# Vector Database Indexing Benchmark

## Weaviate HNSW vs Exact Baseline

This report presents a comparative analysis of approximate
and exact vector indexing strategies under realistic dataset scale.

## Dataset Description

- Dataset: MuskumPillerum / General-Knowledge
- Document count: ~1,000 documents
- Query count: ~1,000 questions
- Domain: Open-domain general knowledge

Empty and invalid entries were filtered during preprocessing to ensure robust embedding generation.
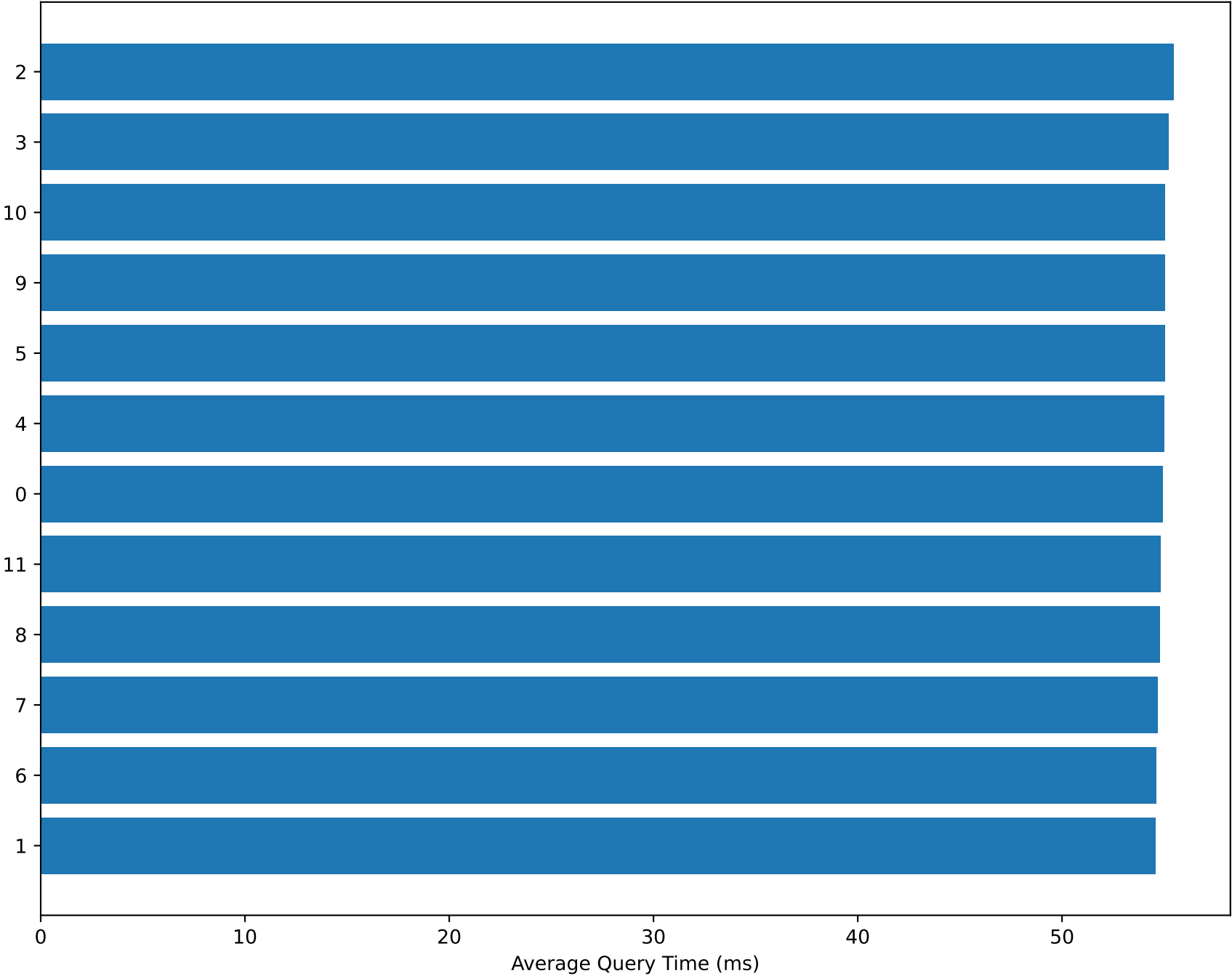
1. Exact Baseline Index
   - Brute-force cosine similarity
   - Provides upper bound on recall
   - Computationally expensive
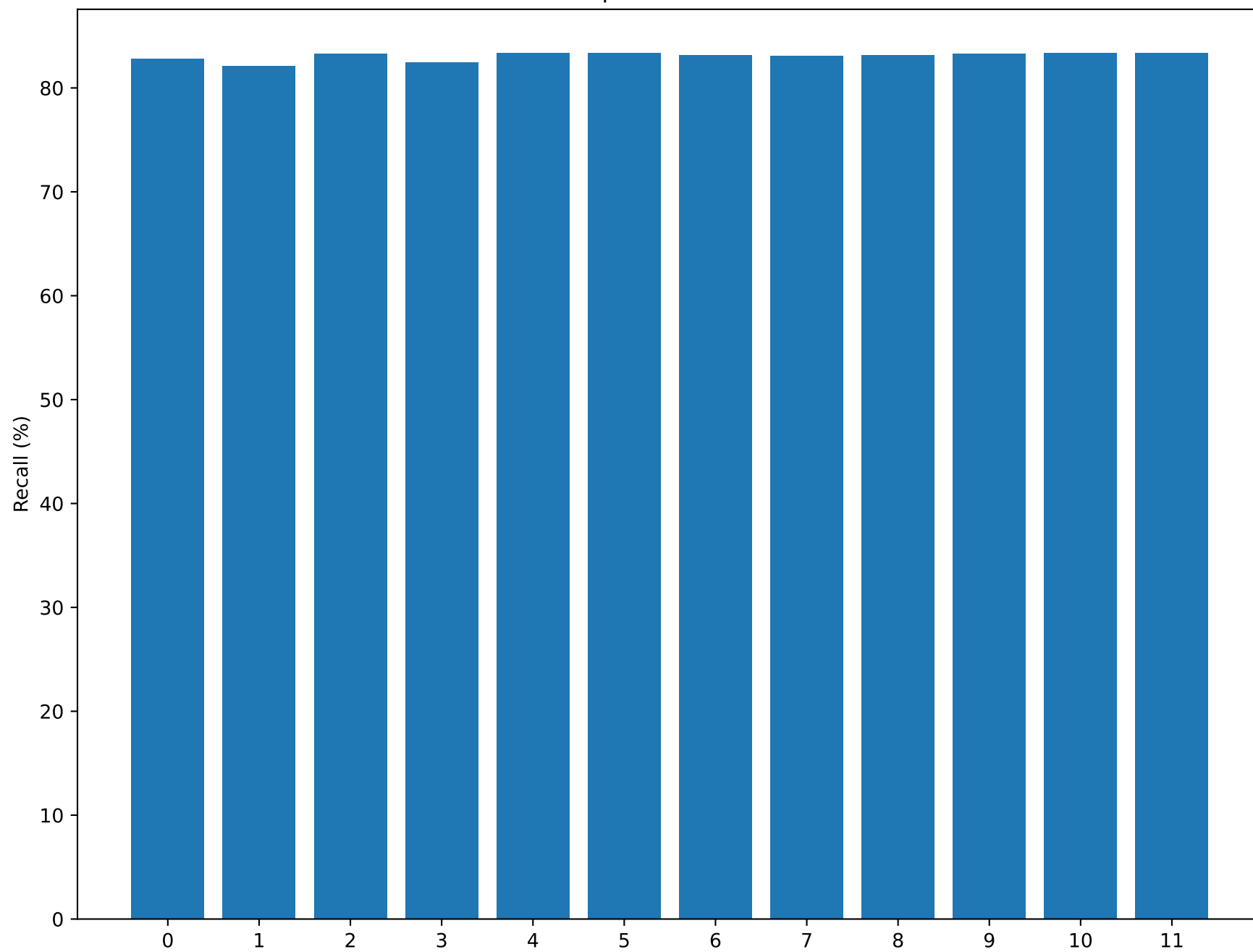
## Indexing Strategies Evaluated

2. HNSW (Weaviate)
   - Approximate nearest neighbour graph
   - Tunable parameters: ef, efConstruction, maxConnections
   - Optimised for low-latency retrieval

Query Latency Across Indexing Strategies

Recall Comparison Across Indexes

## Observations and Trade-offs

- Exact search achieves maximum recall at significantly higher latency.
- HNSW configurations offer substantial latency reductions with moderate recall tradeoffs.
- Parameter tuning enables balanced trade-offs depending on workload.

These results demonstrate the necessity of approximate indexing strategies for scalable LLM-backed retrieval systems.