



CockroachDB at the Edge of Resilience, Security, and AI

Alicia Lu, Product Manager
Biplav Saraf, Product Manager



PC Computing Era

Mobile Computing Era

AI Integrated Computing Era

1980s/90s

2010s/20s

2020s/30s



Traditional DB



NoSQL/Distributed SQL DB



Vector/Distributed SQL DB

Purpose of Databases has remained **largely the same**,
though the way they **function/operate has changed**

CockroachDB is expanding its deployment model options



Today, we'll cover:

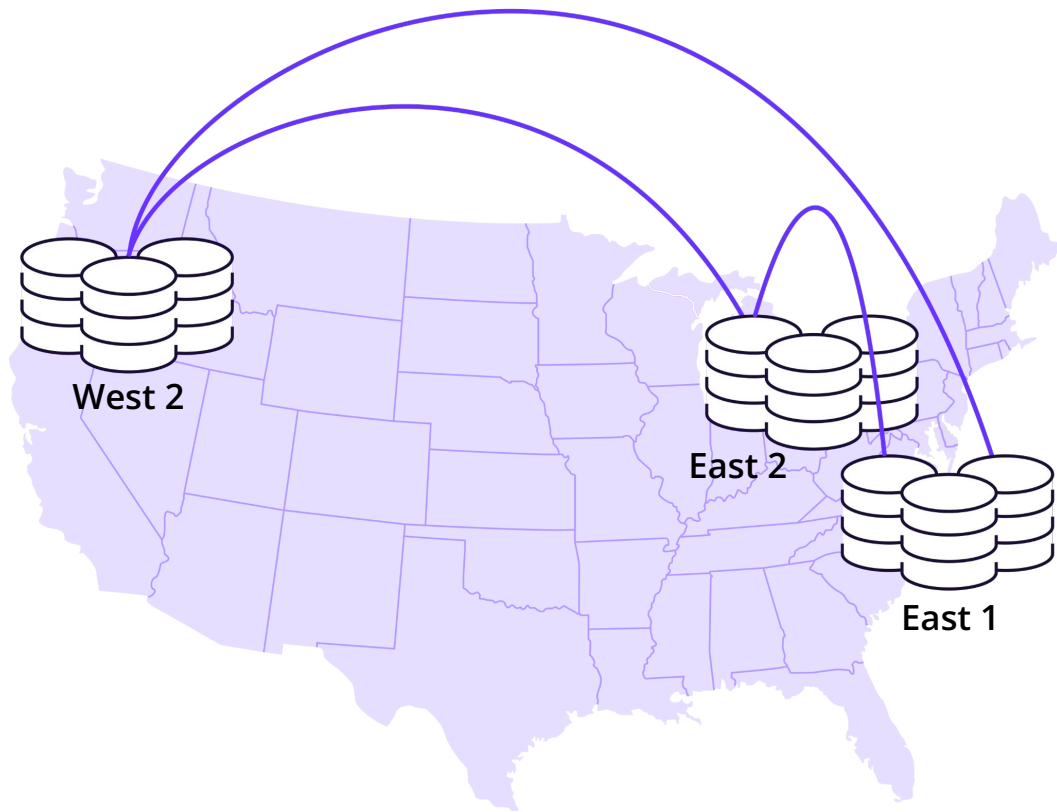


- 1) Cross-Cluster Replication in CockroachDB to meet two data center resiliency
- 2) CockroachDB Security at Scale
- 3) Architecturing CockroachDB for AI Workloads



Resiliency beyond Raft: Cross-Cluster Replication in CockroachDB

CockroachDB enables resilient, fully consistent clusters



Running critical, distributed workloads with serializable isolation

Achieving zero downtime region resiliency with 3+ data centers

Resiliency powered by native Raft replication and distributed SQL

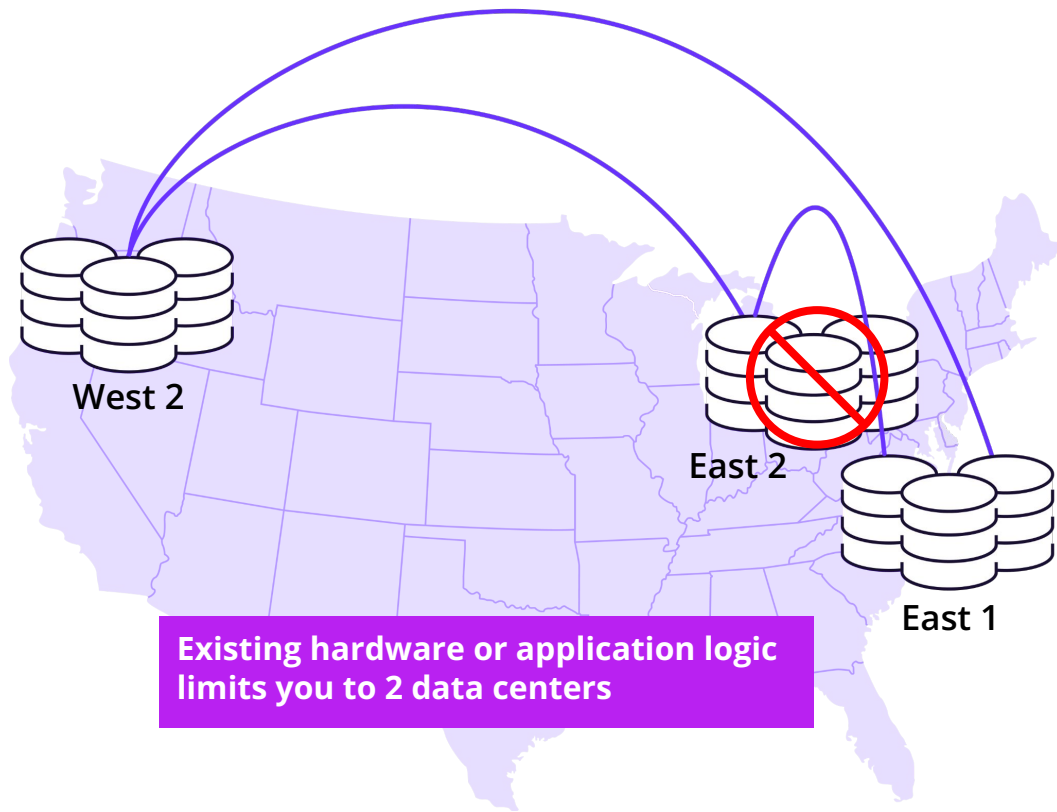


CockroachDB's native consensus-based Raft replication

- Full data consistency across nodes, zones & regions
- Zero data loss during failures within a cluster

Raft Replication

But what if you have other requirements?



