

# The Developer Data Platform, Built for AI

Couchbase Capella Architecture



# Contents

<b>INTRODUCTION</b>	<b>4</b>
The developer data platform for critical applications in our AI world	4
<b>TRANSACTIONAL SERVICES</b>	<b>5</b>
Core database design	5
JSON document data model	6
Data access methods	7
Organizing concepts for documents	7
Deployment design concepts	8
Database Services	8
Distributed design	9
As-a-service-aspects	10
<b>MOBILE EDGE SERVICES</b>	<b>16</b>
About Capella App Services	16
About “Offline-First” applications	16
App Services architecture	17
App Endpoint connection points	18
User journey	19
Prepare	19
Connect	20
Operate	20
<b>AI SERVICES</b>	<b>21</b>
Within Capella	21
Architectural Overview	22
Data Layer	23

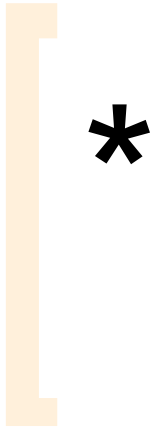


<b>ANALYTICS SERVICE</b>	<b>24</b>
Key Features	25
Ecosystem Integrations	26
Management	27
Operations	28
Security	28
<b>CLOSING</b>	<b>29</b>



# INTRODUCTION

---



The rapid advancement of artificial intelligence (AI) has ushered in the era of new applications and AI agents capable of performing complex tasks across various industries. However, developing these sophisticated applications presents significant data challenges, including managing unstructured data, ensuring low-latency access to large language models (LLMs), and maintaining data security and privacy. Couchbase Capella™ offers a comprehensive solution to these and other foundational data challenges for modern applications. It provides a unified platform that integrates data management for transactions, operational analytics, mobile edge solutions, AI model hosting, vectorization, and AI agent orchestration and governance. Additionally, it is a data platform that can bridge the worlds of generative and predictive AI for agents to make effective and reliable predictions that can drive business outcomes. This paper provides a deep dive into the Capella architecture and how it facilitates the development of modern applications and AI agents, detailing the functionalities of key services and highlighting the anticipated benefits for enterprises.

## The developer data platform for critical applications in our AI world

Couchbase Capella has four key service areas that teams can leverage together to design, build, deploy, scale, and evolve modern applications and AI agents. Those areas include:

- **Transactional Services** – Capella Database-as-a-Service (DBaaS) is the fastest, easiest, and most affordable way to start with Couchbase. At its foundation is a fully managed version of Couchbase Server with a flexible JSON data model, built-in caching, automated data distribution and sharding, with easy horizontal scale-out and advanced security. The architecture is built to ensure high-performance operations, support ACID transactions, provide data model and data access flexibility, and support distributed cluster networks for global high availability.
- **Mobile Edge Services** – Capella App Services is a fully managed backend designed for mobile, IoT, and edge applications to guarantee they are always on, regardless of web connectivity and speed. Developers use App Services to access and sync data between Capella and edge devices and to authenticate and manage mobile and edge application users.
- **AI Services** – Capella integrates cutting-edge AI capabilities directly into the platform, offering features like integrated model hosting and configuration, complete AI memory, automated unstructured data processing and vectorization to support RAG use cases, which improve latency, security, and costs. Additionally, it provides an AI agent catalog that lets developers build and deploy AI agents more efficiently, with critical tools for agent monitoring and governance.



- **Analytic Services** – Couchbase Capella Analytics is a JSON-native NoSQL analytical data store with a massively parallel processing (MPP) engine that enables real-time data analysis by connecting operational and analytical workloads within a single platform. This unified approach eliminates the need for complex ETL processes, allowing organizations to gain timely insights and incorporate new metrics into the operational application.

## TRANSACTIONAL SERVICES

---

### Core database design

Couchbase developed its platform following three guiding principles: memory and network-centric architecture, workload isolation, and an asynchronous approach to everything.

#### MEMORY AND NETWORK-CENTRIC ARCHITECTURE FOR SPEED AND LOW LATENCY

- The most used data and indexes are transparently cached in-memory for fast reads.
- Writes are performed in-memory and replicated or persisted synchronously or asynchronously. Transaction guarantees can be used to ensure consistency, but may introduce lags in performance.
- Internal Database Change Protocol (DCP) streams data mutations from memory to memory at network speeds to support replication, indexing, and mobile synchronization.

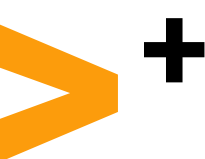
#### MULTI-MODEL DATA ACCESS BLENDS JSON FLEXIBILITY WITH KEY-VALUE SPEED

Couchbase is a pioneer multi-model database that offers multiple data access methods to gain read and update access to its foundational JSON and key-value storage structures. Many other NoSQL systems have only one access method that is bound to their physical storage design structures on disk to minimize access latency.

Couchbase has multiple data access models including key-value, SQL++ query, full-text search, vector search, and event-driven interactions. In the Couchbase design, every access model can simultaneously utilize the cluster's data.

#### WORKLOAD ISOLATION AND ASYNCHRONOUS PROCESSING

All databases perform different tasks in support of an application. These tasks include persisting, indexing, querying, aggregating, and searching data. Each of these workloads has slightly different performance and resource requirements. Couchbase's Multi-Dimensional Scaling (MDS) isolates these workloads from one another at both the process and node levels. MDS allows these workloads to be scaled independently from one another and their resources to be optimized as necessary. It allows the database to be performance-matched to the performance needs of the application, and the database to its available infrastructure.



COUCHBASE IS A PIONEER  
MULTI-MODEL DATABASE  
THAT OFFERS MULTIPLE DATA  
ACCESS METHODS TO GAIN  
READ AND UPDATE ACCESS  
TO ITS FOUNDATIONAL JSON  
AND KEY-VALUE STORAGE  
STRUCTURES.





A SINGLE JSON DOCUMENT'S  
STRUCTURE OFFERS EVEN  
MORE FLEXIBILITY FOR THE  
DEVELOPER BEYOND THE  
DYNAMIC NATURE OF SCOPES  
AND COLLECTIONS.

For cloud deployments, it is advantageous from a cost perspective to redline infrastructure instances before adding them, and to avoid idle and underutilized node instances. Couchbase transparently manages the topology, process management, statistics gathering, high availability, and data movement between these services.

Traditional databases increase latency and block application operations while running synchronous operations (for example, while persisting data to disk or maintaining indexes). Couchbase allows write operations to happen at memory and network speeds while asynchronously processing replication, persistence, and index management. Spikes in write operations don't block read or query operations, while background processes will persist data as fast as possible without slowing down the rest of the system. ACID transactions are available to the developer to ensure durability and consistency while data is in flight. Multiple transaction options allow the developer to decide when and where to increase latency in exchange for durability and consistency of transactions. Somewhat higher latency can be anticipated when multi-document and cross-collection transactions are implemented.

## JSON document data model

The JSON data model supports basic and complex data types, including numbers, strings, nested objects, and arrays. JSON provides rapid serialization and deserialization, is native to JavaScript, and is the most common REST API data format. Consequently, JSON is extremely convenient for web application programming.

A document often represents a single instance of an application object (or nested objects). It can also be considered analogous to a row in a relational table, with the document attributes acting similarly to a column. Couchbase provides greater flexibility than the rigid schemas of relational databases by allowing JSON documents with varied schemas and nested structures. Developers may express many-to-many relationships without requiring a reference or junction table. Subcomponents of documents can be accessed and updated directly, and multiple document schemas can be aggregated into a virtual table with a single query.

### JSON DOCUMENT FLEXIBILITY

In the Couchbase document model, a schema is the result of an application's structuring of its documents and their containment structures such as buckets, scopes, and collections. Schemas can be defined by application developers and managed by applications. This is in contrast to the relational model where the database (and the database administrator) manages the schema. Couchbase created the bucket-scope-collection-document organizational hierarchy (further explained below) to allow maximum flexibility in defining application data metamodels. A single JSON document's structure offers even more flexibility for the developer beyond the dynamic nature of scopes and collections. A JSON document's structure consists of its inner arrangement of attribute-value pairs. How the documents are designed or updated over time is up to the application developer. They can be normalized, denormalized, or a hybrid depending on the needs and evolution of the application. Using JSON, the developer can avoid the lengthy schema design, testing, and deployment cycles of traditional RDBMS-based systems.





COUCHBASE OFFERS A FLEXIBLE MULTI-LEVEL DATA CONTAINMENT AND ORGANIZATION STRUCTURE TO ORGANIZE DOCUMENTS, OPTIMIZE CLUSTER PERFORMANCE, AND FACILITATE HORIZONTAL SCALING.



