GENERATIVE AI: Data Architecture with Google Cloud- AlloyDB & AI

Prasad Venkatachar Sr. Director Products & Solutions @Pliops

COMPUTE, MEMORY, S AND STORAGE SUMMIT

Solutions, Architectures, and Community VIRTUAL EVENT, MAY 21-22, 2024





- Data Platform Considerations
- Gen AI Enterprise Verticals & Horizontal
- Google Cloud AlloyDB Omni Solution
- How do I use it for Business Applications
 - Operational Store, Analytics, Gen Al
- Retrieval Augmented Generation Concepts & Implementation
- AlloyDB AI Features



Next Data Platform Consideration



"71% of respondents in the Data and AI Trends Report plan to use databases integrated with gen AI capabilities."



Gen Al Opportunity Everywhere





One Platform:

Transactions/Analytics/Gen Al

Gen Al

Real Time

E-Commerce Applications



4X PostgreSQL Performance: Transactional Workloads

- 2X Higher Transaction served with 350000 AlloyDB Omni compared to 300000 PostgreSQL
- Up to 4X Transaction served from PostgreSQL to AlloyDB Omni with Pliops & Lenovo
- Serve more Web and Mobile users transaction Requests to meet demand
- Seamlessly scale database while maintaining high performance.





Application User Experience:

Significant Average & Tail Latency Reduction





Analytics: AlloyDB Columnar Engine **Implementation & Benefits**

Real-time business insights



Row vs Columnar Execution



'≤ COMPUTE, MEMORY, द AND STORAGE SUMMIT

Al VS Gen Al Life Cycle

RAG Use cases & Advantages

AlloyDB AI – Vector Database

RAG Architecture

AlloyDB AI Implementation

- Add google_ml_integrations & PgVector
- Register & Integrate the Embedding model endpoint -Preview Mode
- Generate Embedding using embedding() function within AlloyDB
- Store the vector form of data in AlloyDB

CREATE EXTENSION google_ml_integration;

CREATE EXTENSION vector;

GRANT EXECUTE ON FUNCTION emedding to user;

SELECT embedding ('textemedding-geck@001', 'AlloyDB is high performance Postgres database for Enterprises to build Operational, Analytics & Generative AI Applications);

Vector Search

Find Most Similar Embeddings

AlloyDB

Database

14 | ©2024 SNIA. All Rights Reserved.

Source: Deep Learning. Ai

Postgres Vector Indexes

IVFFlat (Inverted File with Flat Compression)

Number/size of the lists

Search: Number of lists to be verified

Hierarchical Navigable Small Worlds(HNSW)

m: Maximum number of connections per layer

Ef_construction : Size of Dynamic list for construction graph

ScaNN Index Benefits over HNSW

Method	ScaNN for AlloyDB	HNSW
Index size	Tree overheads are typically smaller. Quantization further reduces size.	Graphs have more edge connections therefore bigger index overhead.
Index time	Tree training and indexing are faster. The operation is of O(#vectors * #centroids).	Graph construction fundamentally requires many vector–vector comparisons.
Memory access	Vectors in tree leaves are stored contiguously. Memory access is more continuous and friendly to SIMD acceleration.	Graph walk beam search requires random memory accesses. NGT partially solves this by duplicating data in nodes but leads to bloated indices.
Latency	Constant search speed across queries. Typically configured to search a fixed number of leaves.	Varies more per query. Easy queries terminate early as nothing is left in the working set of beam search. But tail latency for harder queries could be higher.

ScaNN : Technology Preview Mode

AlloyDB AI Enhancements

LLM Foundation Models: Vertex AI, Gemma, Open AI

* Technology Preview Mode

Unified Infra & Data Strategy : AlloyDB Omni Solution

Please take a moment to rate this session.

Your feedback is important to us.

COMPUTE, MEMORY,

Solutions, Architectures, and Community VIRTUAL EVENT, MAY 21-22, 2024