

Propuesta para garantizar la calidad/integración de la replicación

Para garantizar la calidad e integridad de la replicación de datos en todo el proceso, es importante implementar buenas prácticas y controles en varios puntos clave. Aquí se presentan algunas propuestas y hallazgos que se pueden considerar:

1. Validación de datos en el Origen y Destino

- **Comparación de datos fuente y replicados:** Realiza una comparación entre los datos originales en MySQL y los replicados en BigQuery para asegurarse de que no haya discrepancias. Algunas técnicas a implementar son:
 - i. **Conteo de filas:** Verificar que el número de registros en MySQL coincida con el número de registros en las tablas de BigQuery después de la replicación.
 - ii. **Pruebas de consistencia:** Crear queries específicas para verificar que los datos críticos (por ejemplo, claves primarias y campos de referencia) no presenten anomalías o faltantes en BigQuery.
- **Integridad referencial:** Asegurarse que las relaciones entre tablas (clave primaria y clave foránea) en MySQL se mantengan después de la replicación. Esto es importante para evitar pérdida o corrupción de datos relacionados.

2. Monitoreo en tiempo real de la replicación

- **Alertas y monitoreo continuo:** Implementar sistemas de monitoreo con herramientas como Cloud Monitoring para observar el estado del proceso de replicación. Configurar alertas en caso de fallas o retrasos en la replicación, como:
 - i. Fallos en la conectividad entre Datastream y MySQL o entre Datastream y BigQuery.
 - ii. Retrasos inusuales en el flujo de datos de CDC.

3. Control de consistencia en el CDC

- **Orden de eventos:** Asegurarse de que Datastream esté configurado para procesar los eventos de CDC en el orden correcto, especialmente cuando se trata de eventos de actualización o eliminación.
- **Respaldo de cambios:** Implementar un sistema de respaldo para los eventos de CDC. Si alguna parte de la cadena de replicación falla (Dataflow, Datastream, BigQuery), es recomendable contar con un registro adicional de las transacciones.

4. Pruebas automatizadas y validaciones de datos

- **Pruebas unitarias y de integración:** Desarrollar pruebas que verifiquen la integridad de los datos en cada paso del flujo de replicación (entrada desde GCS, procesamiento en MySQL y replicación en BigQuery). Automatizar estas pruebas para ejecutarlas periódicamente.

5. Gestión de errores

- **Idempotencia:** Asegurarse de que la replicación sea idempotente, es decir, que la misma operación pueda ejecutarse varias veces sin causar inconsistencias en los datos. Esto es crucial en sistemas de replicación donde los eventos de CDC pueden procesarse más de una vez.
- **Reintentos automáticos:** Configura políticas de reintento en caso de fallas temporales en la replicación. Esto permite que, si un lote de datos falla, el sistema pueda volver a intentarlo sin intervención manual.

Resumen de Propuestas

1. **Validación de datos:** Comparar el origen (MySQL) y el destino (BigQuery).
2. **Monitoreo en tiempo real:** Usar Cloud Monitoring para detectar errores o retrasos.
3. **Control de consistencia:** Asegurar el orden de eventos y monitorear la latencia de replicación.
4. **Pruebas automatizadas:** Ejecutar pruebas unitarias e integrales para validar la calidad de los datos.
5. **Gestión de errores:** Configurar reintentos automáticos y asegurar la idempotencia de la replicación.

Estas prácticas ayudarán a garantizar la calidad, integridad y confiabilidad de la replicación de datos de manera continua.