## Data access methods

Managing JSON data is at the core of Couchbase's document database capabilities, and there are several ways for applications to access the data.

Access Method	Description
Key-value	An application provides a document ID (the key), and Couchbase returns the associated JSON or binary object. The inverse occurs with a write or update request.
Query	SQL-based query syntax, similar to what is used with relational databases, interacts with JSON data and returns matching JSON results. Comprehensive DML, DQL, and DDL syntax supports nested data and nonuniform schema.
Full-text search	Using text analyzers with tokenization and language awareness, a search is done for a variety of field and boolean matching functions. Search returns document IDs, relevance scoring, and optional context data.
Vector search	Enables fast and efficient retrieval of similar items by comparing high-dimensional vector representations of data, for tasks like image recognition, natural language processing, and retrieval-augmented generation (RAG).
Event processing	Custom JavaScript functions are executed within the database based on timers or data changes. Accessing and updating data, writing out to a log, or calling out to an external system are all supported.

## Organizing concepts for documents

Couchbase offers a flexible multi-level data containment and organization structure to organize documents, optimize cluster performance, and facilitate horizontal scaling. This data containment model consists of four levels: buckets, scopes, collections, and documents. This model maps easily to familiar RDBMS constructs of databases, schema, tables, and rows.

- **Buckets** – The topmost container in Couchbase is the bucket. One or many buckets can be defined and assigned to a Capella database.
- **Scopes** – Scopes are an intermediate data organization structure similar to a relational database schema. Scopes are defined by the collections of documents they contain or can access.
- **Collections** – Collections are categorical or logically organized groups of documents. The premise of collections is to behave as traditional table structures.





Most group-oriented access activities are processed at the collection level to minimize full-database operations, simplify replication logic, and streamline indexing options.

- **Documents** – Documents are stored within buckets, but can also be organized within scopes and collections.

## Deployment design concepts

Services and nodes are key elements of the database architecture.

- **Services** – The core of Couchbase is the Data Service that feeds and supports all the other systems and data access methods. Multiple services that offer different types of data access or processing include Query, Indexing, Backup, Search, and Eventing. A service is an isolated set of processes dedicated to particular tasks. For example, indexing, full-text search, and query are each managed as separate services. One or more services can be run on one or more nodes as needed.
- **Nodes** – Capella nodes are virtual machines that host single instances of Couchbase Server within a cloud service provider. Nodes can be added or removed easily through the Capella Control Plane and data is then automatically redistributed evenly across all nodes.
- **Database** – A database consists of one or more nodes running Couchbase Server. Nodes can be added or removed from a cluster. Replication of data occurs between nodes, and cross data center replication (XDCR) occurs between different clusters that are geographically distributed.

## Database Services

Each service has its own resource quotas, and where applicable, related indexing and inter-node communication capabilities. This provides several very flexible methods to scale services when needed. In addition to scaling up to larger machines or scaling out to more nodes, Couchbase also provides the ability to scale specific services independently from one another using Multi-Dimensional Scaling. MDS is the foundation for Couchbase workload isolation and is covered in more detail below.

Couchbase is different from other platforms where a monolithic set of services are installed on every node in a cluster. Instead, Couchbase uses a core data capability that feeds all the other services and a shared-nothing architecture that allows developer control over workload isolation. Small-scale environments can share the same workloads across one or more nodes, while higher scale and performance can be achieved with dedicated nodes to handle specific workloads. This provides the ultimate in scale-out flexibility. The cluster can be scaled in or out and its service topology can be changed on demand with zero interruption or change to the application.

Applications communicate directly with each service through a common SDK that is always aware of the topology of the cluster and how services are configured.

- **Data Service** – The Data Service, or key-value (KV) engine, is the foundation for storing data and must run on at least one node of every database. It is responsible for caching, persisting, and serving data to applications and other services. The







cache provides consistent low latency for individual document read and write operations and streams documents to other services via Database Change Protocol. Due to their simplicity, KV operations execute with extremely low (often sub-millisecond) latency. The KV store is accessed using simple CRUD (create, read, update, delete) APIs, and provides the simplest interface when accessing documents using their IDs.

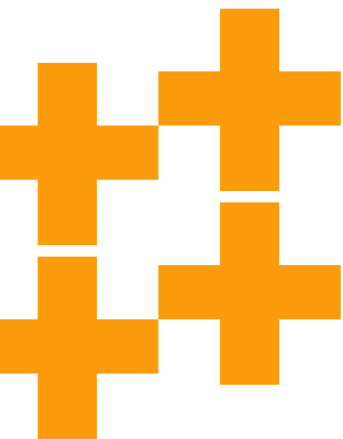
- **Query Service** – An engine for processing SQL++ queries. SQL++ combines the flexibility of JSON with the expressive power of SQL. It provides a rich set of features and familiar data definition language (DDL), data manipulation language (DML), and query language statements, but can operate in the face of NoSQL database features such as key-value storage, multi-valued attributes, and nested objects. Also, users can define ACID transactions within SQL++ for one or more documents across collections and nodes. Transactions in SQL++ have adopted a nearly identical syntax to SQL for relational databases. The Query Service uses a cost-based query optimizer, patented in 2021, to take advantage of indexes that are available.
- **Index Service** – Indexing is an important part of making queries run efficiently and self-update as data mutates. This service supports multi-index types and includes an Index Advisor that recommends specific indexes to build based upon query statements and data structure.
- **Search Service** – An engine for performing full-text and vector searches on stored JSON data. Users can create and query inverted indexes for searching of free-form text within a document or vector indexes for matching by semantically similarity across a wide number of dimensions. Customers using the Search Service often can eliminate the need for a third-party search tool.
- **Eventing Service** – Eventing supports custom server-side functions (written in JavaScript) that are automatically triggered using an event-condition-action model. These functions receive data from the DCP stream and execute code when triggered by data mutations. This service offers a feature like the change data capture (CDC) found in event handlers, and also offers a feature similar to the multi-channel data streaming found in solutions such as Apache Kafka.

## Distributed design

Capella's distributed nature makes high availability, scaling, and disaster recovery easier.

- **Data distribution** – Couchbase automatically partitions and replicates data into vBuckets (synonymous to shards) to automatically distribute data across nodes. This helps enable data replication, failover, and dynamic database reconfiguration. Because vBuckets do not have a fixed physical location on nodes, they are mapped to nodes in a cluster map. Through the Couchbase SDK, the application automatically accesses data without needing to know the exact location of the data.
- **Data transport via DCP** – As data mutates, in-memory replication is used to maintain data updates within Capella and to external services such as Apache Spark or Kafka that are fed from the DCP stream.



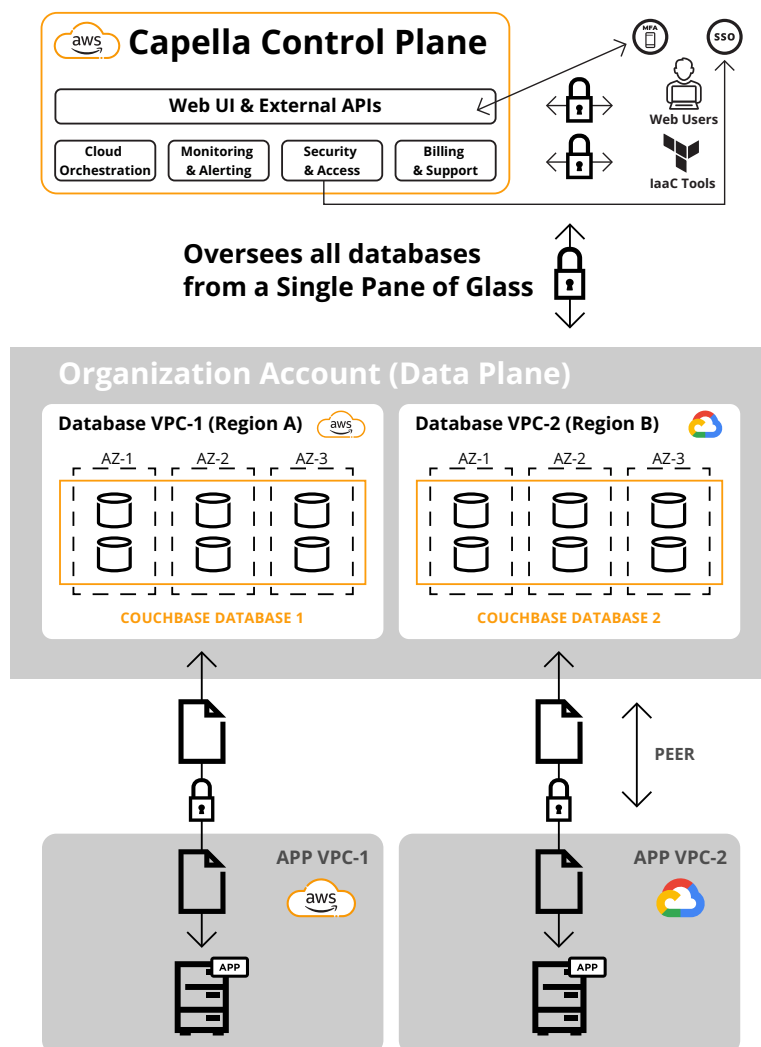


- **Multi-dimensional scaling** – You can improve performance and throughput for systems by independently scaling services to match workloads. Scale-out and scale-up are the two scalability models typical for databases, and Couchbase takes advantage of both. You can combine and mix these models in a single database to maximize throughput and minimize latencies.
- **Failover** – If a node fails in Capella, replicated data on other nodes are promoted to active. A new node is then automatically provisioned and data is rebalanced across all the nodes.

