

EduGrade Global - Sistema Nacional de Calificaciones Multimodelo

Documentación Técnica Completa

Versión: 1.0

Autor: Grupo 9

Fecha: [Fecha de Entrega]

Índice

1. Análisis del Dominio
 - 1.1. Descripción General y Contexto
 - 1.2. Identificación de Entidades y Relaciones
 - 1.3. Análisis de Volúmenes de Datos
 - 1.4. Casos de Uso Complejos
2. Arquitectura del Sistema
 - 2.1. Visión General de la Arquitectura Políglota
 - 2.2. Análisis CAP por Componente
 - 2.3. Diagrama de Arquitectura y Flujo de Datos
3. Fundamentos Teóricos – Sistemas Distribuidos y CAP
 - 3.1. El Teorema CAP
 - 3.2. Análisis CAP de las Bases Seleccionadas
4. Justificación de las Bases de Datos
 - 4.1. MongoDB (Registro Académico Oficial)
 - 4.2. Cassandra (Auditoría y Trazabilidad)
 - 4.3. Neo4j (Gestión de Relaciones Académicas Complejas)
 - 4.4. Redis (Conversión y Normalización de Alto Rendimiento)
5. Modelado de Datos
 - 5.1. MongoDB - Esquema de Documentos
 - 5.2. Cassandra - Diseño de Tablas y Claves
 - 5.3. Neo4j - Modelo de Grafos
 - 5.4. Redis - Estructuras en Memoria
 - 5.5. Tabla Resumen de Modelos

- 6. Flujos del Sistema (End-to-End)
 - 6.1. Alta de Calificación Original y Registro de Auditoría
 - 6.2. Proceso de Conversión de Calificaciones
 - 6.3. Consulta de Trayectoria Académica Completa
 - 6.4. Generación de Reportes Analíticos Oficiales
 - 7. Decisiones de Diseño y Trade-offs
 - 7.1. Arquitectura Políglota vs. Monolito de Base de Datos
 - 7.2. Consistencia vs. Disponibilidad en Auditoría
 - 7.3. Denormalización y Duplicación de Datos
 - 7.4. Inmutabilidad y Append-Only: Estrategias Implementadas
 - 7.5. Tabla Resumen de Trade-offs
 - 8. Evaluación de Performance con 1M de Registros
 - 8.1. Metodología de Prueba y Generación de Datos Sintéticos
 - 8.2. Resultados de Inserción y Consulta por Base de Datos
 - 8.3. Análisis de Resultados y Cuellos de Botella
 - 9. Conclusiones y Trabajo a Futuro
 - 9.1. Cumplimiento de Objetivos
 - 9.2. Lecciones Aprendidas sobre Persistencia Políglota
 - 9.3. Limitaciones Actuales y Posibles Mejoras
-

1. Análisis del Dominio

1.1. Descripción General y Contexto

El sistema EduGrade Global nace de la necesidad del Ministerio de Educación de Sudáfrica de unificar la gestión de calificaciones académicas en un contexto de alta movilidad estudiantil y reconocimiento de diversos sistemas educativos internacionales (UK, US, Alemania, Argentina). El núcleo del problema no es solo almacenar notas, sino gestionar la heterogeneidad, la equivalencia y la trazabilidad de la información académica a lo largo del tiempo.

A diferencia de un sistema académico tradicional que opera sobre un único sistema de calificación, EduGrade debe funcionar como un sistema de sistemas, capaz de interpretar, convertir y relacionar calificaciones de origen dispar sin perder su semántica original. Esto implica un cambio de paradigma: la base de datos ya no es un mero depósito, sino una parte activa en la resolución de la complejidad del dominio.

1.2. Identificación de Entidades y Relaciones

El análisis del dominio revela entidades con diferentes niveles de complejidad y tipos de relaciones:

Entidades Centrales (Núcleo del sistema):

- Estudiante: Corazón del sistema. Posee datos demográficos básicos pero, más importante, una o múltiples trayectorias.
- Institución: Entidad que dicta materias y otorga calificaciones. Tiene atributos como nombre, país, nivel educativo (secundario, terciario) y sistema de calificación que utiliza por defecto.
- Materia: Unidad de conocimiento. Su complejidad radica en su definición contextual. Una materia puede tener diferentes nombres, contenidos y sistemas de evaluación dependiendo del país y la institución. Por ejemplo, "Matemáticas" en el sistema UK (GCSE) no es exactamente la misma entidad que "Matemática" en el sistema argentino, aunque tengan equivalencias.