Running critical, distributed workloads with serializable isolation

Achieving zero downtime region resiliency with 3+ data centers

You want to survive a region failure with 2 regions

You want to achieve low, single-region write latency

You want to migrate data between clusters

You want more flexibility with deployment topologies

Introducing new tools to meet custom requirements



CockroachDB's native consensus-based Raft replication

- Full data consistency across nodes, zones & regions
- Zero data loss during failures within a cluster

Raft Replication



[GA] Physical Cluster Replication

- Transactionally consistent, full cluster replication
- Ensures data protection and disaster recovery with a passive standby cluster



[GA] Logical Data Replication

- Eventually consistent, table and database level replication
- Enables high availability and real-time analysis with active clusters

**Cross-Cluster
Replication Tools**



Physical Cluster Replication

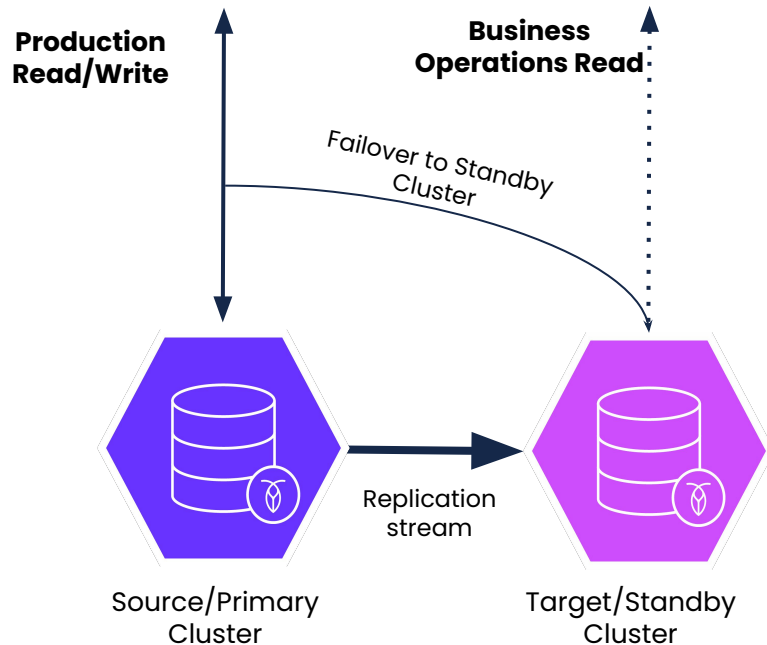
Generally Available in 24.1

Limited Access on Cockroach Cloud

What is Physical Cluster Replication (PCR)?



Physical Cluster Replication continuously replicates all data from a *primary* CockroachDB cluster to an independent *standby* CockroachDB cluster



Use cases



Region and Cloud Survivability

Use PCR to survive an regional, cluster, or cloud outage with low RPO and RTO by cutting over to the standby.

Workload Isolation

Use PCR to offload non-critical operations, such as running analytics queries, to the standby cluster (24.3)

Data Migration

Use PCR to migrate data seamlessly between CockroachDB clusters.

- Cluster upgrade protection
- Populate dev clusters in real time

Physical Cluster Replication (PCR) overview

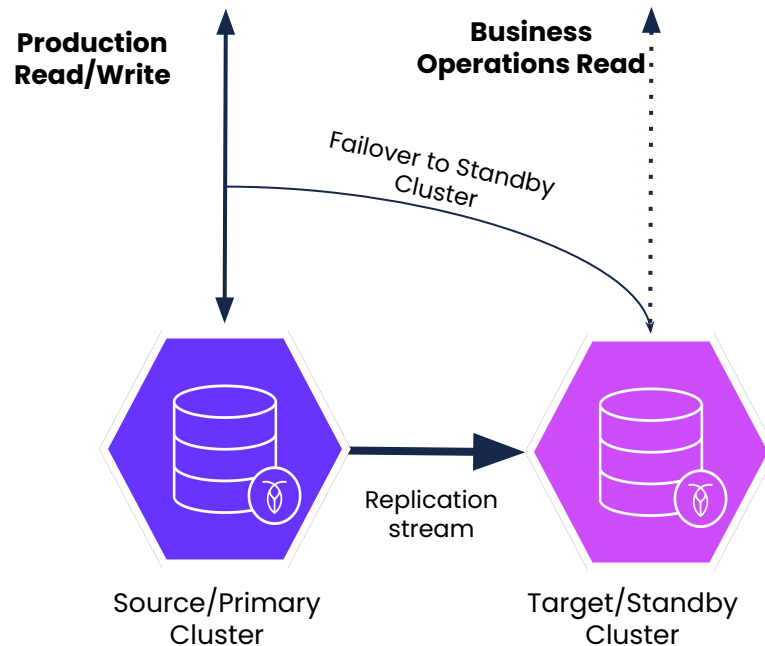


Cluster-level

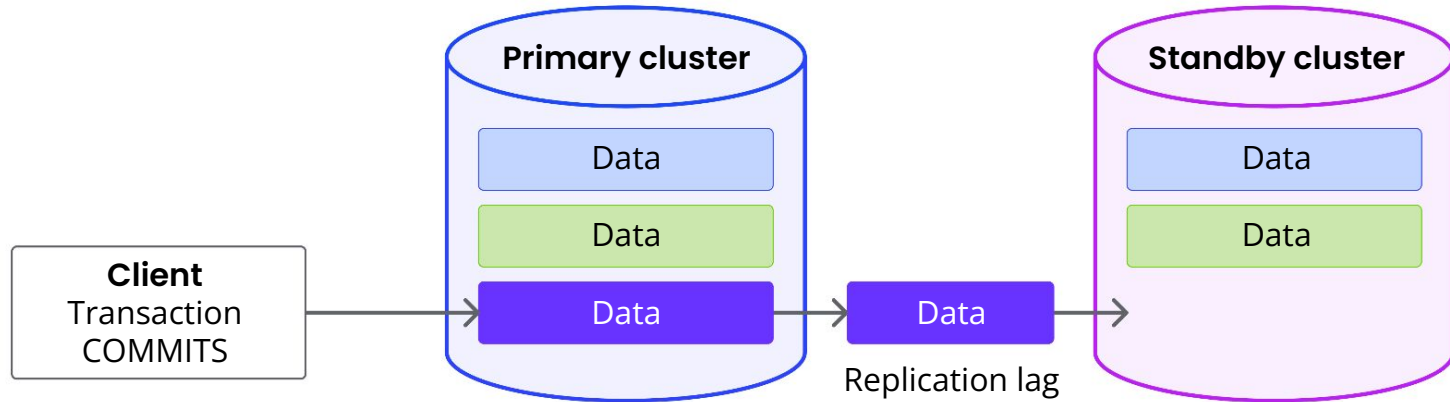
Replication happens at cluster-level, i.e., it is unaware of database, table, row boundaries, etc.

Transactionally Consistent

You see a transactionally consistent state when you “failover” to the standby, including schema changes, zone configs, users, privileges, etc. No conflicts.

