

Bases de Datos SQL Distribuidas: Google Spanner y CockroachDB

En el ámbito de las bases de datos distribuidas, Google Spanner y CockroachDB se destacan por sus capacidades avanzadas para manejar transacciones distribuidas con consistencia fuerte y escalabilidad horizontal.

Google Spanner

Uso de TrueTime para consistencia global

[Google Spanner](#) es una base de datos SQL distribuida que ofrece consistencia fuerte a escala global, una hazaña lograda mediante el uso de TrueTime. TrueTime es un sistema de sincronización de relojes desarrollado por Google que proporciona límites de tiempo precisos y bien definidos, minimizando la incertidumbre en la sincronización de relojes entre los servidores distribuidos.

TrueTime se basa en una combinación de relojes físicos y relojes de red atómicos que se sincronizan periódicamente mediante protocolos de sincronización de tiempo como GPS y NTP. Cada llamada a TrueTime devuelve un intervalo de tiempo `[earliest, latest]`, donde earliest representa el tiempo más temprano posible y latest el tiempo más tardío posible. Este intervalo ayuda a manejar la incertidumbre en la sincronización de relojes.

Para asegurar la consistencia global, Spanner utiliza estos intervalos de tiempo en sus transacciones. Cada transacción en Spanner se asocia con una marca de tiempo TrueTime, garantizando que todas las réplicas de los datos estén sincronizadas y proporcionando una ordenación total de las operaciones. Esto permite que Spanner ofrezca transacciones ACID distribuidas con consistencia fuerte y sin sacrificar la disponibilidad.

Particionamiento automático y replicación geográfica

Spanner automatiza el particionamiento y la replicación de datos, lo que facilita la gestión y escalabilidad de grandes volúmenes de datos distribuidos geográficamente. Los datos en Spanner se organizan en fragmentos llamados "particiones" que se distribuyen automáticamente entre diferentes nodos y centros de datos en todo el mundo.

El particionamiento automático se basa en la distribución de las claves primarias de las tablas. Spanner monitoriza continuamente el acceso y la carga de trabajo para ajustar dinámicamente las particiones, moviendo datos entre nodos y centros de datos según sea necesario para equilibrar la carga y optimizar el rendimiento.

La replicación geográfica asegura que los datos estén disponibles en múltiples ubicaciones. Spanner utiliza una replicación síncrona para mantener la consistencia fuerte entre las réplicas. Cada escritura se replica en todas las réplicas antes de confirmar la transacción al cliente. Esto garantiza que los datos estén disponibles y consistentes, incluso en caso de fallos en una o más réplicas.