## AS-A-SERVICE-ASPECTS

### Architecture

Capella's core architectural aspect is the split of the web UI Control Plane designed for management and the Data Plane for data storage.



## CONTROL PLANE

The Control Plane is a web UI that manages the cloud orchestration Infrastructure-as-a-Service (IaaS), monitoring, alerting, security, access, billing, and support capabilities. It's the access point for your organization's users and also allows access to infrastructure-as-code (IaC) tools such as Terraform. Backend services can be accessed by REST-based tools via a management API.

## DATA PLANE

The Data Plane is where you manage your Capella databases. A database resides in a single region (distributed across multiple availability zones) within a single cloud service provider (CSP), but the Control Plane can control multiple databases across various cloud service providers. Furthermore, data can be replicated between databases, with that replication configured within the Control Plane. From a security perspective, the Data Plane has no internet access unless IPs are specifically allowed or a connection is established through VPC peering or an AWS PrivateLink connection. Disk data, backups, and all traffic is encrypted.

## Management

### USERS, PROJECTS, AND ROLE-BASED ACCESS CONTROL (RBAC)

An organization is the top-level organizational element for managing users, projects, and databases. By default, organizational members do not have data access, but instead have access to the Control Plane. A project is used to organize groups of databases. People must be added to projects and assigned access to databases within a project. Organizations can also add SSO groups (teams) with a project role to access databases within a project.

### SINGLE SIGN-ON (SSO)

SSO is a convenient way to maintain users by making use of their existing corporate credentials. Authentication is delegated to the SSO provider (Azure AD and Okta are supported). Provisioning new users via SSO eliminates the overhead of having to send invites.

## Deployment

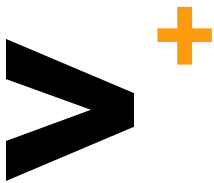
### DATABASE DEPLOYMENT

When deploying a database, you must choose a CSP, a region, and a CIDR block. (A default is provided, but can be changed before deploying.) Couchbase is constantly adding regions, and up-to-date regions can be found on the Couchbase Documentation pages for [AWS](#), [GCP](#), and [Azure](#). You can select the version of the Couchbase Server you would like to use and assign services to nodes. These services can be changed after deployment. You can also choose your support plan and availability mode (single or multi-availability zones), and can choose to purchase credits on a prepaid or pay-as-you-go basis.

### STORAGE ENGINE

Capella supports two different backend storage mechanisms, which are set per bucket. A single Capella database can have a mix of Couchstore and Magma buckets.





- **Couchstore** – Couchstore is the default bucket storage engine that has been in use for more than 10 years. It's optimized for high performance with large datasets while using fewer system resources. (The minimum bucket size for the Couchstore backend is 100 MiB.) If you have a small dataset that can fit in-memory, then you should consider using Couchstore.
- **Magma** – Capella's latest storage engine is designed for high performance with very large datasets that don't fit in-memory. It's ideal for use cases that rely primarily on disk access. The performance of disk access will be as good as the underlying disk subsystems. Magma can work with very low amounts of memory for large datasets (e.g., for a node holding 5 TiB of data, Magma can be used with only 64 GiB RAM). It's especially suited for datasets that won't fit into available memory.

You can learn more about Capella's storage engines in our Couchbase [documentation](#).

## Development

Couchbase provides several tools for developers:

### PLAYGROUND

This tool is integrated into Capella and comes with an SDK tutorial and a SQL++ tutorial. The SDK playground offers examples of multiple SDK languages. SQL++ gives examples of the Couchbase query language. Both tutorials are designed to guide a new developer through several chapters from basics to more advanced concepts.

### QUERY WORKBENCH

The Query Workbench allows users to access data via SQL++ and see data in JSON and tabular formats. The tool provides a built-in index advice feature that tells users what indexes are needed to optimize queries. Inverted search indexes can be created to support search, and JavaScript-based user-defined functions can be used to manipulate data.

### COUCHBASE SHELL

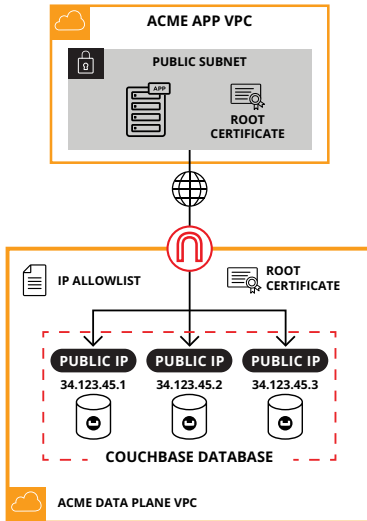
Couchbase Shell (cbsh) is a modern, productive shell that provides CLI access to Capella. It can be obtained via [GitHub](#).



## Connecting

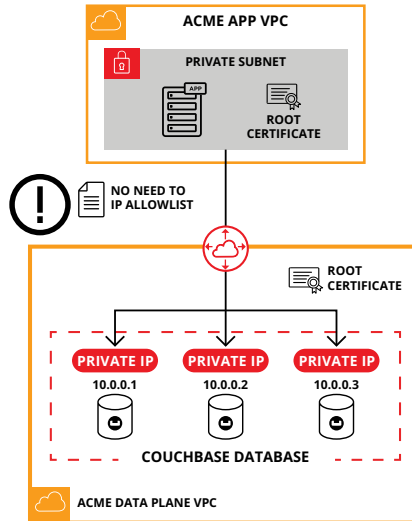
### OPTION 1: PUBLIC CONNECTION

- Connect to Node **Public** IP
- Needs IP Allowlisting
- Traffic traverses public Internet



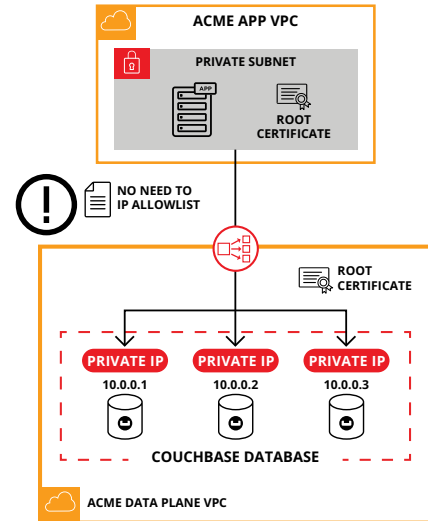
### OPTION 2: VPC PEERING

- Connect to Node **Private** IP
- Traffic contained within CSP backbone
- Joined network between the 2 VPCs



### OPTION 3: PRIVATE LINK

- Connect to Node **Private** IP
- Traffic contained within CSP backbone
- Unidirectional and very secured



## COUCHBASE SDKS

Capella works with the latest versions of all supported Couchbase [SDKs](#). Developers can choose from over 10 SDKs of their favorite programming languages.

## CONNECTORS

To exchange data with other platforms, we offer various big data [connectors](#) for products like Kafka, Spark, and Elasticsearch.

## REST API

Couchbase provides a series of [RESTful APIs](#) that enable you to integrate with Capella to perform operations such as:

- Onboarding and offboarding users
- Managing the lifecycle of a cluster
- Getting monitoring information for a cluster

## APPLICATION CONNECTION

For connecting applications to Capella, you have several options:

- **Public connection** – This is the simplest option and requires the use of IP addresses for encrypted data to traverse the public internet. Public connections should not be used for production environments.
- **VPC peering** – The CSP backbone contains both the connection and traffic. This option reduces networking costs compared to a public connection.
- **AWS private endpoint** – Allows Capella to be offered as a private service and functions as if it were hosted directly within a team's Amazon VPC. This allows access to a specific service or application, and only private endpoints can initiate a connection.