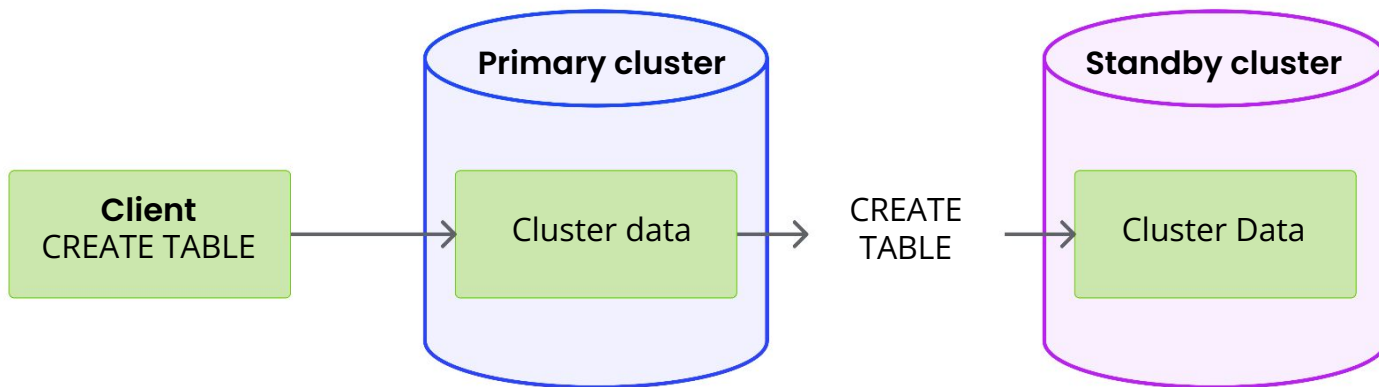


Life of a transaction in PCR

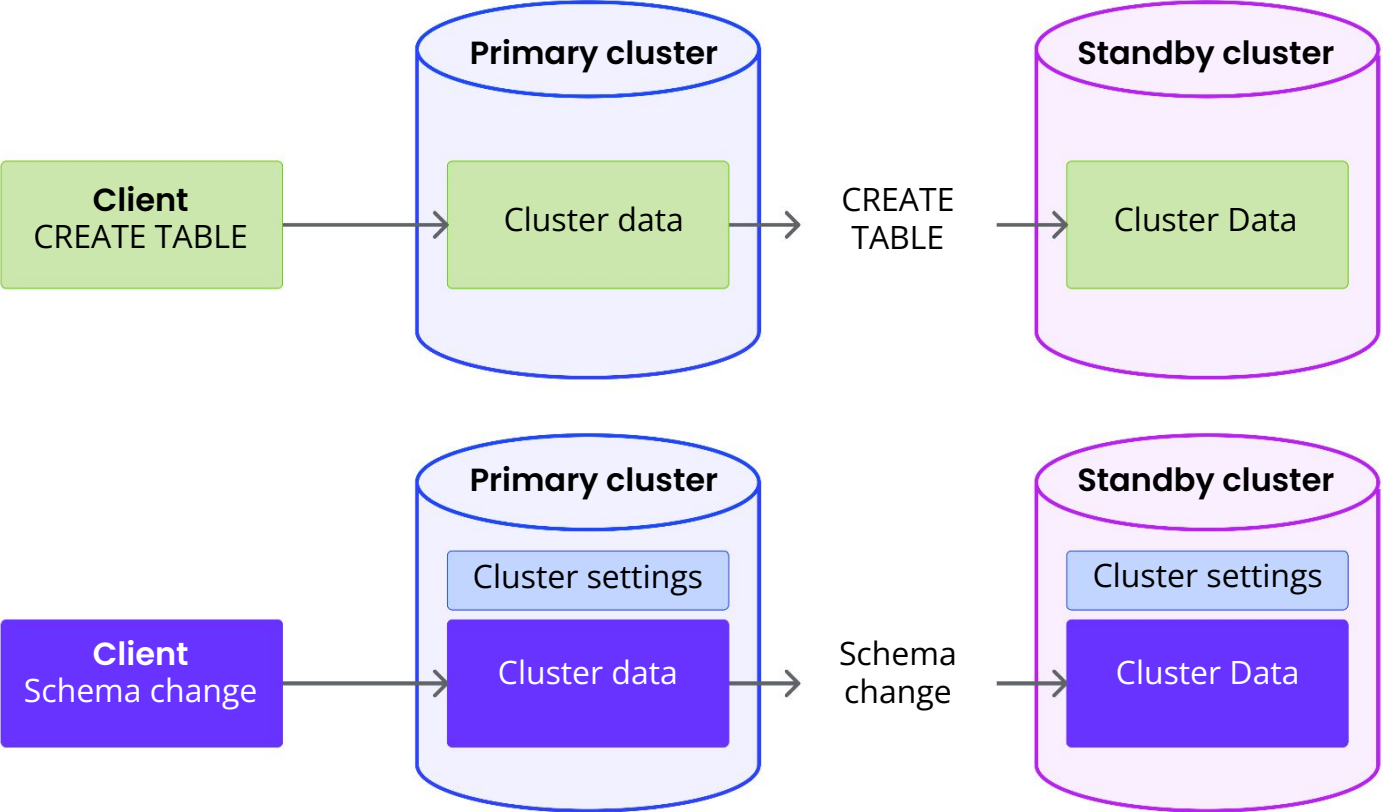


All data is accounted for in PCR during replication.

PCR replicates the entire cluster



PCR replicates the entire cluster



How does PCR handle disasters?



Node Outage, Network Partition

PCR will replan node connections between clusters and continue replication without any manual intervention required.



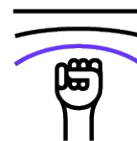
Cluster, Data-center, Region, Cloud Outage

Users can cutover to the standby cluster with a ~30 second RPO. Users will manually redirect application traffic to the standby cluster.



Human Error, dropping a table

PCR lets you cutover to a **time in the past**, ensuring that the target cluster is a consistent copy of the source cluster at the chosen timestamp.



Cyber-attack

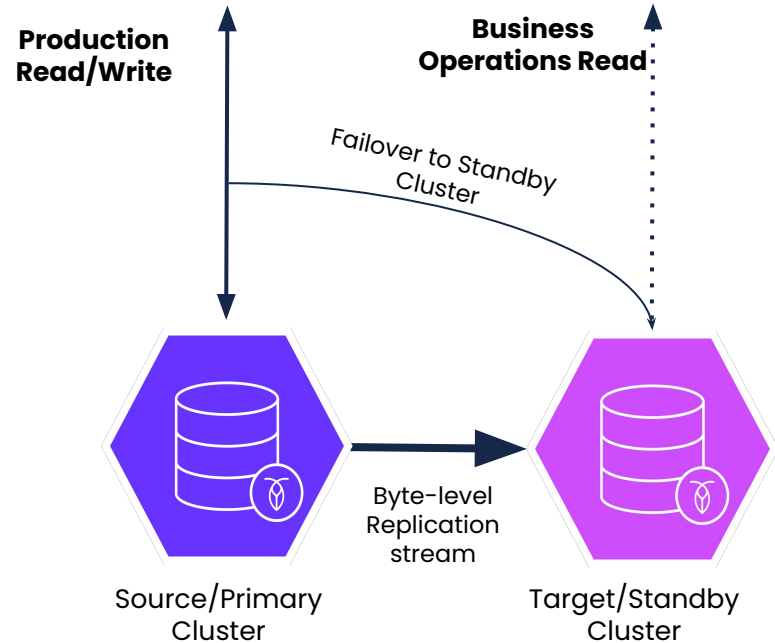
PCR lets you cutover to a **time in the past**, ensuring that the target cluster is a consistent copy of the source cluster at the chosen timestamp.

Physical Cluster Replication (PCR) overview



RPO is ~30 seconds for most applications, depending on the workload traffic and corresponding replication lag.

RTO is ~one minute limited by the time it takes to identify the incident, decide and complete the cutover from the primary to standby database, & promote the standby to primary within your load balancer.