CockroachDB

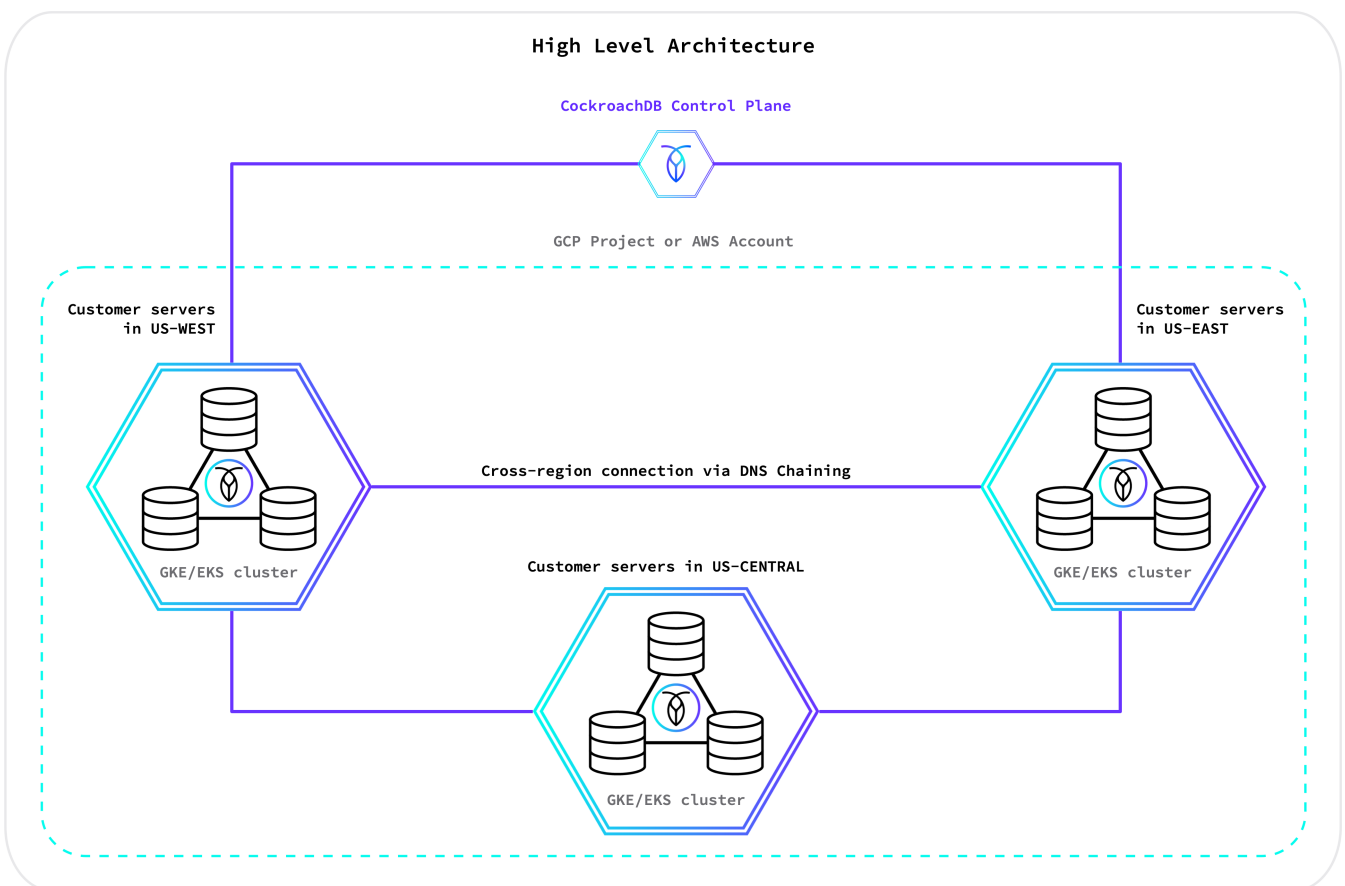
[CockroachDB](#) utiliza el protocolo de consenso Raft para asegurar la consistencia y la disponibilidad en su sistema distribuido. Raft es un algoritmo diseñado para ser comprensible y fácil de implementar, proporcionando las mismas garantías de consistencia que Paxos pero con una estructura más clara.

Raft divide el problema del consenso en tres subproblemas principales: elección de líder, replicación de logs y seguridad. En un clúster de CockroachDB, uno de los nodos actúa como líder, coordinando las operaciones de escritura y replicando las entradas de logs a los nodos seguidores. El líder asegura que todas las réplicas apliquen las entradas de logs en el mismo orden, garantizando la consistencia de los datos.

- Elección de líder: Si el líder actual falla, los nodos restantes llevan a cabo una elección para seleccionar un nuevo líder.
- Replicación de logs: El líder recibe todas las solicitudes de escritura, las añade a su log y las replica a los seguidores. Una vez que una mayoría de seguidores confirma la escritura, el líder aplica la entrada de log a su estado y responde al cliente.
- Seguridad: Raft asegura que cualquier entrada de log comprometida por un líder anterior será preservada por futuros líderes, manteniendo la consistencia y la integridad de los datos.

La imagen muestra una arquitectura de alto nivel para CockroachDB desplegada en múltiples regiones.

Ver <https://www.cockroachlabs.com/docs/cockroachcloud/architecture>



A continuación, se explican los componentes y cómo interactúan entre sí:

1. CockroachDB Control Plane:

- Es la capa de control de CockroachDB, responsable de la administración y orquestación de los clusters. Esta capa interactúa con los proyectos de GCP (Google Cloud Platform) o cuentas de AWS (Amazon Web Services) donde se despliegan los clusters de CockroachDB.

2. GKE/EKS Clusters (Clústeres de GKE/EKS):

- Los clusters se despliegan en Google Kubernetes Engine (GKE) o Amazon Elastic Kubernetes Service (EKS). Cada cluster contiene instancias de CockroachDB.