CAPELLA PROVIDES  
A WIDE VARIETY OF  
METRICS TO MONITOR  
AND HELP OPTIMIZE  
PERFORMANCE.

## Operations

### SCALING

Capella makes it easy to evolve your configuration. You can add or remove nodes at any time and change the amount of RAM, vCPUs, disk space, and type of disk volume (general purpose or high performance). Changes are made with a few clicks, and Capella automatically rebalances data to the new configuration.

### MULTI-DIMENSIONAL SCALING

MDS allows you to further optimize your configuration as application needs evolve over time. MDS allows workloads to be scaled independently and hardware usage to be optimized to help drive down total cost of ownership.

### GEO-REPLICATION

Capella's cross data center replication (XDCR) technology replicates data between databases in different regions. XDCR provides an easy way to replicate active data in-memory to multiple geographically diverse data centers either for disaster recovery or for high availability. It can be set up on a per-bucket or per-collection basis and can be unidirectional or bidirectional. It also provides built-in conflict resolution if the same document was mutated on a separate database before it was replicated.

### BACKUP AND RESTORE

A robust scheduled backup and retention policy is recommended as part of an overall disaster recovery plan for production data. Backup is done on a per-bucket level and can be scheduled on monthly, weekly, or daily cycles, or on demand. You can set up both full and incremental backups. Restoring data can happen at the bucket or collection level with filtering options. Additionally, data can be restored into Capella from a self-managed Couchbase Server environment.

### MONITORING AND ALERTING

Capella provides a wide variety of metrics to monitor and help optimize performance. Metrics can be viewed within a Capella Control Plane configurable dashboard or incorporated into your Prometheus tool. Capella has a set of conditions that generate alerts and provide actionable suggestions when a threshold is hit (e.g., a suggestion to increase RAM or disk size).

### MAINTENANCE AND UPGRADES

If you want to upgrade your databases, those processes can be scheduled within the Control Plane. Security updates are forced immediately. Notifications about upgrades are sent via email and within the UI.

## Security

Couchbase supports the most critical and sensitive workloads for industry-leading businesses every day. Capella's security architecture is based on industry best practices for security and four key pillars:

- **Verify** – Role-based access controls ensure only authorized users or applications have access to data.



- **Enforce** – Enforcement of least privilege access is applied to all credentials and secrets, ensuring strict access controls to sensitive data and actions.
- **Monitor** – To prevent potential breaches, Capella implements a managed cloud intrusion detection system that involves 24x7 monitoring.
- **Modernize** – Capella is built using modern DBaaS principles and secure development practices.

Couchbase has achieved many of the most important security compliances, including SOC 2 Type II, HIPAA, GDPR, PCI DSS, ISO 27001, and more. You can learn more and get detailed security whitepapers and information about compliance at our [Trust Center](#).

## INFRASTRUCTURE SECURITY

The foundation of security in a cloud database is a hardened environment that removes nonessential software, roles, and ports while leveraging an IaaS provider's alerting and auditing services. Trusted and immutable operating system (OS) images are used to protect the OS, with verification upon deployment and ongoing scanning for vulnerabilities afterward. Additionally, end-to-end configuration is automated via templates to ensure consistency. Monitoring is also in place to identify potential misconfigurations.

## NETWORKING SECURITY

By default, the Data Plane only allows clusters to connect to trusted IP addresses that have been defined within the Control Plane. Any attempted connection from an IP address not in a cluster's list of allowed IP addresses will be denied. With VPC peering, traffic never crosses the public internet, which reduces threat vectors and DDoS attacks. If you're using AWS for your data plane, you can further enhance security by using PrivateLink to contain traffic within the CSP backbone with unidirectional access. Alternatively, you can set up Capella as a private service that functions as if it were hosted directly within a team's Amazon VPC. This allows access to a specific service or application, and only private endpoints can initiate a connection.

## ACCESS SECURITY

To bolster security access, Capella is designed so that the Control Plane and Data Plane live in separate VPCs. Access to data is separate from access to the Control Plane, and specific credentials must be established for application access. All admins, users, and applications must authenticate in order to gain access and then be authorized with specific access rights. Multi-factor authentication is possible and recommended.

## DATA SECURITY

Data is encrypted in transit and at rest. In transit, data is encrypted via TLS, which cannot be turned off. If you want to extend data storage encryption within the database, this can be done at the field level within JSON documents. Also, backup data is written to encrypted disks using the cloud provider's native encryption process. Capella creates, manages, and controls cryptographic keys using a CSP's key management system (e.g., KMS for AWS).





## VULNERABILITY MANAGEMENT

Capella protects against threats like brute force attacks, rate-limiting attacks, cross-site request forgery, and more. Capella maintains centralized logs securely and alerts Couchbase site reliability engineers (SREs) of operational concerns should they arise. To reduce potential vulnerabilities, patching is automated and includes monitoring alerts and management reviews. Couchbase has established a formal Incident Response Policy to inform you in the event of a security-related event.

## MOBILE EDGE SERVICES

---

### About Capella App Services

Capella App Services provides a hosted gateway and WebSockets-based protocol for bidirectional data synchronization between Capella and embedded apps on smartphones, tablets, IoT devices, and custom embedded devices. It works in tandem with Couchbase Lite, the embeddable version of the Couchbase database. Wherever Couchbase Lite runs, App Services can securely sync the data it captures to Capella buckets and other embedded devices.

App Services also manages secure data access with role-based access control, providing authentication for mobile users. These key capabilities in Capella are offered as a ready-to-use service for mobile and IoT developers, making it faster and easier than ever to build highly performant and reliable applications.

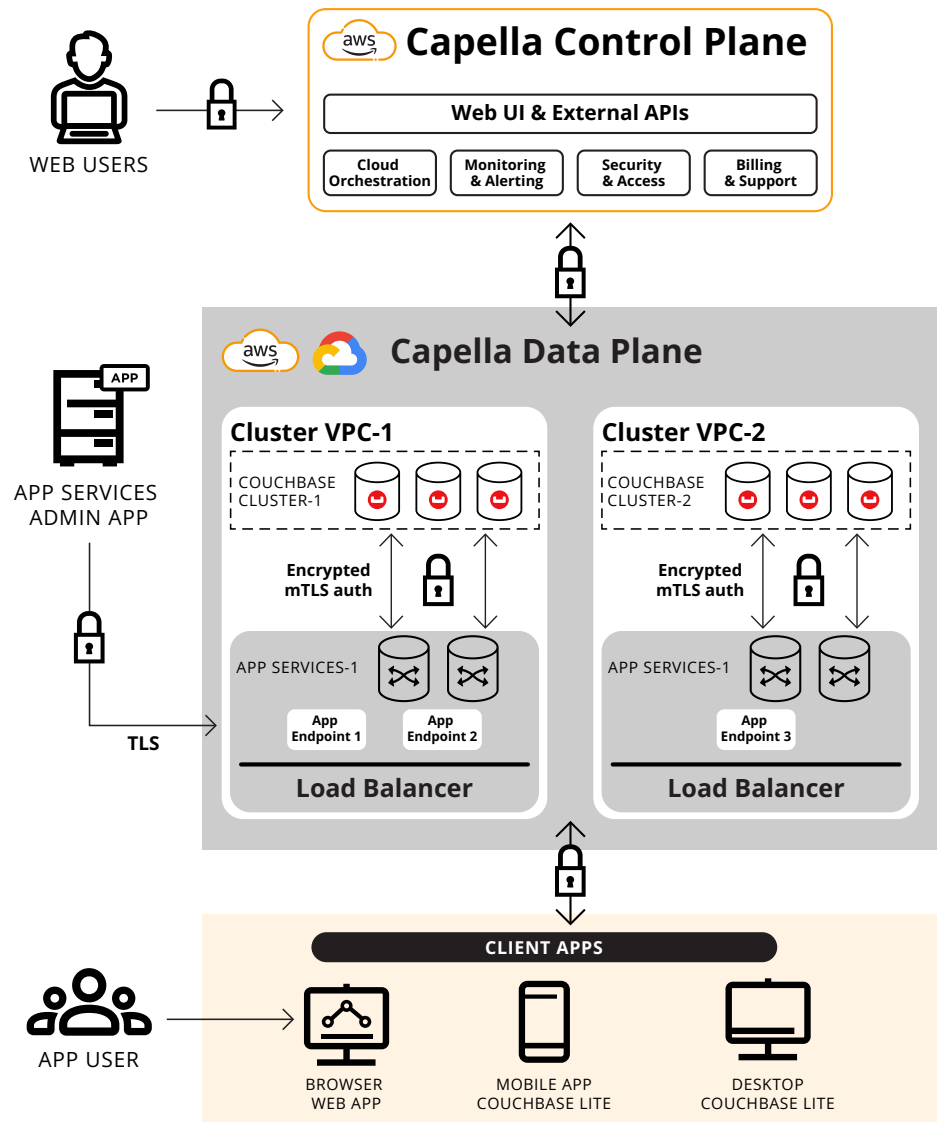
### About “Offline-First” applications

For applications that need to operate in areas with slow or no internet, embedded Couchbase Lite provides on-device data storage and processing. This allows apps to work all of the time, whether online or offline (hence the term offline-first). App Services sync is smart enough to know when connectivity is interrupted. When it's restored, App Services can automatically start syncing from where it left off even after long periods of time.

