

## Abstract

- Efficient **Retrieval-Augmented Generation (RAG)** relies heavily on fast, scalable vector search. While **Milvus**<sup>[1]</sup> is a leading solution, its performance on Arm CPUs is often limited by suboptimal utilization of advanced SIMD (Single Instruction, Multiple Data) vector instructions.
- Our work presents a robust solution: Arm Scalable Vector Extension - **SVE**<sup>[3]</sup> optimized distance kernels integrated directly into **Knowhere**<sup>[2]</sup>, Milvus' core vector search engine.
- This innovation delivers up to **2x faster** retrieval for critical index types, directly translating to significantly higher query throughput, lower retrieval latency, and greater hardware efficiency for RAG workloads. All optimizations are fully upstreamed to the open-source codebase and are readily usable on SVE-capable Arm CPUs.

## Methodology

Milvus<sup>[1]</sup> vector search, fundamental to RAG, relies on index traversal, candidate selection, and **distance computation**. The latter often become the bottleneck for query performance on Arm CPUs, primarily due to underutilized SIMD capabilities for common index types, e.g., IVF\_FLAT, IVF\_SQ8, HNSW, and AutoIndex.

To address this critical limitation for RAG, we have implemented **Arm SVE-optimized L2 and Inner Product distance kernels directly within Knowhere**<sup>[2]</sup>, Milvus' core vector search engine. This targeted integration, **leveraging SVE's advanced vector processing, naturally aligns with the compute-intensive nature of distance evaluation in RAG**.