- En la imagen se muestran tres clusters distribuidos en diferentes regiones de Estados Unidos: US-WEST, US-CENTRAL y US-EAST.
3. Customer servers (Servidores del cliente):
- Los servidores del cliente se encuentran en las diferentes regiones mencionadas (US-WEST, US-CENTRAL y US-EAST) y están conectados a los clusters de CockroachDB desplegados en esas regiones.
4. Cross-region Connection via DNS Chaining (Conexión entre regiones mediante encadenamiento de DNS):
- Este mecanismo permite la conexión y comunicación entre los diferentes clusters en distintas regiones. El encadenamiento de DNS facilita la resolución de nombres y el enrutamiento del tráfico entre las regiones para asegurar la coherencia y disponibilidad de los datos a nivel global.

Funcionamiento general:

- Los clusters de CockroachDB están desplegados en múltiples regiones (US-WEST, US-CENTRAL y US-EAST) utilizando servicios gestionados de Kubernetes (GKE/EKS).
- El plano de control de CockroachDB gestiona estos clusters y se comunica con los proyectos de GCP o cuentas de AWS.
- Los servidores del cliente en cada región acceden a los clusters locales de CockroachDB, lo que permite una baja latencia de acceso a los datos.
- La conexión entre regiones se maneja mediante encadenamiento de DNS, asegurando que los datos sean accesibles y consistentes en todas las regiones, proporcionando alta disponibilidad y tolerancia a fallos.

Esta arquitectura permite a CockroachDB operar de manera distribuida y escalable a nivel global, asegurando que los datos estén siempre disponibles y sean consistentes en múltiples ubicaciones geográficas.

Transacciones distribuidas y consistencia fuerte

CockroachDB soporta transacciones ACID distribuidas, garantizando consistencia fuerte en un entorno distribuido. Utiliza una combinación de técnicas, incluyendo la serialización de transacciones y el uso de locks para asegurar la integridad de las operaciones.

Las transacciones en CockroachDB se manejan a través de un sistema de control de concurrencia multiversión (MVCC). Cada transacción recibe una marca de tiempo lógica, y las lecturas y escrituras se realizan en versiones específicas de los datos, permitiendo la ejecución concurrente de transacciones sin conflictos.

CockroachDB también implementa el protocolo de dos fases (2PC) para coordinar las transacciones distribuidas. En la primera fase, los nodos participantes preparan la transacción, asegurando que todos los cambios puedan ser aplicados. En la segunda fase, los nodos comprometen los cambios si todos los participantes están de acuerdo, garantizando la atomicidad y la durabilidad de la transacción.

Escalabilidad horizontal y tolerancia a fallos

CockroachDB está diseñado para escalar horizontalmente, permitiendo la adición de nodos sin interrupciones y con una mínima reconfiguración. Los datos se distribuyen automáticamente entre los nodos mediante particionamiento basado en rangos, lo que facilita la gestión de grandes volúmenes de datos y la distribución de la carga de trabajo.

La tolerancia a fallos en CockroachDB se logra mediante la replicación de datos y la robustez del protocolo Raft. Cada rango de datos se replica en varios nodos, y el protocolo Raft asegura que las operaciones de escritura se confirmen solo cuando una mayoría de réplicas ha aceptado la entrada de log. Esto permite que CockroachDB continúe operando incluso si uno o más nodos fallan.

En caso de fallo de un nodo, los nodos restantes pueden seguir proporcionando acceso a los datos replicados. Además, CockroachDB puede redistribuir automáticamente los datos desde los nodos fallidos a los nodos activos, asegurando que la carga de trabajo se mantenga equilibrada y que los datos estén siempre disponibles.

Prácticas avanzadas en optimización y escalabilidad de sistemas distribuidos

La optimización y escalabilidad en sistemas distribuidos son aspectos críticos que determinan el rendimiento y la capacidad del sistema para manejar cargas crecientes. A continuación, se exploran técnicas avanzadas para el escalado horizontal y vertical, la optimización de consultas y el uso eficiente de índices distribuidos, y el caching distribuido con herramientas como Redis y Memcached.

Técnicas de escalado horizontal vs. vertical

- **Escalado horizontal (Scale-Out):** Implica añadir más nodos al sistema para distribuir la carga de trabajo. Esta técnica es fundamental para sistemas distribuidos que requieren alta disponibilidad y tolerancia a fallos. Permite manejar grandes volúmenes de tráfico al distribuir las solicitudes entre múltiples servidores. Sin embargo, introduce desafíos en la gestión de la consistencia de datos y la coordinación entre nodos.
- **Escalado vertical (Scale-Up):** Consiste en añadir más recursos (CPU, memoria, almacenamiento) a un solo nodo. Aunque puede ser más simple de implementar, tiene limitaciones en términos de hasta dónde se puede escalar un solo servidor antes de que se vuelva ineficiente o costoso. Es útil para mejorar el rendimiento de aplicaciones que no son fácilmente distribuibles o que requieren operaciones intensivas en un único punto de procesamiento.