More importantly, when multiple Couchbase Lite clients are close to one another, but have no internet, they can still sync data via peer-to-peer. This feature enables offline-first collaboration without the need for any central control point or internet connection.



## App Services architecture



When you create an App Service and associate it with a Couchbase Server cluster, you are effectively extending or enabling it for data sync. A Couchbase cluster can only be linked to one App Service.

When an App Service is created, a cluster of Sync Gateway nodes is deployed behind the scenes in the same VPC network as the corresponding server cluster. Communication between the App Services cluster and the central Capella cluster is secured using TLS and x.509 certificate-based authentication. The App Services cluster is fronted by a load balancer that balances incoming client requests across the App Services nodes.

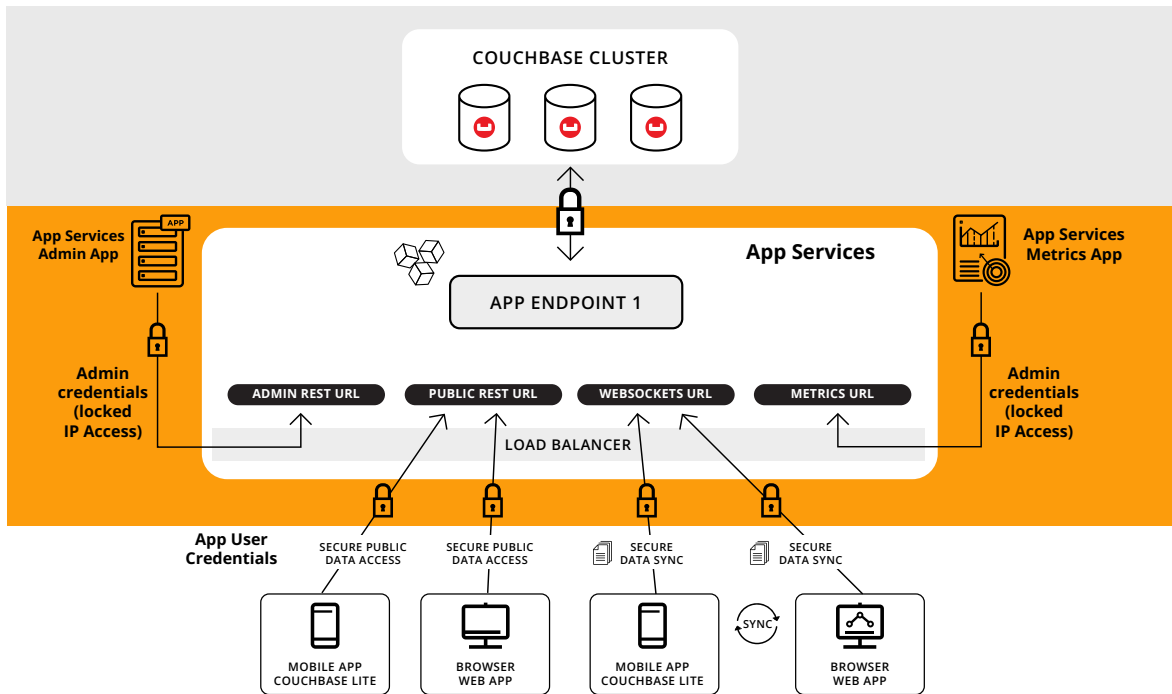
An App Service can handle multiple client applications, each represented by an App Endpoint. Conceptually, an App Endpoint represents the instance of your application on the App Service. Each App Endpoint is backed by a Capella bucket. If you have multiple applications, each will have its own App Endpoint.



Mobile, desktop, and web client apps can access and sync data by connecting to the corresponding App Endpoint.

## App Endpoint connection points

There are multiple options for connecting clients to an App Endpoint. Your choice depends on the type of application and use case.



### SECURE WEBSOCKETS PUBLIC URL

Offline-first sync is the ability for apps to run in offline mode in the face of temporary or extended network disruptions and to sync data with the backend servers when connectivity is restored. Mobile, desktop, and embedded apps powered by Couchbase Lite can locally store and access data in disconnected mode and sync data with App Services when there is connectivity. With the internet being inherently unreliable, the use cases for offline-first data sync are vast and varied.

### SECURE PUBLIC REST API

Applications can also access data securely over a public RESTful API endpoint. This is useful when there is reliable network connectivity and no need for offline storage, or when the apps are running on hardware that doesn't have local storage for running a local embedded database like Couchbase Lite.

### SECURE ADMIN REST API

Administration applications can be granted authenticated access to the admin REST API in order to programmatically create and manage users, roles, and sessions. Admin apps are typically hosted in the cloud backend. An example of an admin app is a login service that handles custom authentication and is responsible for registering users via the secure admin REST API following successful user authentication.

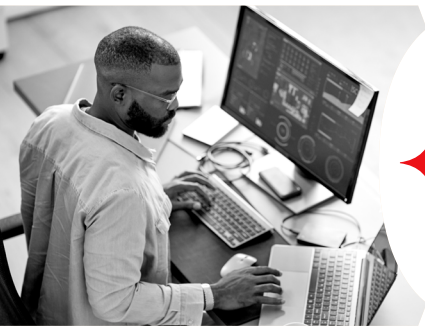
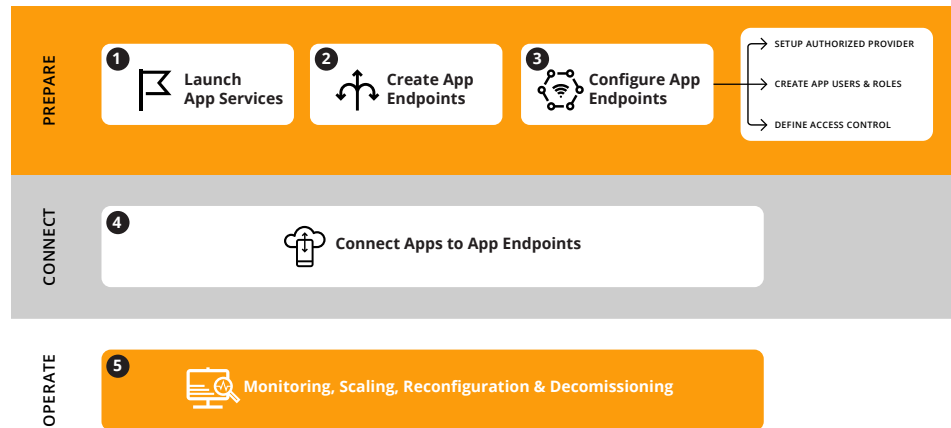


## SECURE METRICS REST API

Monitoring frameworks like Prometheus can access stats exposed via the metrics REST endpoint. In addition, App Services also supports a dashboard of common operational stats.

## User journey

Prerequisite: App Services requires a Couchbase Capella server cluster. Follow [these steps](#) to create a Capella cluster and set up a bucket.



## Prepare

### LAUNCH APP SERVICES

When you create an App Service and associate it with a server cluster, you are effectively enabling it for data sync. When creating an App Service, you give it a name, designate an associated Capella cluster, then choose the deployment configuration that includes the correct number of nodes and type of computer (RAM/core).

### CREATE APP ENDPOINTS

App Endpoints represent the instance of your application on an App Service. You can create multiple App Endpoints on an App Service, each backed by a unique bucket in the corresponding Couchbase Server cluster. By default, all documents in the corresponding bucket are imported by the App Endpoint.

### CONFIGURE APP ENDPOINTS

When the App Endpoint is created, it is set up in offline mode. This allows users to complete the security configuration of the App Endpoint before exposing it to applications.