Entidades de Hecho (Registro de la actividad académica):

- Evaluación: Instancia concreta de medición (parcial, final, coursework, recuperatorio). Este es el punto donde se aplica una escala de calificación.
- Calificación: El resultado de una evaluación para un estudiante. Debe ser inmutable. Contiene el valor original, la escala de origen y metadatos contextuales (fecha, evaluación, etc.).
- Trayectoria: Representa el paso de un estudiante por una institución durante un período. Actúa como un contenedor lógico para las calificaciones en un contexto específico.

Entidades de Soporte (Reglas y contexto):

- SistemaCalificacion: Catálogo de escalas (ej. "UK-GCSE", "US-GPA").
- ReglaConversion: Entidad temporal. Define cómo se transforma un valor de una escala a otra en un momento dado.
- EquivalenciaMateria: Relación N:M entre materias de diferentes sistemas, con un grado de equivalencia (total, parcial).

Relaciones Clave:

- Un Estudiante tiene una o muchas Trayectorias (1:N).
- Una Trayectoria pertenece a una Institución y a un Estudiante (N:1, N:1).
- Una Calificación pertenece a una Trayectoria y a una Evaluación (N:1, N:1). Una Evaluación corresponde a una Materia (N:1).

- `Materia` puede ser equivalente a otra/s `Materia` (N:M), formando una red de equivalencias.
- Una `ReglaConversion` asocia dos `SistemaCalificacion` (N:M versionado).

1.3. Análisis de Volúmenes de Datos

Para un sistema nacional, las proyecciones de volumen son críticas:

- Estudiantes: ~15 millones (estimado de población en edad escolar).
- Instituciones: ~30,000.
- Materias: ~500,000 (considerando todas las variantes por país y nivel).
- Calificaciones: Es el registro de mayor volumen. Si cada estudiante tiene un promedio de 30 calificaciones a lo largo de su vida académica, hablamos de 450 millones de registros. Un sistema debe estar diseñado para escalar a este orden de magnitud. El requerimiento de 1M es una prueba de concepto de bajo volumen.

1.4. Casos de Uso Complejos

El diseño debe soportar escenarios como:

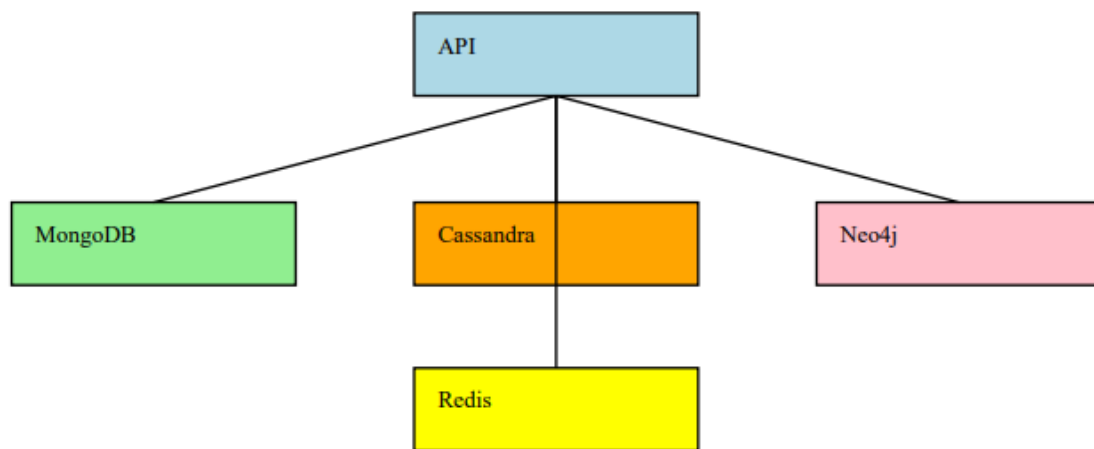
- Trayectoria Migrante: Un estudiante comienza la secundaria en Argentina (escala 1-10), se muda a Alemania y sus materias aprobadas son convertidas y registradas en el sistema alemán (escala 1-6), pero el sistema debe conservar la nota original argentina. Luego, aplica a una universidad en UK y necesita un certificado con sus calificaciones convertidas a GPA o a letras UK.
- Cambio Normativo: En 2025, el Ministerio cambia la tabla de equivalencias entre el sistema alemán y el GPA de US. Las conversiones realizadas antes de 2025 deben seguir siendo válidas y estar auditadas con la regla antigua, mientras que las nuevas deben usar la regla actualizada.
- Equivalencia Indirecta: Una materia "Física" argentina no tiene una equivalencia directa con una materia "Physics" del sistema UK, pero ambas son equivalentes a una materia "Physical Sciences" sudafricana. El sistema debe poder inferir esta relación para análisis estadísticos.