Optimización de consultas y uso eficiente de índices distribuidos

- **Optimización de consultas:** Involucra la reestructuración de consultas SQL para reducir la carga de procesamiento. Esto puede incluir el uso de subconsultas, eliminación de consultas redundantes, y la reescritura de consultas para aprovechar mejor los índices disponibles. En sistemas distribuidos, también es crucial minimizar la transferencia de datos entre nodos mediante consultas bien diseñadas que reduzcan el volumen de datos procesados y transmitidos.
- **Índices distribuidos:** Utilizar índices distribuidos de manera eficiente puede mejorar significativamente el rendimiento de las consultas. Los índices permiten un acceso rápido a los datos sin tener que escanear tablas completas. En sistemas distribuidos, es importante diseñar índices que minimicen la contención y el balance de carga. Índices particionados, por ejemplo, pueden mejorar la distribución de datos y reducir los cuellos de botella.

Caching distribuido (Redis, Memcached)

Redis: Es una base de datos en memoria que puede utilizarse como cache distribuido. Redis soporta estructuras de datos complejas y ofrece persistencia opcional, replicación, y alta disponibilidad mediante

Redis Sentinel y Redis Cluster. Su uso es ideal para aplicaciones que requieren un acceso rápido a datos frecuentemente consultados.

Memcached: Es una solución de caching distribuido simple y eficiente para datos en memoria. Memcached es ideal para aplicaciones que requieren una latencia extremadamente baja y donde los datos pueden ser volátiles. No ofrece las mismas garantías de durabilidad que Redis, pero su simplicidad y rendimiento lo hacen adecuado para muchos casos de uso.

Puedes revisar: [Redis and Memcached Explained](#).

Seguridad y privacidad en bases de datos distribuidas

La seguridad y privacidad son esenciales en sistemas distribuidos, donde los datos están repartidos en múltiples nodos y ubicaciones geográficas. A continuación, se detallan las prácticas avanzadas en autenticación, autorización, encriptación de datos y gestión de identidades.

Autenticación y autorización distribuidas

- **Autenticación:** Es el proceso de verificar la identidad de los usuarios y servicios que acceden al sistema. En sistemas distribuidos, es común utilizar OAuth, OpenID Connect y Kerberos para gestionar la autenticación. Estos protocolos permiten la autenticación segura y escalable en múltiples servicios y aplicaciones.
- **Autorización:** Una vez autenticados, los usuarios necesitan permisos específicos para realizar acciones en el sistema. El control de acceso basado en roles (RBAC) y el control de acceso basado en atributos (ABAC) son enfoques comunes para gestionar permisos en sistemas distribuidos. Estas técnicas permiten definir políticas detalladas que controlan qué recursos pueden ser accedidos por quién y bajo qué condiciones.

Encriptación de datos en tránsito y en reposo

- **Encriptación en tránsito:** Protege los datos mientras se transmiten entre nodos. TLS (Transport Layer Security) es el estándar para asegurar las comunicaciones en tránsito, asegurando que los datos no sean interceptados ni alterados durante la transmisión.
- **Encriptación en reposo:** Protege los datos almacenados en discos y bases de datos. Utilizar algoritmos de encriptación fuerte como AES (Advanced Encryption Standard) es esencial para proteger la integridad y privacidad de los datos almacenados. Las claves de encriptación deben gestionarse de manera segura, utilizando soluciones como Hardware Security Modules (HSMs) y servicios de gestión de claves (KMS).

Gestión de identidades y acceso

- **Gestión de identidades:** Implica la creación, mantenimiento y gestión de identidades de usuarios y servicios. Herramientas como LDAP (Lightweight Directory Access Protocol) y sistemas de gestión de identidades (IDM) proporcionan un marco centralizado para gestionar identidades en sistemas distribuidos.
- **Acceso condicional:** Implementar políticas de acceso condicional permite controlar el acceso a los recursos basándose en el contexto del usuario, como su ubicación, dispositivo utilizado y otros factores de riesgo. Esto proporciona una capa adicional de seguridad que adapta el acceso según las circunstancias específicas.