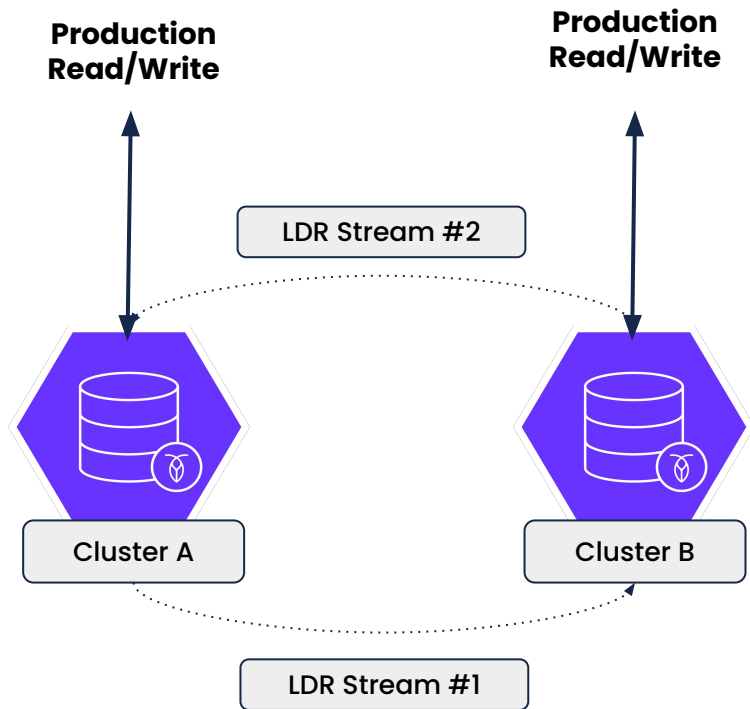




Logical Data Replication

GA in 25.2

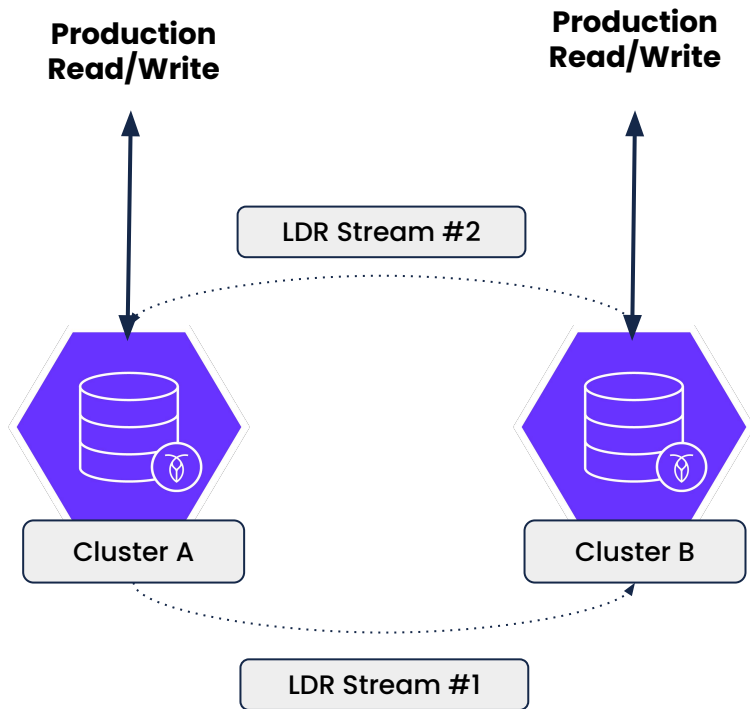
What is Logical Data Replication (LDR)?



LDR provides **ACID consistency** on **each cluster** and **asynchronous replication** between clusters so that users can achieve **eventual consistency** between 2+ clusters.

Reads and writes can occur concurrently on both clusters allowing customers to survive a region failure while providing low latency read and writes.

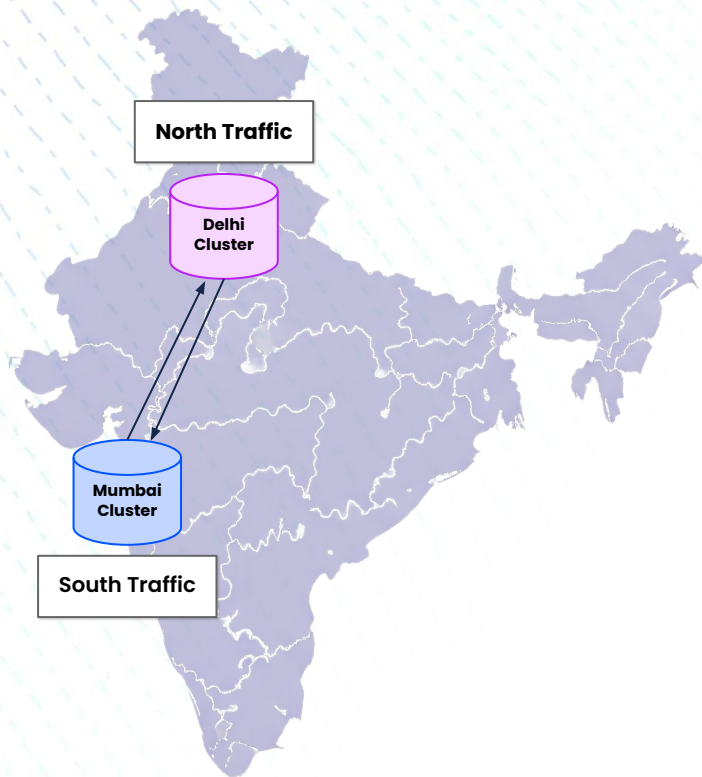
What is Logical Data Replication (LDR)?



LDR replicates **table subsets** of a cluster.

LDR uses **last writer wins** conflict resolution based on the MVCC timestamps of the replicating write.

Achieving high availability and low latency with 2 data centers

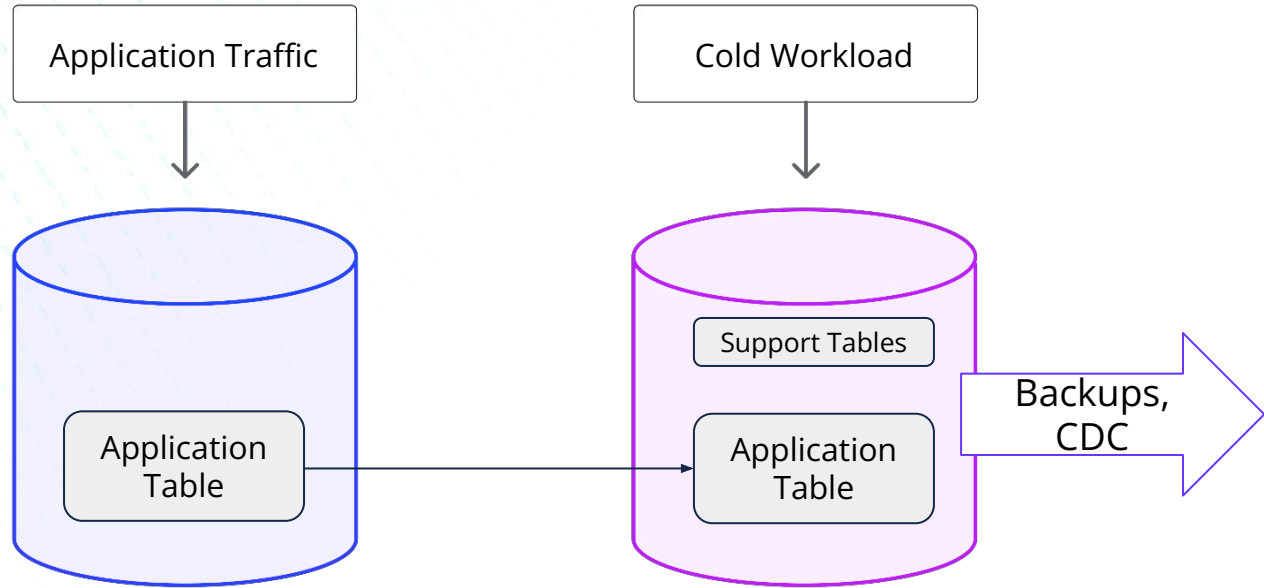


- Both clusters can receive application reads and writes
- Both clusters have low, single-region write latency with transactionally consistent writes via Raft consensus replication
- In a data center or cluster outage, operators can redirect application traffic to the surviving cluster with little downtime

Workload isolation



LDR enables users to isolate critical application workloads from non-critical application workloads, enabling users to isolate resources and size clusters for certain workloads.