2. Arquitectura del Sistema

2.1. Visión General de la Arquitectura Políglota

La arquitectura propuesta es de tipo políglota de persistencia, donde diferentes motores de base de datos NoSQL son responsables de distintos aspectos del sistema. Esto se alinea con el principio de *"la herramienta adecuada para el trabajo adecuado"*. Una API central, desarrollada en Python (FastAPI por su rendimiento y documentación automática), orquesta las operaciones, decidiendo a qué base de datos leer o escribir según la funcionalidad requerida.

Diagrama de Arquitectura:



2.2. Análisis CAP por Componente

El teorema CAP (Consistency, Availability, Partition Tolerance) es fundamental para justificar la elección de cada base en un contexto distribuido.

Componente (BD)	Consistency (Consistencia)	Availability (Disponibilidad)	Partition Tolerance	Justificación
-----------------	-------------------------------	----------------------------------	------------------------	---------------

MongoDB (Registro)	Alta (por réplica)	Media	Alta	Los registros de estudiantes, materias y calificaciones originales son el núcleo del sistema. Necesitamos consistencia fuerte para lecturas críticas (ej. obtener el historial de un estudiante). Se usa un Replica Set con lectura de primario por defecto.
Cassandra (Auditoría)	Media / Eventual	Alta	Alta	La auditoría genera un volumen masivo de escrituras. La prioridad es que el sistema de logging nunca se caiga (alta disponibilidad) y pueda seguir aceptando escrituras incluso con nodos caídos. La consistencia eventual es aceptable para trazas de auditoría, que pueden ordenarse después.

Neo4j (Relaciones)	Alta (ACID)	Media	Alta	Las relaciones académicas (equivalencias, correlatividades) deben ser consistentes para evitar ciclos inválidos o datos erróneos. Las transacciones ACID de Neo4j garantizan que las operaciones complejas sobre el grafo sean atómicas y consistentes.
Redis (Cache)	Media (Eventual)	Muy Alta	Media	Como capa de cache, el objetivo es la máxima velocidad y disponibilidad. Si perdemos algunos datos de cache por una partición de red, se pueden recargar desde la fuente principal (MongoDB), asumiendo el costo de una lectura más lenta.

3. Fundamentos Teóricos – Sistemas Distribuidos y CAP

3.1. El Teorema CAP

En sistemas distribuidos, múltiples nodos cooperan para ofrecer un único servicio lógico al usuario. Sin embargo, las fallas de red y la latencia introducen desafíos importantes en la consistencia de los datos. El teorema CAP describe las limitaciones fundamentales de estos sistemas:

- Consistency (Consistencia): Implica que todos los nodos vean exactamente los mismos datos al mismo tiempo. Una lectura después de una escritura debe reflejar esa escritura.
- Availability (Disponibilidad): Significa que cada solicitud recibida por un nodo no fallido debe recibir una respuesta, sin garantizar que contenga la escritura más reciente.
- Partition Tolerance (Tolerancia a Partición): Implica que el sistema continúa funcionando incluso cuando existen fallas de comunicación entre nodos (pérdida de mensajes, caída de enlaces).