Monitoreo y mantenimiento

La monitorización y el mantenimiento de sistemas distribuidos son cruciales para asegurar su funcionamiento continuo y eficiente. A continuación, se discuten herramientas avanzadas de monitoreo y logging, estrategias de mantenimiento sin tiempo de inactividad, y gestión de backups y recuperación ante desastres.

- Herramientas de monitoreo y logging para sistemas distribuidos (Prometheus, Grafana)**
- **Prometheus**: Es una herramienta de monitoreo y alertas de código abierto que se utiliza ampliamente en sistemas distribuidos. Prometheus recopila y almacena métricas en una base de datos de series temporales, permitiendo un monitoreo detallado del rendimiento y el estado del sistema. Las alertas configurables permiten la detección proactiva de problemas antes de que afecten a los usuarios.
- **Grafana**: Complementa a Prometheus proporcionando una plataforma de visualización de datos poderosa y flexible. Grafana permite crear dashboards interactivos y personalizables que facilitan la visualización y el análisis de métricas en tiempo real. La integración de Grafana con Prometheus y otras fuentes de datos proporciona una visión holística del estado del sistema.

Estrategias de mantenimiento y actualización sin tiempo de inactividad

- Rolling Updates: Permiten actualizar componentes del sistema de manera gradual, uno a uno, sin interrumpir el servicio. Esta técnica es crucial para mantener la disponibilidad del sistema mientras se despliegan nuevas versiones de software o se aplican parches de seguridad.
- Blue-Green Deployments: Involucran tener dos entornos idénticos (azul y verde) donde uno está en producción y el otro es una copia exacta. Las actualizaciones se aplican al entorno inactivo y, una vez verificadas, el tráfico se redirige al nuevo entorno, minimizando el tiempo de inactividad.
- Canary Releases: Consisten en desplegar nuevas versiones de software a un subconjunto pequeño de usuarios para monitorear el comportamiento y detectar problemas antes de un despliegue completo. Esto permite identificar y resolver problemas potenciales con un impacto limitado en los usuarios.

Gestión de backups y recuperación ante desastres

- Backups regulares: Realizar copias de seguridad periódicas es esencial para proteger los datos contra pérdidas y corrupciones. Los backups deben ser almacenados en ubicaciones seguras y, preferiblemente, distribuidas geográficamente para asegurar la disponibilidad en caso de desastres regionales.
- Recuperación ante desastres (DR): Implica planificar y preparar procedimientos para restaurar el sistema a su estado operativo después de un desastre. Esto incluye tener planes detallados de recuperación, pruebas regulares de recuperación y la utilización de sitios de recuperación (DR sites) que puedan tomar el relevo en caso de fallo del sitio principal.

Desarrollo de aplicaciones distribuidas

El desarrollo de aplicaciones distribuidas requiere un enfoque diferente al del desarrollo de aplicaciones monolíticas tradicionales. A continuación, se analizan prácticas avanzadas en el diseño y desarrollo de microservicios, la orquestación y gestión de contenedores, y patrones de diseño específicos para aplicaciones distribuidas.

Diseño y desarrollo de microservicios

- **Microservicios:** Son un enfoque arquitectónico donde una aplicación se construye como un conjunto de servicios pequeños e independientes que se comunican entre sí. Cada microservicio es responsable de una funcionalidad específica y puede ser desarrollado, desplegado y escalado de manera independiente.
- **Descomposición de servicios:** Implica dividir una aplicación monolítica en componentes más pequeños y manejables. Esta descomposición debe hacerse de manera que cada servicio tenga una única responsabilidad y que las interacciones entre servicios sean minimizadas para reducir el acoplamiento.

Orquestación y gestión de contenedores (Docker, Kubernetes)

- **Docker:** Es una plataforma de contenedorización que permite empaquetar aplicaciones y sus dependencias en contenedores ligeros y portátiles. Docker facilita el despliegue consistente de aplicaciones en diferentes entornos, asegurando que el software se ejecute de la misma manera en cualquier lugar.
- **Kubernetes:** Es una plataforma de orquestación de contenedores que automatiza la gestión, el despliegue, el escalado y la operación de aplicaciones en contenedores. Kubernetes permite definir y gestionar infraestructuras complejas de microservicios, proporcionando capacidades avanzadas de recuperación ante fallos, escalabilidad automática y balanceo de carga.