Cross-Cluster Replication in CockroachDB



	Physical Cluster Replication (PCR)	Logical Data Replication (LDR)
What is it?	<ul style="list-style-type: none">• One-way physical replication• Replicates the entire source cluster	<ul style="list-style-type: none">• Bi-directional logical replication• Replicate table subsets of a cluster
Why would you use it?	<ul style="list-style-type: none">• Surviving region/DC outages with 2 DCs• Surviving cluster control plane outages• Data migration between clusters• Create flexible deployment topologies	
Operator Experience	<ul style="list-style-type: none">• Started by SQL statements, in the DB• Jobs infrastructure integration• Admission Control for protecting foreground workloads	



**Modern applications
demand trust:
secure everywhere
and at scale**

Security starts where the data lives



Contain, recover, evolve: resilience against threats

CockroachDB goals



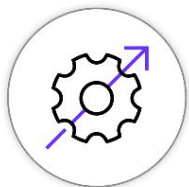
Anticipate

Expect the unexpected



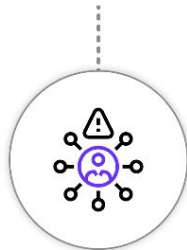
Withstand

Keep calm and keep operating



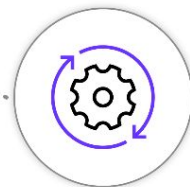
Recover

Get back up faster



Contain

Limit the damage quickly

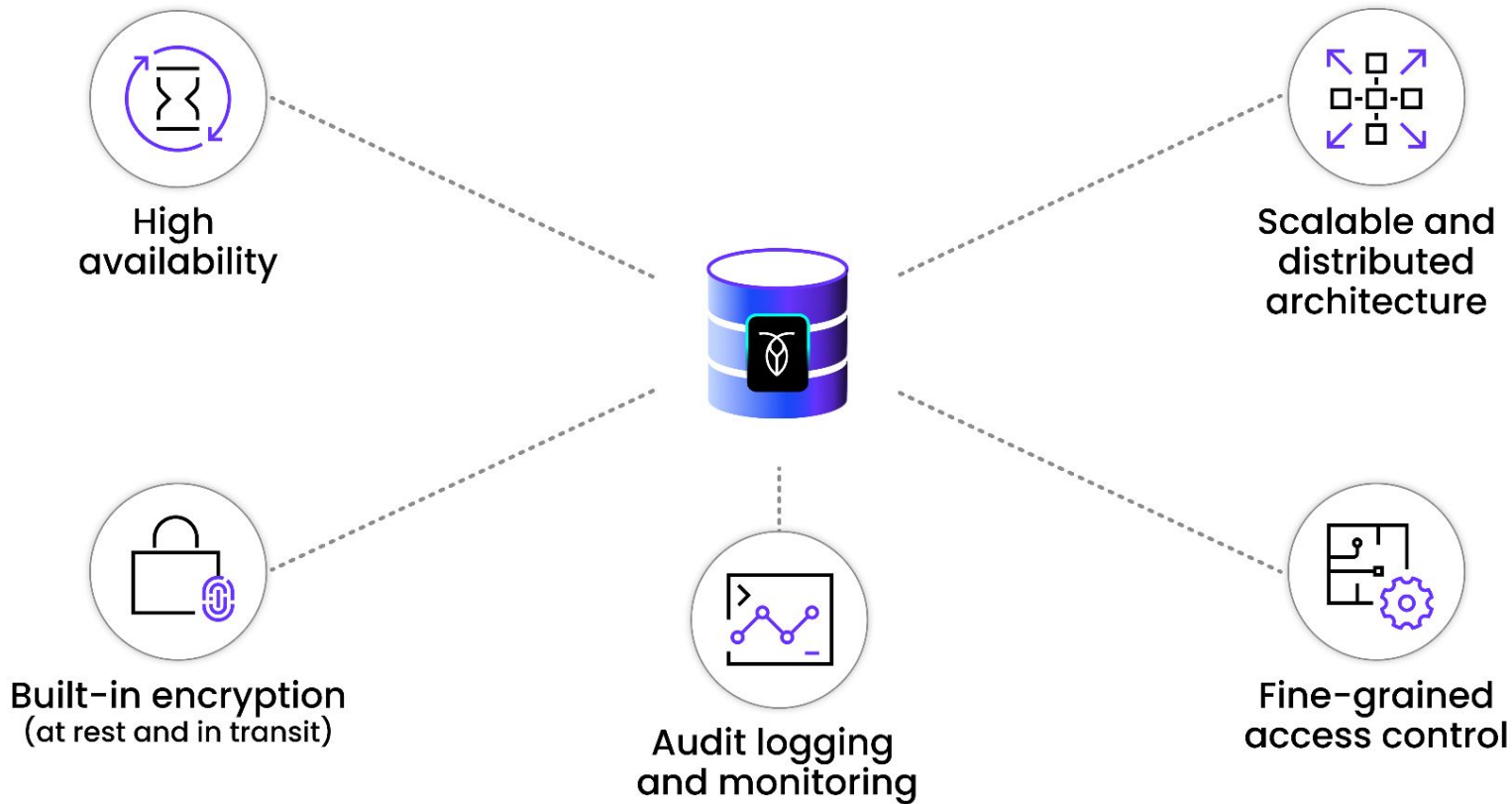


Evolve

Always stay ahead

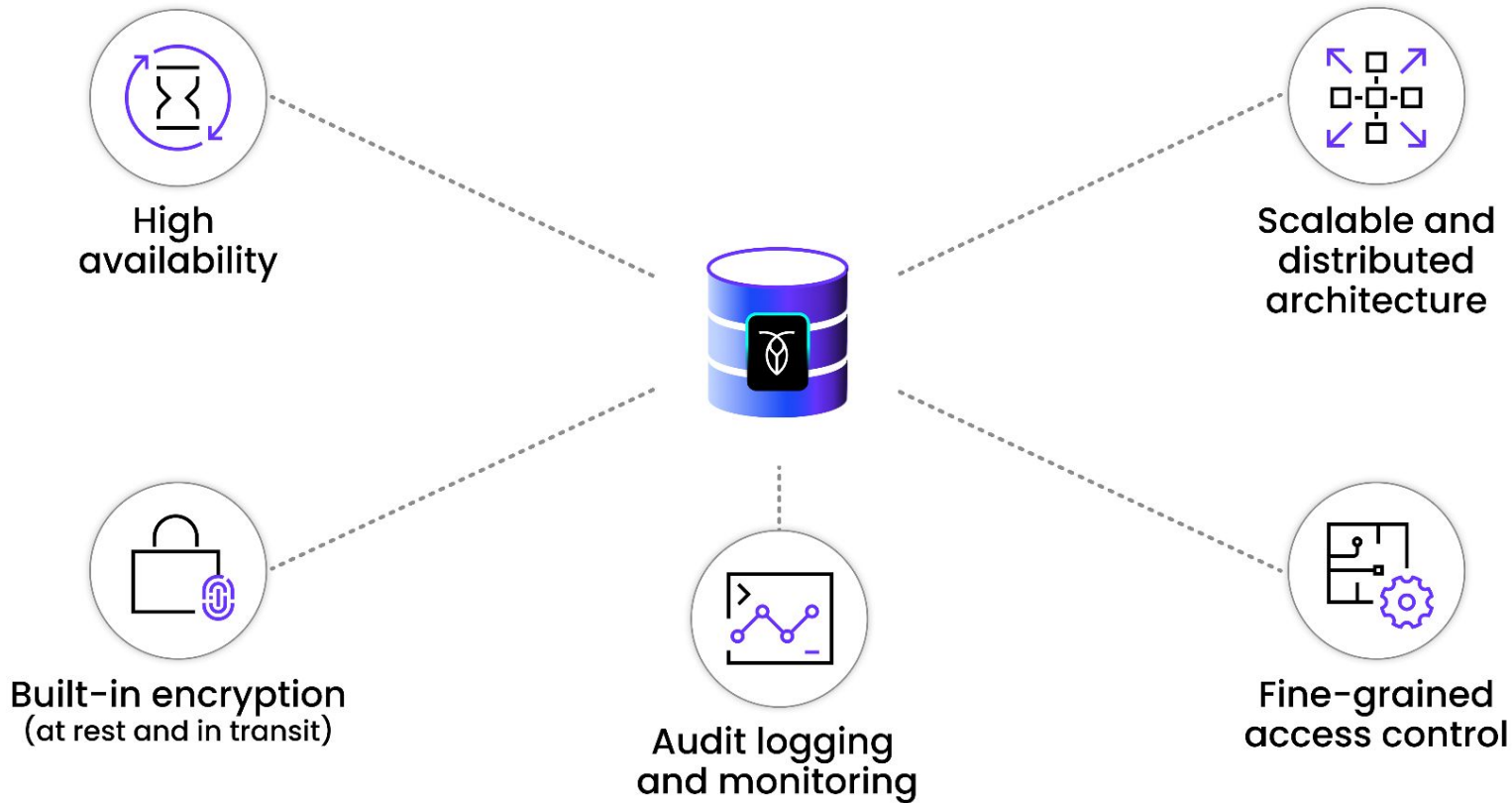
Security and resilience built into the database