### AUTHENTICATION PROVIDER

Authentication providers define how users are authenticated with the App Services. A default auth provider of basic auth is selected for you during App Endpoint creation. So you can skip this config if the default option works for your application.



Capella supports the following modes of authentication:

- **Basic Auth** – This is where the app users are authenticated using username and password credentials that are Base64 encoded and passed in as part of the authorization header of an HTTP request.
- **Open ID Connect (OIDC)** – App users are authenticated against a third-party identity provider that is registered with App Endpoint. This is implemented using [OIDC Implicit flow](#).
- **Anonymous** – In this mode, we allow unauthenticated read-only access to data. This mode can be useful when your app is only dealing with public static data.

## USER MANAGEMENT

With the exception of “Anonymous” mode, all client-side access must be authenticated with suitable user credentials. The choice of how users (and roles) are created depends on the authentication provider that is configured.

- **Basic Auth** – Users are created via the Capella web UI or via Admin REST Endpoint.
- **Open ID Connect (OIDC)** – By enabling the “auto-register” option when configuring OIDC provider, users will be automatically created on App Service after successful authentication.

## ACCESS CONTROL

Access control is implemented using the [channel-based access control model](#) of Couchbase Mobile. Access control specifies who has access to what data. This is specified via a JavaScript access control function. Read access control is at the granularity of a document, while write access control is at the granularity of a field.

## Connect

After completing the security setup for the App Endpoint, unpause the App Endpoint to bring it online. Once online, apps can be connected using any connection points discussed earlier.

## Operate

Once your App Service is operational, you can administer the App Service and App Endpoints and change the configuration to meet the evolving needs of the apps.

## MONITORING

Metrics dashboards provide insights into resource utilization of the App Service as well as the operational state of the App Endpoints. These include stats such as the number of documents read/written, error counts, number of active replications, etc.

## ACTIVITY LOG

All key system events of type info, warning, and error are recorded in the activity center. Users are also alerted to key events that may need attention, such as significantly high memory utilization over an extended period of time.

## ON-DEMAND SCALING

To keep up with the evolving needs of the app, users can scale App Services horizontally and/or vertically by changing the number of nodes and/or compute type.



Software development is evolving beyond traditional databases, with AI models now playing a central role in how applications interact with and process data. Historically, applications primarily retrieved, stored, and processed data only through databases, but AI models have introduced new opportunities by dynamically analyzing, transforming, and generating data in real time. This shift significantly increases interaction complexity, requiring a sophisticated data platform to manage data flow efficiently between agents and AI models. To address latency and security risks, integrating AI models closely with databases is becoming essential, ensuring faster responses and minimizing vulnerabilities associated with distributed architectures.

AI agents perceive their surroundings, make decisions, and take actions to achieve specific goals by breaking down complex objectives into smaller, manageable tasks. These agents represent a transformative technology designed to function across digital and physical environments by leveraging orchestration logic, AI techniques, and APIs. Their efficiency can improve over time by learning from past experiences and user preferences. Key components include AI models and frameworks, data access and memory storage, as well as seamless tool integrations. By automating complex processes, enhancing decision-making, and personalizing user experiences, AI agents have the potential to revolutionize operations across industries.

### Within Capella

Capella AI Services are a comprehensive suite tailored to support AI-enabled agent development and deployment workflows. The services are designed to provide developers with control over data throughout the development lifecycle and offer seamless integration with AI models. Key benefits include:

- **Complete memory** – Being a unified memory store for multi-agentic and long-running agentic applications, enabling the agents and applications to maintain context and personalize responses across multiple conversations, essentially acting more like a human with a memory by remembering relevant details from previous interactions.
- **Cataloging** – Enabling discovery, governance, and auditing of agentic tools with a catalog of tools, (e.g., weather API, map API, etc.), code, and functions (get data, summarization) for efficient storage and fast semantic recall.
- **Data access** – Eliminating complexity by integrating multiple data access patterns and support for transactional, analytical, vector, caching, and predictive AI use cases in a single platform by unifying diverse data and formats, ranging from structured (tables w/rows and columns), semi-structured (JSON, HTML code), or unstructured (text, audio, video, etc.).
- **Performance** – Delivering all of this with high performance, scalability, and low-latency responses is crucial for AI-driven workflows requiring real-time processing. Coupled with a “deploy anywhere” paradigm that gives developers access to familiar interfaces, APIs, and SDKs across clouds, on-premises to the edge, and right on devices.

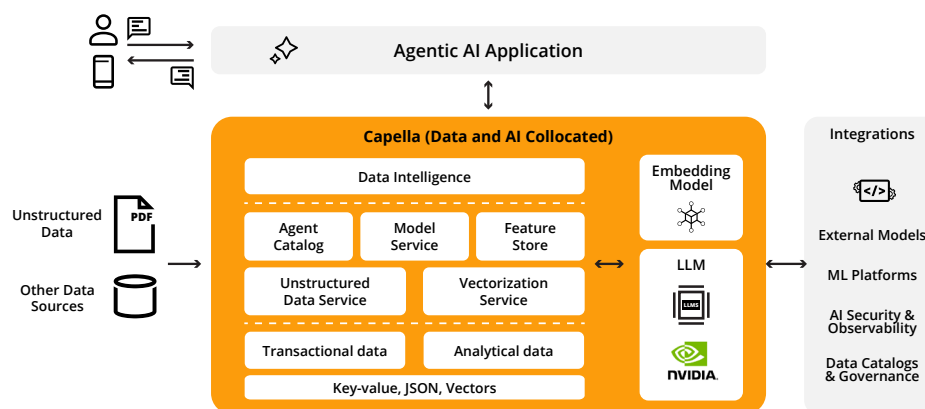


- **Integration** – Easy integration into our cloud service provider partners and the AI ecosystem utilizing key frameworks like LlamaIndex, LangGraph, CrewAI, and more.
- **Security** – Capella's architecture ensures that data remains within the customer's VPC, mitigating risks associated with data exposure during model interactions. This design upholds data privacy and complies with regulatory standards, which is crucial for enterprises handling sensitive information.

## Architectural Overview

The Capella data platform helps customers build AI applications with several important capabilities to support their data needs within the full development and evolution lifecycle. The following breaks down those key areas.

### DATA INTELLIGENCE



- **Capella iQ** – Enhances developer productivity by enabling natural language interactions for tasks such as generating SQL++ queries, creating sample datasets, suggesting optimal indexes, and visualizing data.
- **SDKs and data API** – Provide access to the data platform, including more sophisticated techniques using SDKs across a dozen languages or with simple calls via the data API for simpler low code projects.
- **AI functions** – Enhances developer productivity by embedding AI-driven data analysis and output directly into application workflows using familiar SQL++ syntax, eliminating the need for external tools or custom code. Examples of capabilities include summarization, classification, sentiment analysis, and data masking.

### CORE AI SERVICES

In order to best interact with the AI ecosystem, the AI Services help teams safely interact with models, manage agent development to speed build and LLM interaction traceability, and centralize the storing, managing, and serving of preprocessed ML features to ensure consistency, reusability, and scalability across training and inference pipelines.

- **Model Service** – Host and configure AI models utilizing NVIDIA NIM in order to colocate data and AI in a secure manner where data does not leave a customer's environment. This collocation also reduces latency between the data and the model. Value-added services like semantic caching are designed to reduce risk and





cost. Customers can also utilize the model service work with external models, like those offered by leading cloud service providers.

- **Agent Catalog** – Accelerates agentic application development by offering a centralized repository for tools, metadata, prompts, and audit information for LLM flow, traceability, and governance. It also automates the discovery of relevant agent tools to answer user questions and manages guardrails to ensure that agent exchanges are consistent over time.
- **Machine Learning (ML) Feature Stores** – Serve as a high-performance, distributed repository for storing, managing, and serving ML features in real time. Leveraging Capella's in-memory caching, flexible JSON data model, and multi-dimensional scaling, it enables fast feature retrieval for both batch training and low-latency inference.
- **Unstructured Data Service** – Automates the ingestion, cleaning, chunking, and transformation of unstructured documents into JSON format, preparing them for vectorization while extracting structured information from complex documents to make them queryable within Capella. This process enables AI agents to seamlessly process diverse data types, making the data more manageable and accessible for analysis and decision-making.
- **Vectorization Service** – Automates the creation, storage, and indexing of vector embeddings in real time to improve conversation quality and maintain context for evolving LLMs via natural language prompts.