Patrones de diseño para aplicaciones distribuidas (Circuit Breaker, Bulkhead)

- **Circuit Breaker:** Es un patrón de diseño que previene fallos en cascada en sistemas distribuidos. Cuando un servicio detecta que un componente ha fallado repetidamente, el Circuit Breaker abre el circuito para detener las solicitudes al componente fallido, permitiendo que el sistema se recupere y evitando que los fallos se propaguen.
- **Bulkhead:** Este patrón aísla diferentes partes del sistema en compartimentos independientes, asegurando que el fallo de una parte no afecte a otras. Cada compartimento tiene su propio pool de recursos, lo que previene que un fallo de un componente agote los recursos del sistema y afecte a otros componentes.

Nuevas tendencias en bases de datos distribuidas

La evolución de la tecnología de bases de datos distribuidas ha sido impulsada por la necesidad de manejar grandes volúmenes de datos, garantizar la alta disponibilidad, y proporcionar una escalabilidad efectiva. Recientemente, han emergido varias tendencias innovadoras en este campo, que incluyen bases de datos multi-modelo, la computación sin servidor (serverless), y la integración del aprendizaje automático con bases de datos distribuidas. Este informe explora estas tendencias, analizando sus características, beneficios y desafíos a nivel avanzado.

Bases de datos multi-modelo

Las bases de datos multi-modelo representan una de las innovaciones más significativas en el ámbito de las bases de datos distribuidas. Estas bases de datos son capaces de soportar múltiples modelos de datos (documento, gráfico, clave-valor, columna, etc.) dentro de una única plataforma. Esto permite a las organizaciones gestionar y consultar diferentes tipos de datos de manera más flexible y eficiente.

Características y beneficios

- **Flexibilidad de modelado:** Una base de datos multi-modelo permite a los desarrolladores elegir el modelo de datos más adecuado para cada aplicación específica. Por ejemplo, se puede usar un

modelo de documentos para almacenar datos semi-estructurados y un modelo de gráfico para gestionar relaciones complejas entre datos.

- Consolidación de infraestructura: Al soportar múltiples modelos de datos en una única plataforma, las bases de datos multi-modelo eliminan la necesidad de mantener múltiples sistemas de bases de datos separados. Esto reduce los costos operativos y la complejidad de la gestión de datos.
- Interoperabilidad de consultas: Las consultas pueden cruzar diferentes modelos de datos, proporcionando una visión unificada de la información. Esto es particularmente útil para análisis complejos y aplicaciones que requieren la integración de diversos tipos de datos.

Desafíos

- Complejidad de implementación: Soportar múltiples modelos de datos dentro de una única plataforma puede incrementar la complejidad del sistema, tanto en términos de diseño como de mantenimiento.
- Rendimiento: Optimizar el rendimiento para diferentes modelos de datos simultáneamente puede ser un desafío. Los sistemas deben equilibrar las necesidades de procesamiento y almacenamiento para diferentes tipos de datos sin degradar el rendimiento general.

Computación sin Servidor (Serverless) y bases de datos distribuidas

La computación sin servidor (serverless) es un paradigma donde los desarrolladores pueden construir y ejecutar aplicaciones sin gestionar la infraestructura subyacente. Este modelo ha sido aplicado a las bases de datos distribuidas, ofreciendo ventajas significativas en términos de escalabilidad, gestión y costos.

Características y beneficios

- Escalabilidad automática: Las bases de datos sin servidor escalan automáticamente en respuesta a la demanda. No hay necesidad de aprovisionar o gestionar servidores, lo que permite a las aplicaciones manejar cargas variables sin intervención manual.
- Modelo de pago por uso: Los usuarios solo pagan por los recursos utilizados, lo que puede reducir significativamente los costos en comparación con el aprovisionamiento de servidores permanentes que pueden estar infrautilizados.
- Simplificación de la gestión: La computación sin servidor elimina la necesidad de tareas administrativas como la configuración, el parcheo y el mantenimiento de servidores. Esto permite a los desarrolladores centrarse en el desarrollo de aplicaciones en lugar de en la gestión de infraestructura.

Desafíos

- Latencia: Las bases de datos sin servidor pueden experimentar latencias más altas debido a la naturaleza dinámica del aprovisionamiento de recursos. Esto puede ser un problema para aplicaciones que requieren respuestas en tiempo real.
- Control limitado: Al confiar en un proveedor de servicios en la nube, las organizaciones tienen un control limitado sobre la configuración y la optimización de la infraestructura subyacente. Esto puede ser una desventaja para aplicaciones con requisitos específicos de rendimiento o seguridad.