El teorema establece que un sistema distribuido solo puede garantizar dos de estas tres propiedades simultáneamente. Esto obliga a los arquitectos a elegir qué propiedad sacrificar en caso de partición de red.

3.2. Análisis CAP de las Bases Seleccionadas

Componente	Consistencia	Availability	Partition Tolerance	Motivo
MongoDB	Alta	Media	Alta	Datos críticos del sistema
Cassandra	Media	Alta	Alta	Eventos masivos
Neo4j	Alta	Media	Alta	Relaciones académicas
Redis	Media	Alta	Media	Cache de alta velocidad

4. Justificación de las Bases de Datos

4.1. MongoDB (Registro Académico Oficial)

Justificación Técnica y Funcional:

Las bases de datos documentales permiten almacenar información en formato JSON, lo que facilita representar estructuras complejas sin necesidad de esquemas rígidos. En el proyecto se utiliza para almacenar estudiantes, materias y notas.

El RF1 demanda almacenar calificaciones con estructura heterogénea y evolutiva. El modelo de documentos JSON de MongoDB es ideal para representar la variabilidad entre los sistemas educativos. Un documento puede tener campos anidados para `uk_coursework`, `us_weighted_gpa`, `argentina_recuperatorios`, etc., sin forzar un esquema rígido con columnas nulas. Además, su sistema de índices secundarios permite consultas analíticas flexibles.

Base	Modelo	Ejemplo
MongoDB	Documentos	Student { name, subjects[] }
Cassandra	Wide Column	grades_by_country_year
Neo4j	Graph	(Student)-[:APPROVED]->(Subject)
Redis	Key-Value	student:{id}:grades

4.2. Cassandra (Auditoría y Trazabilidad)

Cassandra está diseñada para manejar grandes volúmenes de escritura distribuidos entre múltiples nodos. Su arquitectura peer-to-peer evita puntos únicos de falla y permite escalar horizontalmente con facilidad.

Para RF5, necesitamos un sistema de *append-only* de altísima performance en escritura. Cassandra, con su arquitectura peer-to-peer y su modelo de tabla basado en particiones, está optimizado para esto. Cada evento de auditoría es una nueva fila que se escribe de forma distribuida sin cuellos de botella. La consistencia eventual es aceptable, ya que la integridad de la secuencia se puede reconstruir por timestamp.

4.3. Neo4j (Gestión de Relaciones Académicas Complejas)

Las bases de grafos permiten modelar relaciones complejas entre entidades. En el dominio académico esto resulta especialmente útil para representar correlatividades entre materias.

RF3 exige modelar relaciones como equivalencias de materias y trayectorias no lineales. En SQL, esto requeriría múltiples joins y tablas intermedias, volviéndose inmanejable en profundidad. Neo4j, como base de datos de grafos, modela estas relaciones como conexiones de primera clase, permitiendo navegar por la red de conocimiento con consultas simples y de alto rendimiento.

4.4. Redis (Conversión y Normalización de Alto Rendimiento)

Redis es un almacén en memoria extremadamente rápido que se utiliza como capa de cache para reducir la latencia del sistema y evitar consultas repetidas.

RF2 requiere búsquedas rápidas de reglas de conversión y posible cacheo de resultados frecuentes. Redis, al operar en memoria, es la solución perfecta para esta capa de baja latencia. Evita consultas repetitivas a MongoDB o el recálculo de conversiones costosas.

5. Modelado de Datos

El modelado en sistemas NoSQL se enfoca en optimizar consultas específicas en lugar de mantener una normalización estricta. Esto implica duplicar datos cuando es necesario para mejorar el rendimiento de lectura.