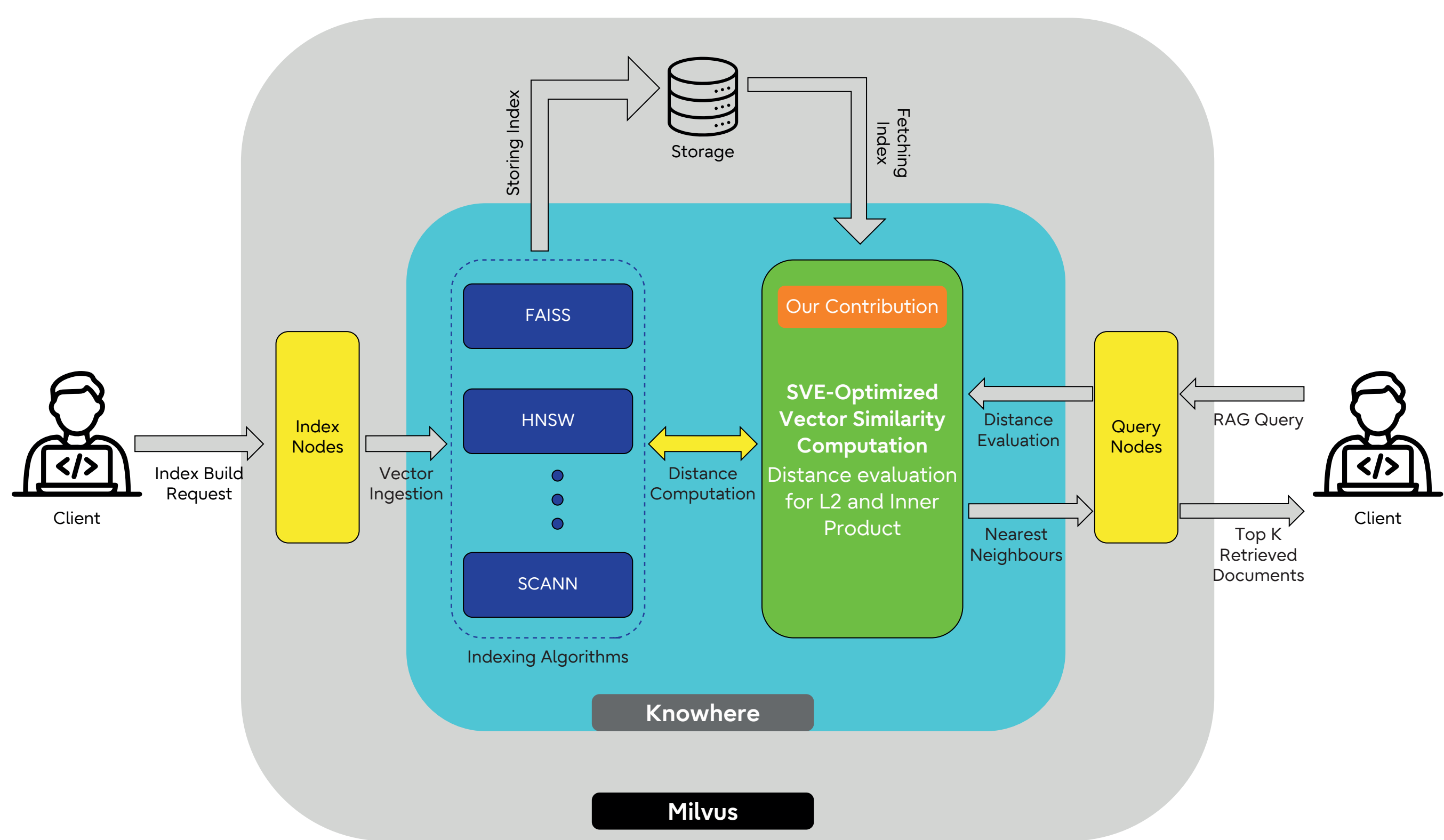


Fig 1: Milvus accelerates RAG retrieval by executing SVE-optimized vector similarity kernels within Knowhere's query and indexing paths [1,2,3].

### SVE-Optimized L2 Distance Computation

#### Inputs:

$x, y$  (FP32 vectors),  $d$  (dimension)

#### Algorithm:

- Initialize vector accumulator  $acc \leftarrow 0$
- Determine SVE vector length  $VL$
- For  $i = 0$  to  $d$  in steps of  $VL$ :
  - (a) Generate predicate mask  $pg$  for valid lanes ( $i < d$ )
  - (b) Load vector elements using predicate
  $x\_vec \leftarrow ldr(x + i, pg)$ 
 $y\_vec \leftarrow ldr(y + i, pg)$
  - (c) Compute vector difference
  $diff \leftarrow x\_vec - y\_vec$
  - (d) Accumulate squared distance using fused multiply-add
  $acc \leftarrow acc + diff \times diff$
- Reduce vector accumulator across lanes
- Return scalar L2 squared distance

#### End Algorithm

**Algorithm 1:** Pseudocode for SVE-optimized<sup>[3]</sup> FP32 L2 distance function. A similar approach has been used to extend SVE support to Inner Product(IP) metric functions and to FP16, BF16, INT8 data types.

This **bottom-up approach elevates data-level parallelism, minimizes per-query compute overhead**, and ensures efficient utilization of Arm vector hardware, all while preserving correctness and portability.

## Results

All benchmarks were performed using **VectorDBBench**<sup>[4]</sup>, the official benchmarking tool for Milvus, on an **AWS Graviton3** instance with **1M vectors**. We evaluated performance on **Cohere Dataset (768D)** and **BioASQ Dataset (1024D)** to capture the performance across varying dimensionalities. For each dataset, we measured **Queries Per Second(QPS)** to evaluate retrieval performance and **Indexing Time** to evaluate indexing efficiency.

We compared the default Arm execution path against the **SVE-optimized**<sup>[3]</sup> implementation, highlighting performance gains in both retrieval and indexing workloads relevant to RAG applications.