Integración de aprendizaje automático con bases de datos distribuidas

La integración del aprendizaje automático (machine learning, ML) con bases de datos distribuidas representa una convergencia poderosa que permite a las organizaciones extraer valor en tiempo real de grandes volúmenes de datos. Esta integración facilita la construcción de modelos de aprendizaje automático que pueden ser entrenados y desplegados directamente en la plataforma de bases de datos.

Características y beneficios

- **Análisis en tiempo real:** La capacidad de entrenar y ejecutar modelos de ML directamente en la base de datos permite realizar análisis en tiempo real sobre datos en movimiento. Esto es crucial para aplicaciones como la detección de fraudes, la personalización de contenidos y la predicción de fallos.
- **Reducción de latencia de datos:** Al integrar ML con bases de datos distribuidas, se elimina la necesidad de mover grandes volúmenes de datos entre diferentes sistemas para el entrenamiento y la inferencia de modelos. Esto reduce la latencia y mejora la eficiencia del procesamiento de datos.
- **Escalabilidad:** Las bases de datos distribuidas proporcionan una infraestructura escalable para el almacenamiento y procesamiento de datos. Esto es ideal para el entrenamiento de modelos de ML que requieren grandes volúmenes de datos y potencia de cálculo distribuida.

Desafíos

- **Complejidad técnica:** Integrar ML con bases de datos distribuidas requiere una infraestructura compleja y una sólida comprensión tanto de las técnicas de ML como de las arquitecturas de bases de datos distribuidas.
- **Rendimiento:** Optimizar el rendimiento del entrenamiento y la inferencia de modelos de ML en un entorno distribuido puede ser un desafío. Es necesario equilibrar la carga de trabajo y asegurar que los recursos estén disponibles cuando sea necesario.
- **Seguridad y privacidad:** El manejo de datos sensibles en bases de datos distribuidas con capacidades de ML plantea importantes preocupaciones de seguridad y privacidad. Es crucial implementar medidas robustas para proteger los datos y garantizar el cumplimiento de las regulaciones.

Ejercicios

1. Bases de Datos SQL Distribuidas: Google Spanner y CockroachDB:

- **Pregunta:** Explica cómo Google Spanner utiliza TrueTime para lograr la consistencia global y compare su enfoque con el uso de Raft en CockroachDB.
- **Ejercicio:** Propón un diseño de base de datos distribuida para una aplicación global de reserva de vuelos utilizando Google Spanner, detallando el uso de particionamiento automático y replicación geográfica.
- **Pregunta:** Analiza cómo Google Spanner y CockroachDB manejan las transacciones distribuidas y garantizan la consistencia fuerte.
- **Ejercicio:** Proporciona un diseño detallado para una aplicación de comercio electrónico global utilizando CockroachDB, explicando cómo manejaría la escalabilidad y la tolerancia a fallos.

2. Optimización y escalabilidad:

- **Pregunta:** Discute las técnicas de escalado horizontal vs. vertical y su aplicabilidad en diferentes escenarios.
- **Ejercicio:** Optimiza una consulta compleja en una base de datos distribuida utilizando índices distribuidos y explique los pasos y decisiones tomadas.
- **Pregunta:** Analiza las diferencias entre escalado horizontal y vertical en el contexto de una base de datos distribuida de alta demanda.
- **Ejercicio:** Optimiza una consulta SQL compleja en una base de datos distribuida utilizando técnicas avanzadas de indexación y caching.

1. Seguridad y privacidad en bases de datos distribuidas:

- Pregunta: Analiza los métodos de autenticación y autorización distribuidas en bases de datos distribuidas.
- Ejercicio: Diseña un esquema de seguridad para una base de datos distribuida que incluya encriptación de datos en tránsito y en reposo, y gestione identidades y accesos.
- Pregunta: Discute las técnicas avanzadas de autenticación y autorización en sistemas distribuidos.
- Ejercicio: Diseña un sistema seguro de gestión de identidades y accesos para una base de datos distribuida, incluyendo la encriptación de datos en tránsito y en reposo.