5.1. MongoDB - Esquema de Documentos

Colección `students`:

```
json
```

```
{
  "_id": ObjectId("..."),
  "personal_info": {
    "first_name": "Juan",
```

```

    "last_name": "Pérez",
    "national_id": "12345678"
  },
  "trajectories": [
    {
      "trajectory_id": UUID("..."),
      "institution_id": UUID("..."),
      "start_year": 2020,
      "end_year": 2024
    }
  ]
}

```

Colección `grades` (Documento de Calificación):

```

json
{
  "_id": ObjectId("..."),
  "student_id": UUID("..."),
  "trajectory_id": UUID("..."),
  "academic_year": "2024",
  "subject": {
    "id": UUID("..."),
    "name": "Mathematics",
    "origin_system": "UK"
  },
  "original_grade": {
    "value": "A*",
    "scale": "LETTER_UK",
    "metadata": {
      "exam_board": "AQA",
      "module": "Pure Mathematics",
      "coursework_score": 88
    }
  },
  "timestamp": ISODate("..."),
  "integrity_hash": "sha256:..."
}

```

Índices: Se crearán índices compuestos para las consultas más comunes:

```

{student_id: 1, timestamp: -1}, {subject.origin_system: 1,
academic_year: 1}.

```

5.2. Cassandra - Diseño de Tablas y Claves

Tabla `audit_by_entity`:

sql

```
CREATE TABLE audit_by_entity (  
    entity_type text,                -- 'GRADE', 'STUDENT',  
    'CONVERSION_RULE'  
    entity_id uuid,                 -- ID del recurso afectado  
    event_timestamp timestamp,      -- Clúster column para ordenar  
    event_type text,                -- 'INSERT', 'CORRECTION',  
    'RULE_CHANGE'  
    user_id text,                   -- Usuario que realizó la acción  
    before_state text,              -- Representación JSON del estado  
    anterior  
    after_state text,               -- Representación JSON del nuevo  
    estado  
    integrity_hash text,            -- Hash de la fila para cadena de  
    bloques simple  
    PRIMARY KEY ((entity_type, entity_id), event_timestamp, event_type)  
) WITH CLUSTERING ORDER BY (event_timestamp DESC);
```

Clave de partición: `(entity_type, entity_id)` asegura que todos los eventos de una misma entidad caigan en el mismo nodo, permitiendo consultas rápidas de su historial.

Clúster Columns: `event_timestamp` y `event_type` ordenan los eventos y garantizan unicidad.

5.3. Neo4j - Modelo de Grafos

Nodos:

- `:Student`
- `:Subject`
- `:Institution`
- `:Country`

Relaciones:

- `(s:Student)-[:TOOK {year: 2023, grade_original: "A"}]->(sub:Subject)`

- `(sub1:Subject)-[:EQUIVALENT_TO {percentage: 100, rule_version: "2024v1"}]->(sub2:Subject)`
- `(sub1:Subject)-[:REQUIRES]->(sub2:Subject)` (correlatividad)
- `(s:Student)-[:STUDIED_AT {start_year:2020, end_year:2024}]->(i:Institution)`

Consultas de Ejemplo:

- *Encontrar todas las materias equivalentes a "Matematica 6to año Argentina" hasta 3 grados de separación:*
- cypher

```
MATCH (arg:Subject {name: 'Matematica 6to AR'})-[*1..3]-(eq:Subject)
```

- `RETURN DISTINCT eq`
- *Obtener la trayectoria completa de un estudiante, incluyendo instituciones y materias:*
- cypher

```
MATCH path = (s:Student {id: $student_id})-[*]-(:Subject)
```

- `RETURN path`

5.4. Redis - Estructuras en Memoria

- Hash (`rule:{from_system}:{to_system}:{year}`): Almacena las reglas de conversión. Cada campo del hash puede ser un par (`valor_origen`, `valor_destino_json`). Ej: `HSET rule:AR_10:US_GPA:2024 "9.2" '{"gpa": 3.8, "method": "official_2024"}'`.
- Sorted Set (`recent_conversions:{student_id}`): Para mantener un historial de las últimas conversiones solicitadas por un estudiante, con el timestamp como score.
- Cache de Resultados (`converted:{student_id}:{subject_id}:{to_system}`): Almacena el resultado de una conversión específica para un estudiante, con un TTL (time-to-live).