CockroachDB features



Security and resilience built into the database

CockroachDB features





**Compliance at
scale: regional,
tenant-aware,
secure**

Row Level Security provides database native advanced security and compliance controls

✦ SHOW DON'T TELL

CockroachDB

Row-Level Security



Rob Reid

Technical Evangelist at
Cockroach Labs



twitter.com/robreid_io



youtube.com/cockroachdb



Architecting for AI: Evolving Database Design for the AI Era

AI memory: the agent's brain

Short-Term Memory (context window)

What it is: *The AI's "working memory" for the current conversation.*

Analogy: *An AI's immediate consciousness.*

Key Traits: *Volatile, limited capacity, and extremely fast.*

AI memory: the agent's brain

Short-Term Memory (context window)

What it is: *The AI's "working memory" for the current conversation.*

Analogy: *An AI's immediate consciousness.*

Key Traits: *Volatile, limited capacity, and extremely fast.*

Long-Term Memory (knowledge base)

What it is: *The permanent store for learning, personalization, and reasoning.*

Episodic: *Remembers specific past interactions. "Last time you asked about billing."*

Semantic: *General knowledge of facts and concepts. "What is a distributed database?"*

Procedural: *Knows how to perform tasks. "How do I book a flight?"*

AI memory: the agent's brain

Short-Term Memory (context window)

What it is: *The AI's "working memory" for the current conversation.*

Analogy: *An AI's immediate consciousness.*

Key Traits: *Volatile, limited capacity, and extremely fast.*

Long-Term Memory (knowledge base)

What it is: *The permanent store for learning, personalization, and reasoning.*

Episodic: *Remembers specific past interactions. "Last time you asked about billing."*

Semantic: *General knowledge of facts and concepts. "What is a distributed database?"*

Procedural: *Knows how to perform tasks. "How do I book a flight?"*

Incidental Memory (governance)

What it is: *The unintentional capture of irrelevant or sensitive information (e.g., personal details shared in a conversation).*

Analogy: *Remembering the color of a passing car you weren't trying to memorize.*

Key Challenge: *This is not a feature, but a governance and security problem. The goal is to prevent storing this data to protect user privacy.*

AI memory: the agent's brain

Short-Term Memory (context window)

What it is: *The AI's "working memory" for the current conversation.*

Analogy: *An AI's immediate consciousness.*

Key Traits: *Volatile, limited capacity, and extremely fast.*

Long-Term Memory (knowledge base)

What it is: *The permanent store for learning, personalization, and reasoning.*

Episodic: *Remembers specific past interactions. "Last time you asked about billing."*

Semantic: *General knowledge of facts and concepts. "What is a distributed database?"*

Procedural: *Knows how to perform tasks. "How do I book a flight?"*

Incidental Memory (governance)

What it is: *The unintentional capture of irrelevant or sensitive information (e.g., personal details shared in a conversation).*

Analogy: *Remembering the color of a passing car you weren't trying to memorize.*

Key Challenge: *This is not a feature, but a governance and security problem. The goal is to prevent storing this data to protect user privacy.*

Contextual Memory (orchestrator)

What it is: *The effective blending of short-term and long-term memory to maintain coherent understanding across interactions.*

Function: *Ensures the AI is consistent and applies past knowledge to current conversations.*

AI memory: the agent's brain

Short-Term Memory (context window)

What it is: *The AI's "working memory" for the current conversation.*

Analogy: *An AI's immediate consciousness.*

Key Traits: *Volatile, limited capacity, and extremely fast.*

Long-Term Memory (knowledge base)

What it is: *The permanent store for learning, personalization, and reasoning.*

Episodic: *Remembers specific past interactions. "Last time you asked about billing."*

Semantic: *General knowledge of facts and concepts. "What is a distributed database?"*

Procedural: *Knows how to perform tasks. "How do I book a flight?"*

Incidental Memory (governance)

What it is: *The unintentional capture of irrelevant or sensitive information (e.g., personal details shared in a conversation).*

Analogy: *Remembering the color of a passing car you weren't trying to memorize.*

Key Challenge: *This is not a feature, but a governance and security problem. The goal is to prevent storing this data to protect user privacy.*

Contextual Memory (orchestrator)

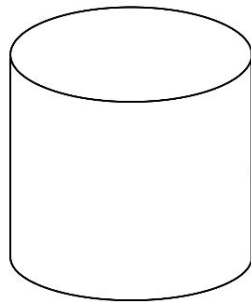
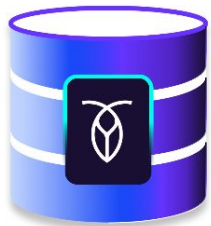
What it is: *The effective blending of short-term and long-term memory to maintain coherent understanding across interactions.*

Function: *Ensures the AI is consistent and applies past knowledge to current conversations.*

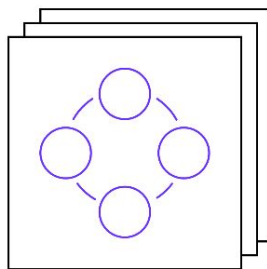
LLMs are Dumb without it's brain!

CockroachDB vector data and indexing powers native AI RAG workloads

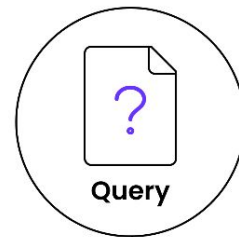
CockroachDB fits here as a knowledge base