2. Monitoreo y mantenimiento:

- Pregunta: Compara las herramientas de monitoreo Prometheus y Grafana y sus usos en sistemas distribuidos.
- Ejercicio: Describe una estrategia de mantenimiento y actualización sin tiempo de inactividad para un sistema de bases de datos distribuidas, incluyendo la gestión de backups y recuperación ante desastres.
- Pregunta: Compara las capacidades de Prometheus y Grafana en el monitoreo de sistemas distribuidos.
- Ejercicio: Describe una estrategia de mantenimiento y actualización para una base de datos distribuida de una empresa global, incluyendo la gestión de backups y recuperación ante desastres.

3. Nuevas tendencias en bases de datos distribuidas:

- Pregunta: Explica el concepto de bases de datos multi-modelo y sus ventajas y desafíos.
- Ejercicio: Proporciona un diseño para integrar aprendizaje automático con una base de datos distribuida y explique cómo manejaría la escalabilidad y la eficiencia del procesamiento de datos.
- Pregunta: Explica cómo las bases de datos multi-modelo pueden beneficiar a las aplicaciones modernas.
- Ejercicio: Diseña una arquitectura de base de datos multi-modelo para una plataforma de análisis de big data, explicando cómo manejaría la integración de diferentes modelos de datos.

4. Comparación de diferentes tecnologías de bases de datos distribuidas:

- Pregunta: Compara y contrasta diferentes tecnologías de bases de datos distribuidas en términos de rendimiento y escalabilidad.
- Ejercicio: Realiza un análisis comparativo de dos bases de datos distribuidas (por ejemplo, Google Spanner vs. CockroachDB) en un escenario de implementación en una empresa de tecnología financiera.

5. Estudio de casos de implementación en empresas y aplicaciones reales:

- Pregunta: Analiza un caso de estudio real de implementación de una base de datos distribuida en una empresa global.
- Ejercicio: Diseña un caso de estudio teórico para una empresa que necesita migrar de una base de datos monolítica a una distribuida, explicando los pasos y las decisiones técnicas involucradas.

In []: *## Tus respuestas*

Ejercicio: Implementa un sistema de consistencia de sesión en Python. Simula varios clientes que realizan operaciones de lectura y escritura, y muestra cómo cada cliente ve una visión consistente de sus propias operaciones, aunque no necesariamente las operaciones de otros clientes.

In []: *## Tu respuesta*

Ejercicio: Implementa un sistema de control de concurrencia optimista en Python. Simula transacciones concurrentes que intentan leer y escribir datos, y utiliza un mecanismo de validación para resolver conflictos.

In []: *## Tu respuesta*

Ejercicio: Implementa un algoritmo simplificado de Raft en Python. Simula un clúster de nodos que eligen un líder y replican entradas de log.

In []: *## Tu respuesta*

Ejercicio: Implementa un sistema de sharding basado en rangos en Python. Distribuye las claves entre varias particiones y proporciona operaciones de escritura y lectura.

In []: *## Tu respuesta*

Ejercicio: Usando la biblioteca cassandra-driver, escribe un script en Python que se conecte a una base de datos Apache Cassandra, cree una tabla, y realice operaciones de inserción y consulta.

In []: *## Tu respuesta*

Ejercicio: Usando cockroachdb y sqlalchemy, escribe un script en Python que se conecte a una base de datos CockroachDB, crea una tabla, y realiza operaciones de inserción y consulta.

In []: *## Tu respuesta*

Ejercicio: Usando google-cloud-spanner, escribe un script en Python que se conecte a una base de datos Google Spanner, crea una tabla, y realiza operaciones de inserción y consulta.

In []: *## Tu respuesta*

Ejercicio: Implementa un mecanismo de caching distribuido usando Redis en Python. Escribe un script que se conecte a un clúster de Redis, almacena datos en el cache y los recupere.

In []: *## Tu respuesta*

Ejercicio: Implementa un sistema de encriptación de datos en tránsito usando TLS en una aplicación Python que se conecta a una base de datos PostgreSQL. Usa psycopg2 con configuraciones de SSL.

In []: *## Tu respuesta*

Ejercicio: Implementa un sistema de monitoreo usando Prometheus y Grafana. Escribe un script en Python que exponga métricas personalizadas y configúralo para que Prometheus las recopile.

In []: *## Tu respuesta*

Ejercicio: Diseña un sistema seguro de gestión de identidades y accesos para una base de datos distribuida en una empresa financiera, incluyendo la encriptación de datos en tránsito y en reposo, y políticas de acceso basadas en roles.

In []: *## Tu respuesta*