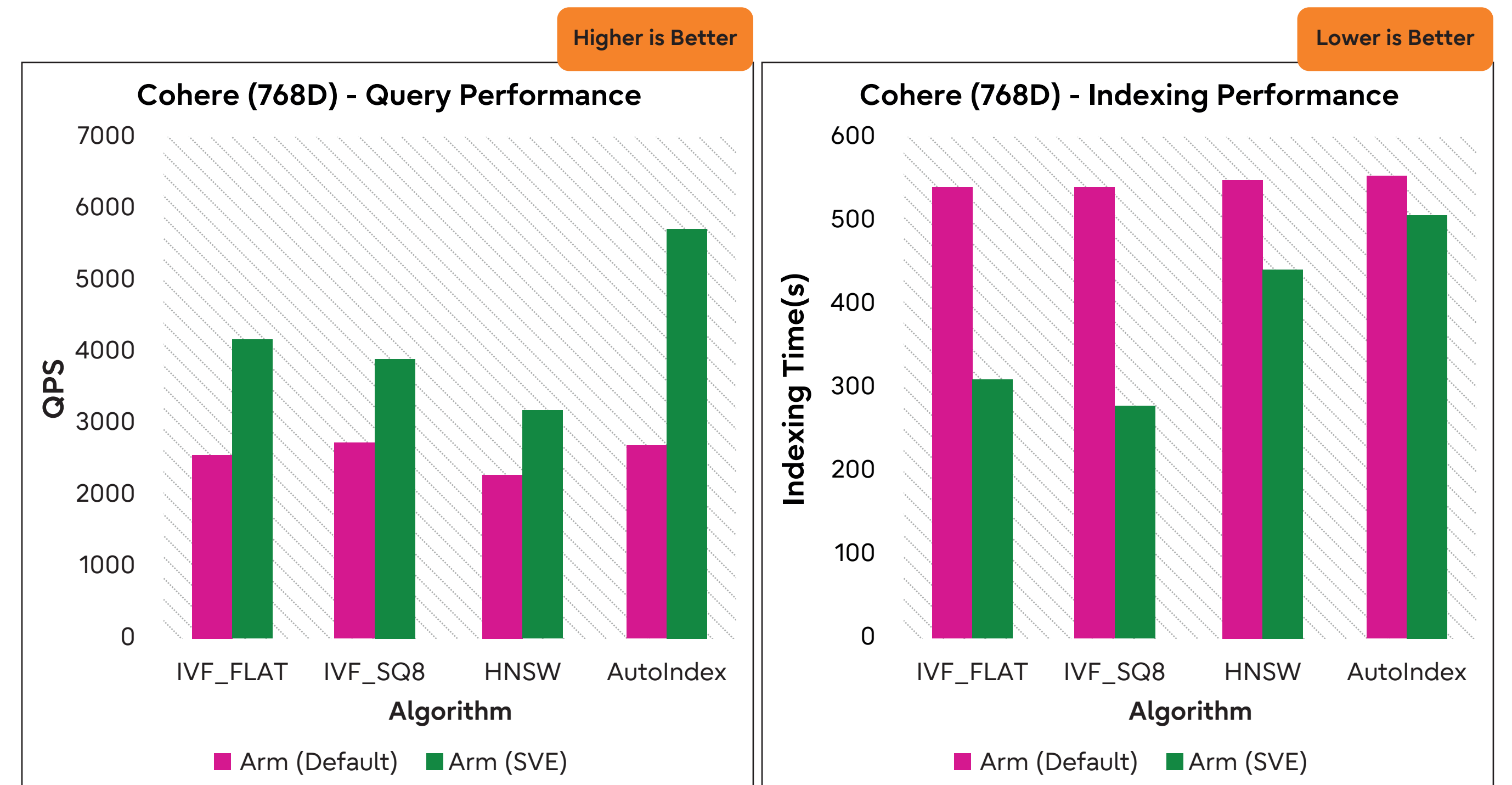


Fig 2: Performance comparison on the Cohere (768D) dataset: (a) query throughput and (b) indexing time between default Arm and SVE-optimized implementations.