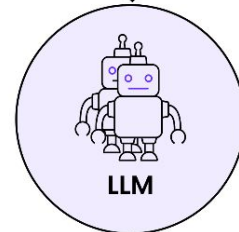
Knowledge base



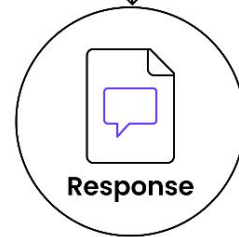
Relevant context



Query



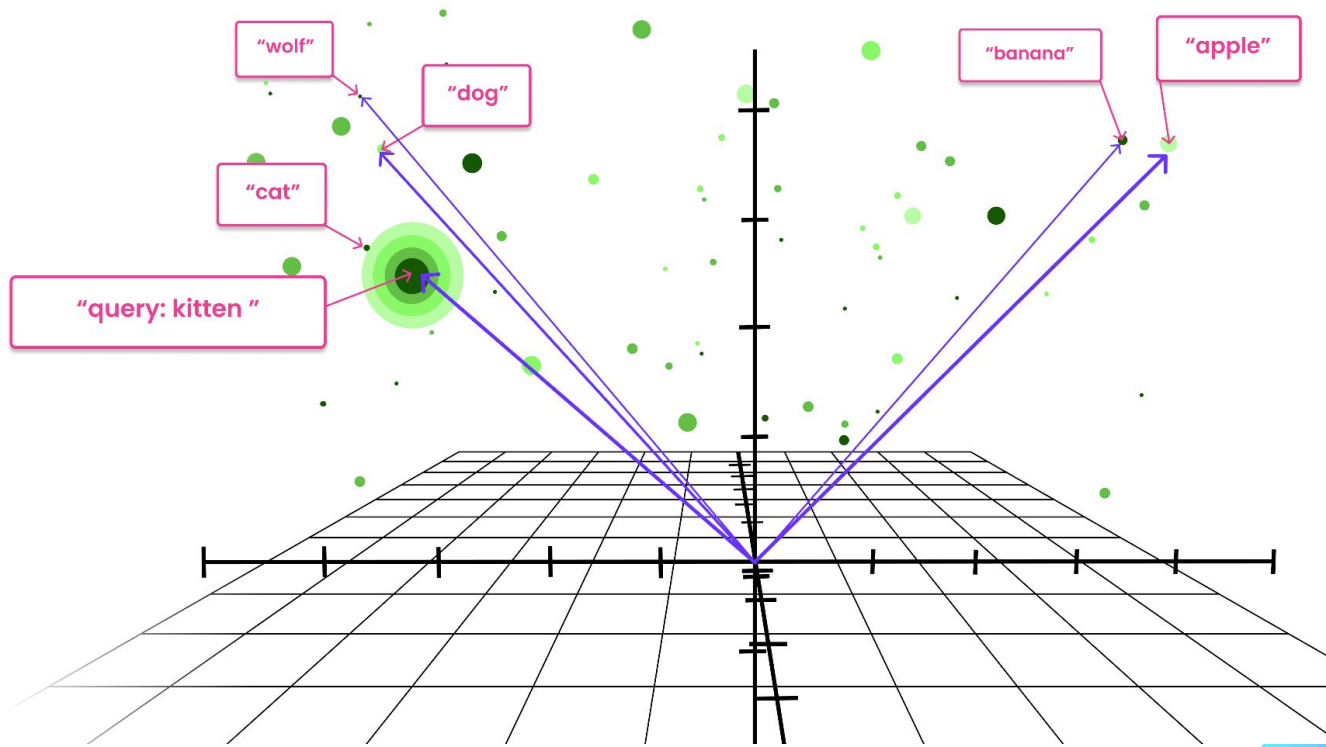
LLM



Response

CockroachDB Vector indexing provides automatically up-to-date + partitioned searches

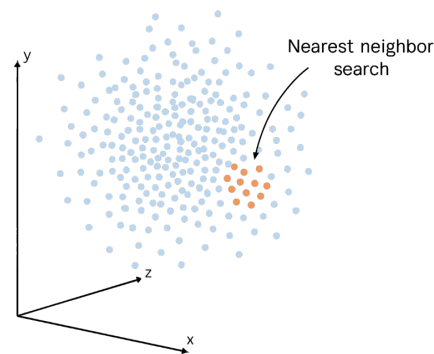
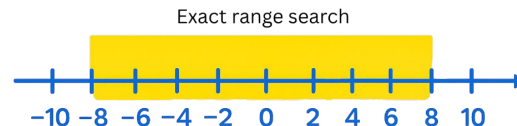
What is a vector?



Why you need a vector index



- Regular indexes work for single-dimensional, ordered data types like strings or numbers
- Embeddings are high-dimensional and don't have a natural sort order
- We care about **semantic similarity**, not exact match
- Brute-force scanning works at small scale, but falls apart with millions of vectors
- Vector indexes enable fast, approximate nearest neighbor (ANN) search **at scale**



The challenge: distributed vector indexing



- No central coordinator
- No large in-memory structures
- Minimal network hops
- Sharding-compatible layout
- No hot-spots
- Incremental updates





The solution: introducing C-SPANN

- Vector indexing algorithm designed specifically for CockroachDB
- Based on **SPANN** and **SPFresh** papers from Microsoft
- **Adapted** for a distributed SQL architecture
- Data quantization done with **RaBitQ** to speed up performance

arXiv:2111.08566v1 [cs.DB] 5 Nov 2021

SPANN: Highly-efficient Billion-scale Approximate Nearest Neighbor Search

Qi Chen¹, Bing Zhou^{1,2}, Haining Wang¹, Mingqiu Li¹, Chuanqi Liu^{1,3},
Zongsheng Li¹, Mao Yang¹, Jingdong Wang^{1,4,5},
Microsoft, Peking University, Tencent, Tsinghua
[qchen1, bingzhou, mingqiu1, janz1, mingqiu@tencent.com,
liz1, hainingwang@pku.edu.cn, liz1, chenq1@outlook.com,
wangjingdong@outlook.com]

Abstract

The in-memory algorithms for approximate nearest neighbor search (ANNS) have achieved great success for fast high-recall search, but are extremely expensive when handling very large-scale datasets. Thus, there is an increasing request for the hybrid ANNS solutions with small memory and inexpensive solid-state drive (SSD). In this paper, we present a simple but efficient memory-disk hybrid indexing and search system, named SPANN, that follows the inverted index technology. It stores the central portion of the posting lists in the memory and the large posting lists in the disk. We guarantee both disk access efficiency (low latency) and high recall by effectively reducing the disk access number and retrieving high-quality posting lists. In the index-building stage, we adopt a hierarchical balanced clustering algorithm to balance the length of posting lists and segment the posting list by adding the points on the cluster of the corresponding clusters. In the search stage, we use query-aware scheme to dynamically prune the access of unnecessary posting lists. Experiment results demonstrate that SPANN is 2-4 faster than the state-of-the-art ANNS solution DiskANN to reach the same recall quality. 80% with same memory cost in three billion-scale datasets. It can reach 90% recall@1 and recall@10 in around one millisecond with only 2GB memory cost. Code is available at <https://github.com/microsoft/SPANN>.

1 Introduction

Vector nearest neighbor search has played an important role in information retrieval area, such as multimedia search and web search, which provides relevant results by searching vectors with minimum distance to the query vector. Exact solutions for K -nearest neighbor search [49, 40] are not applicable in big data scenarios due to substantial computation cost and high queries. Therefore, researchers have proposed many kinds of approximate nearest neighbor search (ANNS) algorithms in the last decade [1, 10, 29, 28, 14, 11, 13, 12, 15, 16, 17, 18, 31, 27, 9, 19, 30, 20, 36]. However, most of the algorithms mainly focus on how to do low latency and high recall search all in memory with different ways [1, 10, 29, 28, 14, 11, 13, 12, 15, 16, 17, 18, 31, 27, 9, 19, 30, 20, 36]. However, most of the algorithms mainly focus on how to do low latency and high recall search all in memory with different ways [1, 10, 29, 28, 14, 11, 13, 12, 15, 16, 17, 18, 31, 27, 9, 19, 30, 20, 36]. However, most of the algorithms mainly focus on how to do low latency and high recall search all in memory with different ways [1, 10, 29, 28, 14, 11, 13, 12, 15, 16, 17, 18, 31, 27, 9, 19, 30, 20, 36].

There are only a few approaches working on the hybrid ANNS solutions, including DiskANN [39] and HNSW-LSH [64]. Both of them are graph based solutions. DiskANN uses Proton Quantization

¹Corresponding author.

²Work done while at Microsoft.

SPFresh: Incremental In-Place Update for Billion-Scale Vector Search

Yuming Xu^{1,2}, Hengyu Liang^{1,2}, Jin Li^{1,2},
Cheng Li¹, Ziyue Yang¹, Fan Yang¹,
Shuotao Xu¹, Qi Chen¹, Qianxi Zhang¹,
Microsoft Research Asia, ²Harvard University

Abstract
Approximate Nearest Neighbor Search (ANNS) on high-dimensional vector data is now widely used in various applications, including information retrieval, question answering, and recommendation. As the amount of vector data grows continuously, it becomes important to support updates to vector index, the enabling technique that drives for efficient and accurate ANNS on vectors.

Because of the curse of high dimensionality, it is often easy to identify the right neighbors of a new vector, a necessary process for index update. To amortize update costs, existing systems maintain a secondary index to accelerate updates, which are merged with the main index to globally re-ranking the entire index periodically. However, this approach has high fragmentation of search latency and accuracy, not to mention that it requires substantial resources and is extremely time-consuming to rebuild.