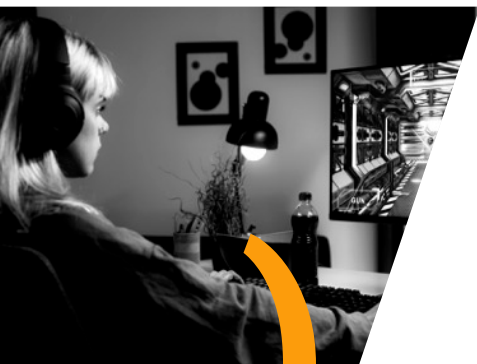
## Data Layer

The data layer is where the data is stored, indexed, and accessed for a variety of purposes and use cases within the other Capella services. Capella supports operational workloads, transactions, JSON document data, vector embeddings, and columnar analytic data with a built-in cache for driving low latency.

### DATA INGESTION AND AI INTEGRATIONS

AI Services make it easy to work within your data and AI ecosystem to develop AI-powered solutions.

- **Unstructured data ingestion** – Automates data preprocessing in order to easily incorporate unstructured data, which has historically been challenging to accomplish. Documents are converted to JSON and can be automatically vectorized and indexed.
- **External models and ML platforms** – Streamlines the integration and configuration of generative AI and predictive AI models to make it easier for developers to work with a variety of AI systems.
- **Security and observability** – Allows for teams to monitor performance, detect anomalies, enforce compliance, and mitigate risks by providing real-time insights into data access, model behavior, and potential vulnerabilities.
- **Data catalogs and governance** – Helps customers ensure data quality, compliance, and accessibility by systematically organizing, tracking, and enforcing policies on data usage, lineage, and security.



For decades, organizations have separated transactional and analytic workloads to optimize performance, but modern JSON-native applications with rapidly changing data expose the limitations of legacy relational analytic tools, which struggle with JSON data and lack the ability to write data back to source systems. Couchbase Capella addresses these challenges as the only JSON-native DBaaS optimized for both real-time transactional and analytical workloads, delivering exceptional performance, scalability, and cost-effectiveness. It simplifies JSON data ingestion and transformation through zero-ETL capabilities and leverages an optimized columnar storage engine to enhance analytic query efficiency. Capella empowers developers with real-time access to analytics, minimizing reliance on centralized data teams, and with Capella iQ, its AI-powered coding assistant, users can perform interactive conversational analytics. Built-in write-back capabilities further enable real-time metrics to be fed back into applications, improving user experiences. By supporting both transactions and real-time analysis in a unified platform, Capella reduces costs and streamlines integration with existing systems.

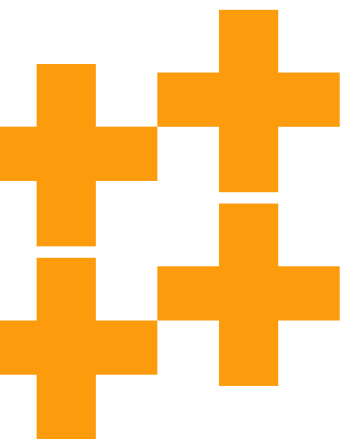
### ARCHITECTURE OVERVIEW

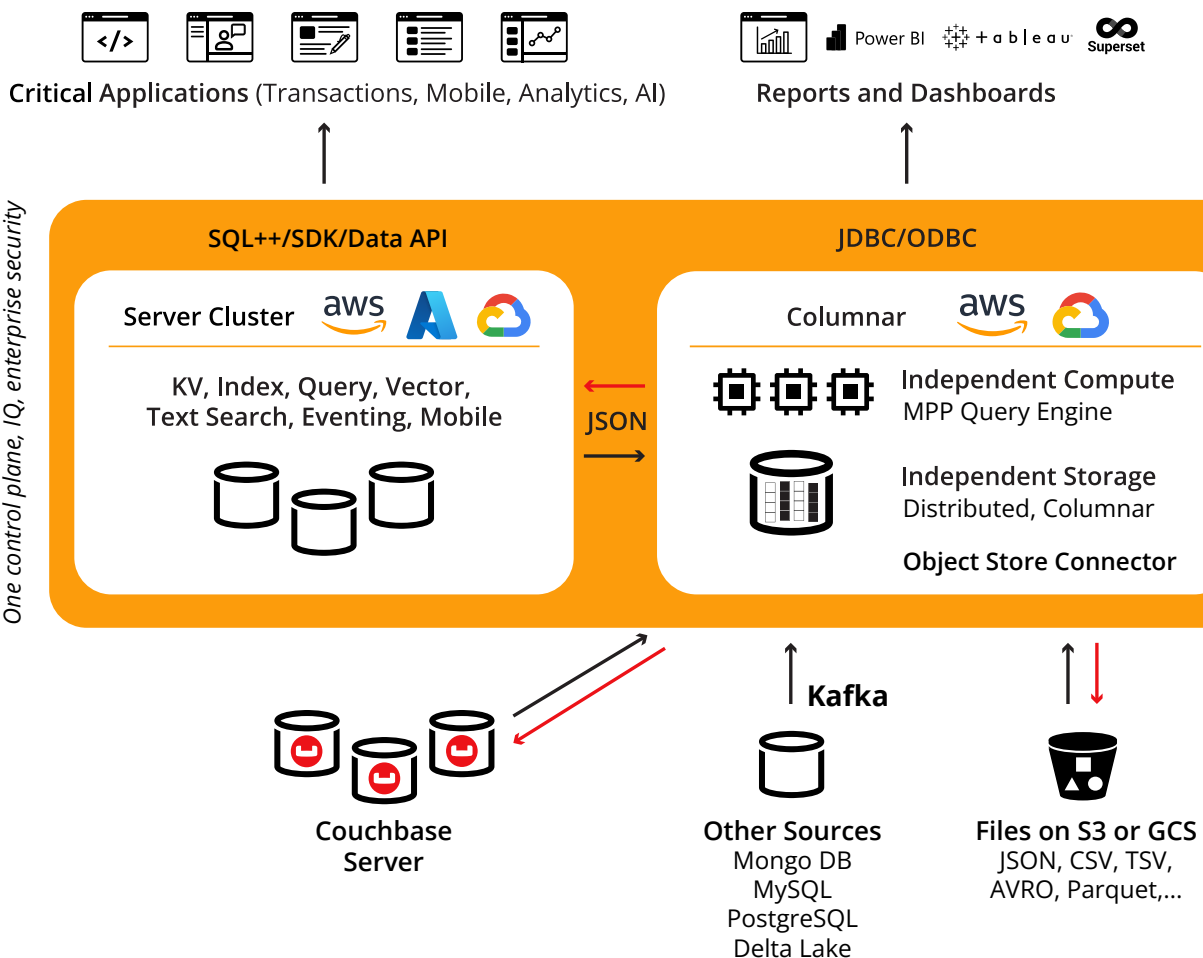
Capella's architecture features a clear separation between the Control Plane and Data Plane, creating a scalable, secure, and efficient environment for managing and utilizing data. The Control Plane serves as the centralized management layer and primary access point for users via the web interface. It orchestrates cloud operations such as provisioning, scaling, and maintaining clusters, while also handling monitoring, alerting, security, billing, and support. Operating across multiple clusters, it enables consistent, efficient management from a single interface.

The Data Plane ensures strong isolation and performance, with each organization assigned a dedicated account and each cluster deployed in its own Virtual Private Cloud (VPC). Clusters span geographic regions and multiple availability zones (AZs) for high availability and fault tolerance. With compute and storage separated, workloads can scale flexibly by independently scaling compute or storage. This happens with no data rebalancing.

Capella also supports easy integration with external systems like Kafka, Amazon S3, and Google Cloud Storage (GCS) for efficient data flow.

Applications connect directly to clusters via public internet, VPC peering, or AWS PrivateLink, with all data in transit fully encrypted for end-to-end security. This architecture delivers a reliable, secure, and high-performance platform that simplifies data management and enables real-time insights.





## Key Features

### COLUMNAR STORAGE FOR JSON

Couchbase Capella Analytics efficiently stores JSON data in a columnar format by automatically inferring schemas as data is moved from memory to disk, eliminating the need for predefined schemas and allowing for seamless handling of dynamic JSON structures. By separating schema from data, Capella significantly reduces overall data size, and stores the data internally in a compact binary format optimized for performance. This process is tightly integrated with the Log-Structured Merge (LSM) component lifecycle, ensuring efficient data organization and management over time. JSON documents are stored column-wise, enabling fast, high-performance analytical queries while maintaining full JSON support without compromise. Columns are compressed and encoded to further minimize storage footprints and reduce I/O overhead, and there is no requirement for all of a collection's data to fit into memory – regardless of the data's origin – ensuring scalability and efficiency for large and diverse datasets.

