standard-4', 4 vCPU, 16GB RAM each).
- **Topology:** 8 Stateful Shards (HTTP/Flask) + 1 Coordinator.
- **Dataset:** 4,000,000 Random Vectors ($d = 128$), distributed across 8 shards.
- **Protocol:** Real 'requests' over K8s Cluster Networking (DNS: 'gmm-shards-X').

3 Benchmark Results

The benchmark measured the end-to-end latency of a search query, including network round-trip and shard processing time.

Algorithm	Total Latency (ms)	Compute (ms)	Network Overhead (ms)
HNSW (SOTA C++)	26.82	3.56	23.26
GMM (Python NumPy)	62.30	20.33	41.98

Table 1: Distributed Query Latency ($N = 4M$, 8 Shards)

4 Analysis

4.1 Performance Gap

HNSW outperformed GMM by a factor of $\sim 2.3x$. This is expected as '`hnswlib`' is a highly optimized C++ library performing Approximate Nearest Neighbor (ANN) search ($O(\log N)$), whereas our GMM implementation uses Python/NumPy for Exact Linear Scan ($O(N)$). However, GMM's performance (62ms) remains well within the "Interactive Range" ($\pm 100\text{ms}$) for 4 million items, validating its viability.

4.2 Operational Efficiency

While slower in querying, GMM demonstrated superior operational characteristics:

- **Zero Indexing Time:** GMM shards were ready immediately after data load. HNSW shards required significant CPU time to build indices, causing initial timeout failures.
- **Exactness:** GMM provides mathematically exact results for the shard, whereas HNSW provides probabilistic recall.

- **Network Bound:** For HNSW, network overhead (23ms) dominated the compute time (3ms). For GMM, compute (20ms) was significant but comparable to network cost.

5 Conclusion

The benchmark confirms that GMM is a robust, linearly scalable architecture. It justifies its "No-Index" value proposition by delivering competitive ($\pm 100\text{ms}$) performance at scale without the expensive pre-computation/maintenance required by HNSW.