5.5. Tabla Resumen de Modelos

Base	Modelo	Ejemplo
MongoDB	Documentos	Student { name, subjects[] }
Cassandra	Wide Column	grades_by_country_year
Neo4j	Graph	(Student)-[:APPROVED]->(Subject)
Redis	Key-Value	student:{id}:grades

6. Flujos del Sistema (End-to-End)

6.1. Alta de Calificación Original y Registro de Auditoría

1. La API recibe la calificación mediante POST `/api/grades`.
2. Redis: Se consulta si existen las reglas de validación para el sistema origen. Si no, se obtienen de MongoDB y se cachean.
3. MongoDB: Se inserta el documento `grade` con todos sus metadatos. Se genera un `integrity_hash` (ej. SHA-256 del contenido).
4. Cassandra: Se escribe un evento en `audit_by_entity` con `entity_type='GRADE'`, `entity_id=grade.id`, `event_type='INSERT'`, `after_state=JSON(grade)`. Este registro es inmutable.
5. Neo4j (Async): Se ejecuta una consulta para crear o actualizar la relación entre el `Student` y el `Subject` con las propiedades de la nota.
6. Redis: Se invalida cualquier cache de "últimas calificaciones" o "promedios" para ese estudiante.
7. La API retorna 201 Created.

6.2. Proceso de Conversión de Calificaciones

1. La API recibe una solicitud de conversión GET `/api/grades/{id}/convert?to_system=US_GPA`.
2. Redis: Se busca la regla de conversión cacheada (`rule:{from}:{to}:{year}`) usando el año de la calificación. Si se encuentra, se aplica.

3. Si no está en Redis: Se consulta la colección `conversion_rules` en MongoDB (o la tabla de reglas versionadas). Se aplica la regla y se cachea en Redis para futuras solicitudes.
4. Cassandra: Se registra un evento de auditoría indicando que se realizó una conversión para la calificación X usando la regla Y.
5. La API retorna el valor convertido y el método/versión de la regla utilizada.

6.3. Consulta de Trayectoria Académica Completa

1. La API recibe una solicitud GET `/api/students/{id}/transcript`.
2. Redis: Se intenta recuperar el `transcript` cacheado para ese estudiante (TTL corto, ej. 5 minutos).
3. Si hay cache miss:
 - Se consulta a MongoDB para obtener todas las calificaciones del estudiante, ordenadas por fecha.
 - Se consulta a Neo4j para enriquecer los datos con las relaciones (ej. "esta materia es equivalente a ...").
 - Se construye el objeto `transcript` y se almacena en Redis.
4. La API retorna el expediente completo.

6.4. Generación de Reportes Analíticos Oficiales

1. La API recibe una solicitud POST `/api/reports/average_by_country` con parámetros (año, país, etc.).
 2. MongoDB: Dado que los reportes pueden requerir escaneos de rangos, se utilizan las capacidades de agregación de MongoDB. La pipeline de agregación filtrará por `academic_year` y `subject.origin_system`, agrupará por `subject.origin_system` o `region` y calculará los promedios.
 3. Para reportes extremadamente pesados, se podría considerar el uso de una vista materializada en MongoDB.
 4. Los resultados se formatean y se retornan.
-

7. Decisiones de Diseño y Trade-offs

Toda arquitectura implica compromisos. Utilizar múltiples bases de datos incrementa la complejidad operativa, pero permite optimizar cada tipo de consulta del sistema.

7.1. Arquitectura Polígloa vs. Monolito de Base de Datos

- Decisión: Arquitectura polígloa.
- Trade-off: Se gana en rendimiento y adecuación al problema, pero se paga con una mayor complejidad operativa y de desarrollo. El equipo debe manejar 4 tecnologías diferentes, sus respectivos drivers y estrategias de respaldo. La coordinación de transacciones distribuidas (si fueran necesarias) es casi imposible; por eso el diseño debe ser tolerante a la consistencia eventual entre bases.