### COMPUTE AND STORAGE SEPARATION

Couchbase Capella Analytics is designed with a clear separation of compute and storage, enabling flexible scaling and efficient resource use. Data is managed by the storage layer and securely stored in object storage like Amazon S3. When in use,

collections and indexes are cached in high-speed NVRAM, with a lazy caching policy that loads data on demand. Collections are partitioned and dynamically assigned to compute units, each responsible for its own partitions. This shared-nothing compute on shared-disk architecture allows Capella instances to scale based on workload needs, offering isolation, fault tolerance, and efficient parallel processing.

### **MASSIVELY PARALLEL PROCESSING**

Analytic Services achieve high-performance query execution through massively parallel processing, with all compute units working in parallel on their data partitions. It supports fast execution of projections, selections, joins, sorting, aggregations, and windowing functions, while applying relational best practices to JSON data. Local algorithms handle large-scale queries efficiently, spilling to disk when needed. Capella's columnar storage reads only necessary data, minimizing I/O and boosting speed. It also supports partitioned parallelism for querying external data, ensuring balanced workloads and rapid integration from external sources.

### **EXTERNAL COLLECTIONS**

Couchbase Capella Analytics support external collections, allowing users to query data from sources like Amazon S3, Google Cloud Storage, and Delta Lake as if it were internal, without duplication or ingestion. (Note: Apache Iceberg support is coming soon.) Multiple file formats are supported, including JSON, CSV/TSV, Parquet, and Avro. With dynamic path filters, users can specify which folders to include or exclude, turning storage structures into live, queryable datasets. This capability enables real-time analytics on external data while maintaining flexibility and cost efficiency.

### **COPY TO TRANSACTIONAL SERVICES AND EXTERNAL STORAGE**

Analytic Services provide powerful capabilities for writing calculated or transformed data back into the Couchbase transactional data store, enabling seamless integration between analytics and operational workflows. Users can copy entire analytics collections directly into the Couchbase Data Service, allowing insights derived from analysis to be persisted for real-time application use. Capella also supports transforming data stored in object stores like Amazon S3 and writing the results into the Data Service in JSON format (with CSV and Parquet coming soon), enabling efficient data refinement and operationalization. Moreover, the results of any analytics query – whether it involves joins across multiple collections, the use of views, built-in or user-defined functions – can be written back into the Data Service, ensuring that valuable analytics outcomes are directly actionable. This write-back capability works seamlessly whether targeting Capella's fully managed environment or a self-managed Couchbase Server cluster, giving organizations flexibility in how they deploy and utilize their data infrastructure. Data can also be written to external collections.

## **Ecosystem Integrations**

### **STREAMING INGESTION**

Couchbase Capella Analytics offer flexible streaming ingestion that allows users to ingest large volumes of data from diverse sources without custom code or connectors. It supports both initial load and continuous CDC ingestion into Columnar collections, with quick setup via the UI or SDK. Securing private connectivity with



Kafka ensures safe data transfer, and ingestion supports Avro, Protobuf, or JSON formats. Capella scales automatically by aligning with Kafka partitions, enabling high throughput. Users can leverage Kafka's ecosystem of open source connectors to bring in data from different sources, with easy no-code integration for existing Kafka deployments and support from AWS or Confluent for new users. Alongside Amazon S3 and Capella's transactional services, Kafka integration offers a seamless, scalable path for ecosystem-wide data ingestion.

### **BI TOOLS AND TABULAR ANALYTICS VIEWS**

Analytic Services integrate smoothly with BI tools like Tableau, Power BI, and Apache Superset through built-in JDBC/ODBC connectivity, enabling real-time data access in familiar platforms. It also supports tabular analytics views, which let users analyze JSON data in a relational format using SQL++ queries. Features like UNNEST simplify nested data, while view signatures and constraints map data types and keys for better usability. These views are not materialized, offering a real-time, efficient lens over JSON document data without duplication.

## **Management**

### **USER ORGANIZATION**

Couchbase Capella Analytics organize access and resources through a structured hierarchy: Organization > Projects > Clusters > Databases > Scopes > Collections. Organization members don't have data access by default; access is granted at the Project level by assigning roles to users or SSO groups (Teams). Each Cluster contains one or more compute nodes and hosts Databases, which group Scopes and Collections – the core units for querying and managing data. This structure ensures secure, scalable, and efficient data management.

### **MONITORING AND ALERTING**

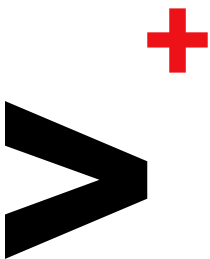
Analytic Services provide built-in observability with a curated set of performance, ingestion, and resource metrics displayed in a user-friendly Monitoring UI, powered by Prometheus – no setup required. Preconfigured alerts notify users via email with actionable insights. For enterprise needs, Prometheus-ready endpoints allow integration with third-party tools for centralized monitoring and root cause analysis.

### **QUERY ANALYSIS**

Couchbase Capella Analytics simplify query performance analysis with pre-built SQL++ functions that provide detailed stats on running and completed queries. Users can easily identify issues by running SQL queries for metrics like elapsed time, CPU, and memory usage, answering questions such as “What are my slowest queries?” or “Which queries exceed 1GB memory?” to optimize performance efficiently.

### **CLUSTER ON/OFF**

Analytic Services offer flexible cluster lifecycle management, allowing clusters to be turned on or off on demand, via scheduled times, or through management APIs. When off, compute is released, connections are terminated, and users are only billed for storage. Clusters can stay off for up to 30 days, after which they auto-restart for maintenance, with a two-day advance notification, helping users reduce costs while staying current.





## Operations

### SCALING

Couchbase Capella Analytics deliver flexible compute and storage scaling for optimal performance and cost efficiency. Users can scale horizontally by adjusting the number of compute nodes to meet workload demands, while vertical scaling (CPU/RAM per node) is fixed after cluster creation. Storage scaling is fully automated, expanding or contracting based on actual usage, with billing based only on what's used.

### BACKUP AND RESTORE

Analytic Services supports reliable backup and restore for disaster recovery and data protection. Each cluster can be backed up via Amazon S3 or Google Cloud Storage on demand or on a schedule, with backup intervals from 1 to 24 hours and retention up to 30 days. Backups can be restored to the same cluster within the same cloud provider, ensuring fast recovery with minimal downtime.

## Security

### SINGLE SIGN-ON

Couchbase Capella Analytics support single sign-on (SSO) to simplify authentication and enhance security by integrating with identity providers like Okta, Azure AD, Ping Identity, Google Workspace, OneLogin, and CyberArk. After setup, users can access Capella without managing separate credentials. The Teams feature lets admins map user groups to specific permission sets, streamlining access control and ensuring appropriate privileges.

### ROLE-BASED ACCESS CONTROL

Analytic Services offer flexible role-based access control (RBAC), allowing admins to create custom roles with specific privileges for fine-grained access control. Roles can be easily assigned to user accounts or privileges granted directly, ensuring users have the access they need. This approach enhances security, scales with user growth, and adapts to organizational needs.

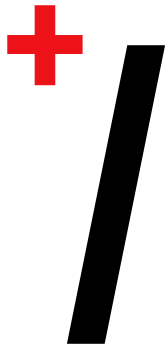


## CLOSING

---

Capella's capabilities make it a versatile choice for modern enterprises, combining the performance of in-memory computing with the flexibility of NoSQL databases for operational applications around the globe and at the edge, with the added power of real-time operational analytics. Capella's Transactional, Mobile, AI, and Analytic Services empowers developers to build scalable, high-performance applications with minimal operational complexity while ensuring security, availability, and efficiency in the cloud.

Agentic AI is the next revolution in how applications are designed, built, operated, and evolved. With increased variability in data input via human natural language, the block box nature of AI models, and the need to track agent responses over time, it has never been more important to ensure your data architecture is supporting the broad needs of the next generation of your applications. Find out more about why Couchbase is the developer data platform for critical applications in our AI world at [www.couchbase.com](http://www.couchbase.com) and via our documentation at [docs.couchbase.com](http://docs.couchbase.com).







Modern customer experiences need a flexible database platform that can power applications spanning from cloud to edge and everything in between. Couchbase's mission is to simplify how developers and architects develop, deploy and run modern applications wherever they are. We have reimaged the database with our fast, flexible and affordable cloud database platform Capella, allowing organizations to quickly build applications that deliver premium experiences to their customers – all with best-in-class price performance. More than 30% of the Fortune 100 trust Couchbase to power their modern applications.

For more information, visit [www.couchbase.com](https://www.couchbase.com) and follow us on X (formerly Twitter) @couchbase.

© 2025 Couchbase. All rights reserved.