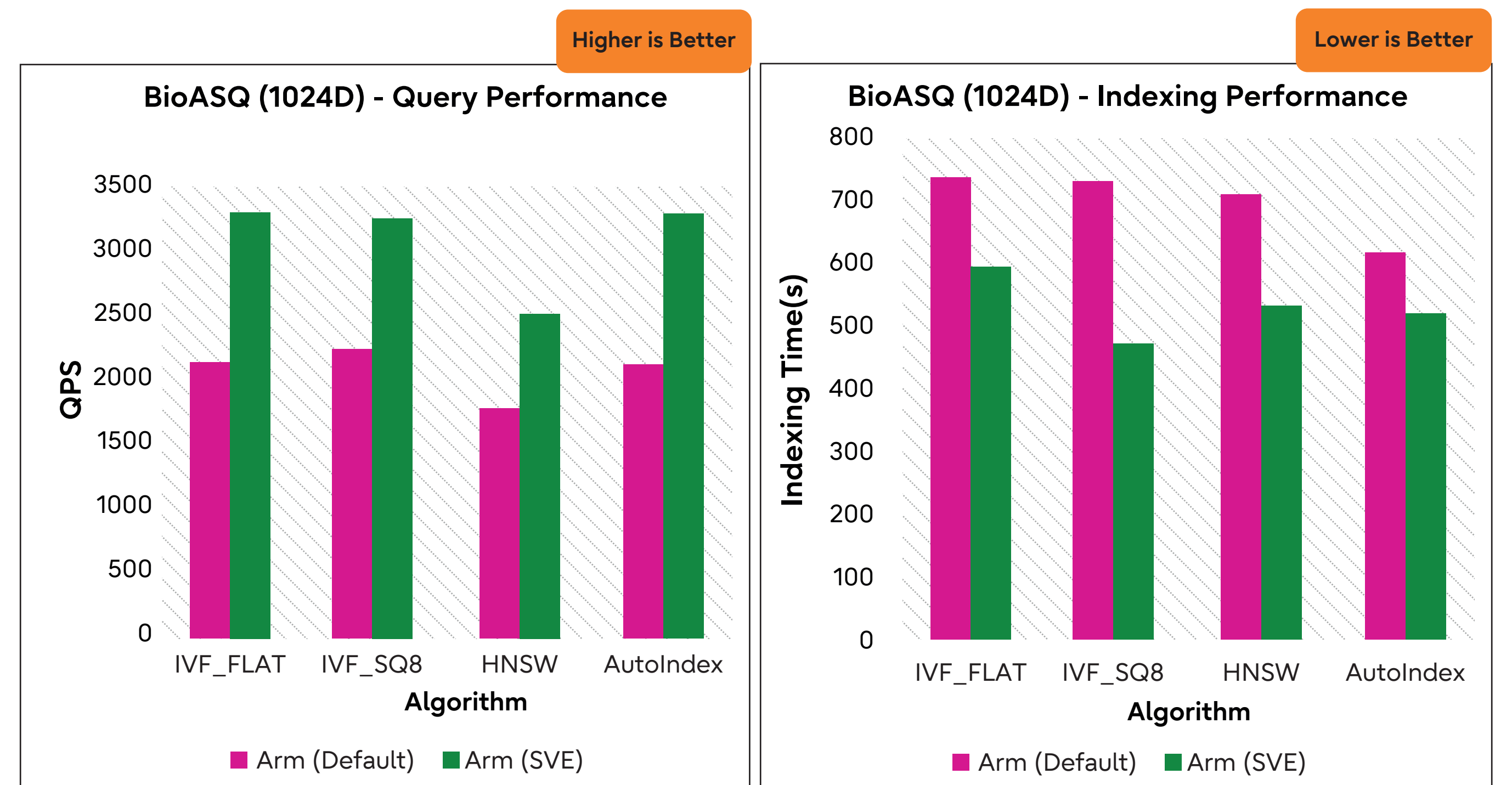


Fig 3: Performance comparison on the BioASQ (1024D) dataset: (a) query throughput and (b) indexing time between default Arm and SVE-optimized implementations.

## Conclusion and Future Work

- We integrated Arm SVE-optimized<sup>[3]</sup> distance kernels into Milvus<sup>[1]</sup>, achieving up to **~2x speedup** (as shown in Figs. 2 and 3) in indexing and query workloads. These optimizations accelerate similarity computation in retrieval-heavy AI and RAG applications.
- While Milvus exposes FP32, FP16, BF16, and binary vectors at the API level, some indexing paths still rely on FP32 internally. We will extend end-to-end SVE support for reduced-precision vectors to eliminate conversions and improve performance on Arm CPUs.

## References

- Milvus: A Vector Database for Similarity Search, [https://milvus.io/]
- Knowhere: Vector Execution Engine, [https://github.com/milvus-io/knowhere]
- Arm SVE Instruction Set Architecture, [https://developer.arm.com/documentation/ddi0596/2021-06/SVE-Instructions]
- VectorDBBench: An Open Source Vector Database Benchmark Tool, [https://github.com/zilliztech/VectorDBBench]

### Links to Pull request



L2-FP32 IP-FP32 L2-FP16 IP-FP16 L2-INT8 IP-INT8 L2-BF16 IP-BF16