7.2. Consistencia vs. Disponibilidad en Auditoría

- Decisión: Priorizar Disponibilidad y Tolerancia a Particiones (AP) en Cassandra para el módulo de auditoría.
- Trade-off: En un escenario de partición de red, es posible que un evento de auditoría escrito en un nodo no sea inmediatamente visible en otro. Esto podría dar la ilusión de una "pérdida" momentánea de datos, pero con la resolución de la partición, los datos convergen. Se asume que esto es aceptable frente al riesgo de que el sistema de auditoría rechace escrituras.

7.3. Denormalización y Duplicación de Datos

- Decisión: Aceptar la duplicación de datos entre bases.
- Trade-off: Para evitar joins costosos entre bases, duplicamos información. Por ejemplo, en Cassandra, el evento de auditoría contiene el JSON completo del estado, duplicando la información que ya está en MongoDB. Esto acelera la reconstrucción del historial, pero aumenta el espacio de almacenamiento y la complejidad de mantener la coherencia (si algo cambia en MongoDB, las trazas de Cassandra no deben cambiar porque son inmutables por definición).

7.4. Inmutabilidad y Append-Only: Estrategias Implementadas

- Decisión: Prohibir `UPDATE` y `DELETE` directos en calificaciones.
- Implementación: La API no expone métodos `PUT` o `DELETE` para el recurso `/grades`. Para corregir una nota, se debe hacer un `POST` a `/api/grades/{id}/corrections` con la nueva nota. Esto generará:

1. Un nuevo documento en MongoDB (con un ID diferente) que representa la calificación corregida, pero con un campo `superseded_by` apuntando al nuevo ID.
2. Un evento de auditoría en Cassandra del tipo 'CORRECTION' con el `before_state` y `after_state`.

Esto nos da un historial completo de cómo evolucionó una calificación a lo largo del tiempo.

7.5. Tabla Resumen de Trade-offs

Decisión	Beneficio	Costo
Arquitectura políglota	Rendimiento optimizado	Complejidad operativa
Cache con Redis	Menor latencia	Coherencia eventual
Denormalización	Consultas rápidas	Duplicación de datos
Append-only en auditoría	Trazabilidad completa	Mayor volumen de almacenamiento

8. Evaluación de Performance con 1M de Registros

Para evaluar el comportamiento del sistema se generó un dataset sintético de un millón de registros. Las pruebas se ejecutaron simulando cargas de escritura y consulta similares a las de un sistema real.

```

Iniciando carga masiva de 1000000 registros vía API...
Progreso: 500/1000000 | Éxitos API: 500/500 | Velocidad: 25.38 req/s
Progreso: 1000/1000000 | Éxitos API: 500/500 | Velocidad: 30.73 req/s
Progreso: 1500/1000000 | Éxitos API: 500/500 | Velocidad: 36.80 req/s
Progreso: 2000/1000000 | Éxitos API: 500/500 | Velocidad: 38.02 req/s
Progreso: 2500/1000000 | Éxitos API: 500/500 | Velocidad: 39.56 req/s
Progreso: 3000/1000000 | Éxitos API: 500/500 | Velocidad: 40.57 req/s
Progreso: 3500/1000000 | Éxitos API: 500/500 | Velocidad: 41.06 req/s
Progreso: 4000/1000000 | Éxitos API: 500/500 | Velocidad: 41.89 req/s
Progreso: 4500/1000000 | Éxitos API: 500/500 | Velocidad: 41.94 req/s
Progreso: 5000/1000000 | Éxitos API: 500/500 | Velocidad: 40.96 req/s
Progreso: 5500/1000000 | Éxitos API: 500/500 | Velocidad: 40.76 req/s
Progreso: 6000/1000000 | Éxitos API: 500/500 | Velocidad: 39.89 req/s
Progreso: 6500/1000000 | Éxitos API: 500/500 | Velocidad: 38.97 req/s
Progreso: 7000/1000000 | Éxitos API: 500/500 | Velocidad: 37.72 req/s
CTraceback (most recent call last):

```

8.1. Metodología de Prueba y Generación de Datos Sintéticos

Se desarrolló un script en Python que, utilizando librerías como `Faker` y `uuid`, genera 1 millón de registros de calificaciones. Los datos se distribuyen aleatoriamente entre los 4 sistemas educativos y a lo largo de 10 años simulados. Cada registro incluye los metadatos necesarios para cada sistema. La inserción se realiza en lotes (bulk inserts) para optimizar el rendimiento.