We introduce SPFresh, a system that supports in-place vector updates. At the heart of SPFresh is LIRE, a lightweight incremental re-ranking protocol to split vector partitions and merge vectors in the nearby partitions to adapt to data distribution shifts. LIRE achieves low overhead vector updates by only re-ranking vectors at the boundary between partitions, where the high-quality vector index for the current partition is deemed small. With LIRE, SPFresh processes update queries along with accuracy reduction based on global rebuild, with only 1% of DRAM less than 10x faster rebuild of the peak compared to the index-of-the-art, a 1-billion-scale disk-based vector index with a 1% of daily vector update rate.

CCS Concepts • Information systems • Information systems; Vector Search; Incremental Update; Billion-scale

Keywords: Vector Search; Incremental Update; Billion-scale

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are made for non-commercial use and that the copies are made without charge to the original publisher. For all other uses, permission should be sought from ACM. Copyright 2021 ACM 978-1-4503-8326-7/21/000000-0000 \$5.00.

ACM Reference Format:

Yuming Xu, Hengyu Liang, Jin Li, Cheng Li, Ziyue Yang, Fan Yang, Shuotao Xu, Qi Chen, Qianxi Zhang, Microsoft Research Asia, Harvard University

1 Introduction

Today deep learning models can embed almost all types of data, including speech, vision, and text information, into multi-dimensional vectors with tens or even hundreds of dimensions. Such vectors are critical for complex semantic understanding tasks [11, 47]. To enable effective vector analysis, vector nearest neighbor search (ANNS) systems have become critical system components for an increasing number of online services like search [5] and recommendation [7].

To satisfy the strict query latency requirement for these online services, vector search systems often resort to approximate nearest neighbor search (ANNS) [11, 22, 34, 31, 15, 16, 40, 42] to boost as many correct results as possible (i.e., query accuracy). At the heart of a large-scale ANNS system is a vector index, a key data structure that organizes high-dimensional vectors efficiently for high-accuracy low-latency vector searches [15, 16, 34, 31, 42].

Like traditional indices, a high-quality vector index requires quick "repartitioning" of vectors based on the vector proximity in a high-dimensional space. The proximity measurements are often implemented with "sharding", which only cuts between a pair of vectors with a slight distance. A search query traverses the dataset based on "sharding" the result set. The quality of the index for different traversal is highly dependent on the quality of sharding, where sharding incurs expensive traversal and storage costs. For high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

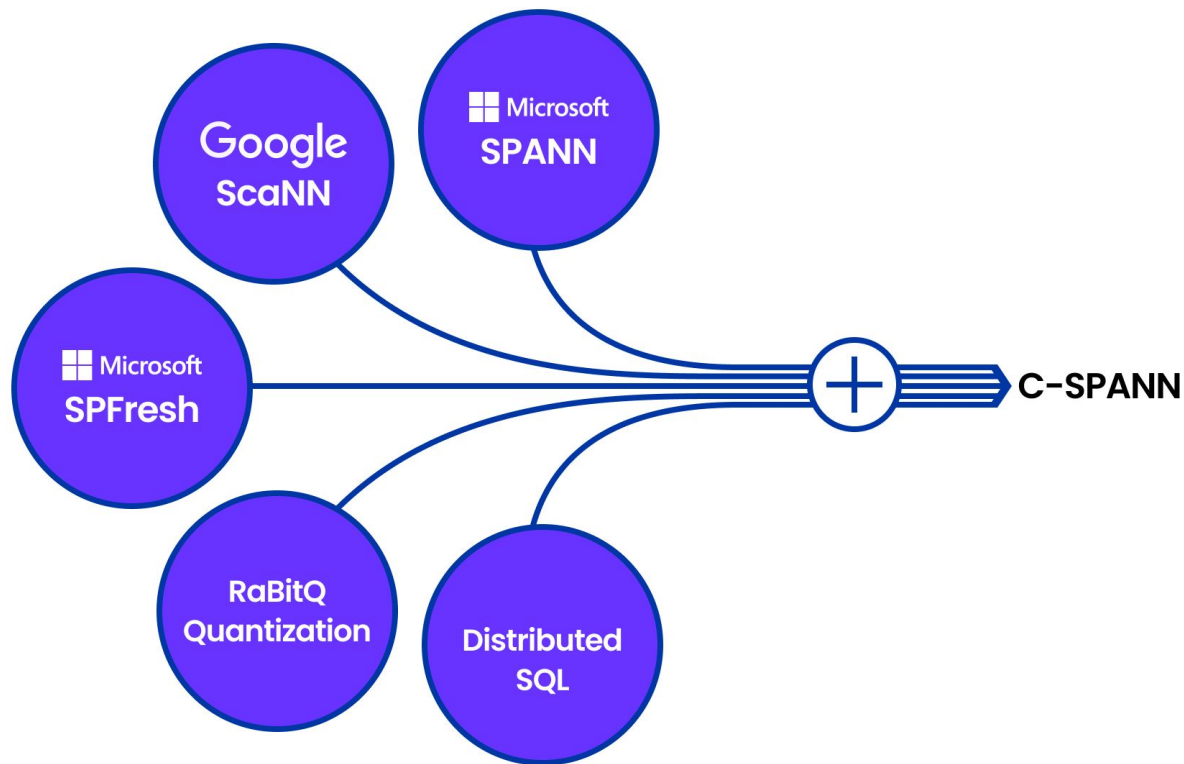
For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

For the high-dimensional data, vector indices require careful construction to produce a sufficient amount of high-quality "sharding" [11, 47].

The solution: introducing C-SPANN



Why Cockroach DB + vector search?



Scalability

- **Horizontal Scale:** Take advantage of Cockroach DB's horizontal scalability to store large amounts of Vector data.
- **Distributed Architecture:** Automatically scales out horizontally as Vector data grows, maintaining performance and reliability without complex manual sharding.





Why CRDB + vector search?

Scalability

- **Horizontal Scale:** Take advantage of Cockroach DB's horizontal scalability to store large amounts of Vector data.
- **Distributed Architecture:** Automatically scales out horizontally as Vector data grows, maintaining performance and reliability without complex manual sharding.

Resilience

- **High Availability:** Industry-standard 99.999% SLA on CockroachDB Cloud Dedicated ensures continuous operation and minimal downtime.
- **Self-healing infrastructure:** Vector data automatically rebalances when a node is lost and heals when you add back a node





Why CRDB + vector search?

Scalability

- **Horizontal Scale:** Take advantage of Cockroach DB's horizontal scalability to store large amounts of Vector data.
- **Distributed Architecture:** Automatically scales out horizontally as Vector data grows, maintaining performance and reliability without complex manual sharding.

Resilience

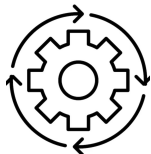
- **High Availability:** Industry-standard 99.999% SLA on CockroachDB Cloud Dedicated ensures continuous operation and minimal downtime.
- **Self-healing infrastructure:** Vector data automatically rebalances when a node is lost and heals when you add back a node

Enterprise Features

- **Integrated Vector and Operational Data:** Combines the strengths of a Vector database and an operational database into a single, resilient and scalable solution, simplifying architecture and reducing costs.
- **Future-Proof Investment:** Supports modernizing legacy applications and meets evolving demands of AI and machine learning workloads.



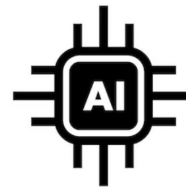
Redefining the limits of operational data using CockroachDB



Modernization



Cloud



AI

CockroachDB Distributed SQL

Delivers state of the art resilience

Scales effortlessly to meet demand

Expands easily everywhere

Simplifies operations, optimize TCO

Can be deployed anywhere and available as self-hosted software or SaaS managed service



Thank you!

Your feedback fuels our roadmap