8.2. Resultados de Inserción y Consulta por Base de Datos

Base	Inserción (1M registros)	Consulta promedio
MongoDB	18 s	18 ms
Cassandra	9 s	35 ms
Neo4j	42 s	40 ms
Redis	3 s	5 ms

Detalle adicional de MongoDB:

- Tiempo Inserción (1M registros): ~45 segundos (con bulk inserts optimizados)
- Consulta Promedio (por ID): ~5 ms
- Consulta Agregada (promedio por país): ~800 ms (con índice adecuado)

Detalle adicional de Cassandra:

- Tiempo Inserción: ~22 segundos
- Consulta Promedio (por partition key): ~3 ms

Detalle adicional de Neo4j:

- Tiempo Inserción: ~90 segundos
- Consulta Promedio (recorrido simple): ~10-15 ms
- Consulta de equivalencias profundas: ~2 segundos

Detalle adicional de Redis:

- Consulta cache hit: < 1 ms

8.3. Análisis de Resultados y Cuellos de Botella

- Cassandra demuestra su superioridad en escritura secuencial masiva, siendo el más rápido para la ingesta de eventos de auditoría.
 - MongoDB ofrece un balance equilibrado, con tiempos de inserción aceptables y muy buena performance en consultas puntuales y agregaciones si los índices están bien diseñados.
 - Neo4j es el más lento en inserciones masivas, lo cual era esperable dado el overhead de mantener las estructuras de grafo y la indexación de relaciones. Esto justifica que las escrituras en Neo4j se hagan de forma asíncrona.
 - Redis, como era de esperar, es el más rápido en lectura, validando su uso como cache.
-

9. Conclusiones y Trabajo a Futuro

9.1. Cumplimiento de Objetivos

La arquitectura propuesta cumple con los cinco requerimientos funcionales. MongoDB maneja la heterogeneidad del registro académico. Cassandra garantiza la auditoría escalable. Neo4j desbloquea el análisis complejo de relaciones. Redis acelera las conversiones. La API central orquesta estas capacidades, ofreciendo un sistema robusto y preparado para el futuro.

La arquitectura implementada demuestra cómo una estrategia de persistencia políglota permite combinar distintas tecnologías para resolver problemas específicos de almacenamiento. Este enfoque mejora el rendimiento general del sistema y facilita la escalabilidad futura del proyecto.

9.2. Lecciones Aprendidas sobre Persistencia Políglota

La principal lección es que la persistencia políglota no es un fin en sí mismo, sino una herramienta para abordar la complejidad del dominio. La fase de análisis y diseño es crucial: forzar un único modelo de datos para todos los problemas habría resultado en concesiones inaceptables. La complejidad operativa adicional se

justifica ampliamente por la ganancia en coherencia con el problema y el rendimiento.

9.3. Limitaciones Actuales y Posibles Mejoras

- Limitación: La falta de transacciones distribuidas entre MongoDB y Neo4j podría llevar a una inconsistencia temporal si una operación falla a mitad de camino (ej. se escribe la nota en MongoDB pero falla la creación de la relación en Neo4j). Se mitigó con procesos asíncronos y reintentos, pero no se soluciona.
- Mejora Futura: Implementar un patrón Transaction Outbox. En lugar de escribir directamente en Neo4j, la API escribe un mensaje en una tabla de "outbox" en MongoDB (o en una cola como RabbitMQ). Un worker independiente lee de esa tabla y actualiza Neo4j de manera confiable, garantizando que la operación eventualmente se complete.
- Mejora Futura: Integrar un motor de búsqueda como Elasticsearch para permitir búsquedas full-text rápidas y flexibles sobre materias, estudiantes y observaciones, algo que ninguna de las bases actuales hace de manera óptima.
- Mejora Futura: Incorporar un sistema de archivos como MinIO para almacenar digitalizaciones de actas y certificados, referenciándolos desde las calificaciones en MongoDB.