

HCL (BIG DATA)

CASO PRACTICO I UD5



ALUMNO CESUR

25/26

Alejandro Muñoz de la Sierra

PROFESOR

Óscar González Núñez

INTRODUCCION

Este proyecto nace de una realidad del sector tecnológico: el gran volumen de datos de internet. Diseñaremos la arquitectura de "FARO", un motor de búsqueda propio. No competiremos con Google de inmediato. Queremos demostrar nuestra capacidad para gestionar la ingestión y consulta de millones de páginas web con eficiencia.

El equipo estableció una premisa clara durante el análisis inicial. El modelo relacional clásico no funciona para este fin. Un buscador rastrea millones de páginas diarias. Cada una tiene una estructura distinta, como blogs, foros o tiendas. Forzar la web en un esquema rígido de tablas con Oracle o MySQL causaría fallos graves. Generaría cuellos de botella y dificultaría el mantenimiento.

Basamos la arquitectura de FARO en MongoDB. Es una base de datos NoSQL orientada a documentos. Su sistema sin esquema fijo usa documentos BSON. Esto permite guardar cada página web con su estructura natural. Los datos crecen de forma exponencial. MongoDB ofrece una escalabilidad horizontal nativa superior a las bases SQL por el mismo coste.

ANÁLISIS DE NECESIDADES Y MODELADO DE DATOS

Aquí detallo las fases de diseño, modelado e implementación de este clúster.

El éxito de una base NoSQL no depende de crear colecciones sin orden e insertar archivos JSON. Aplicamos una regla estricta. Diseñamos el modelo según la lectura de los datos. En SQL normalizamos para evitar repeticiones. En NoSQL desnormalizamos para priorizar la velocidad de lectura.

1.1. Identificación Exhaustiva de Necesidades

Variedad Estructural (Polimorfismo): Guardamos artículos de prensa, productos de tiendas y vídeos. Estos tienen datos distintos como autor, fecha, precio o stock. MongoDB permite documentos con esquemas diferentes en una misma colección sin errores.

Velocidad de Respuesta Crítica: El usuario espera resultados en menos de 200 milisegundos. Evitaremos operaciones complejas de cruce de datos durante la consulta. Estas operaciones usan `\$lookup` en MongoDB.

Datos Enriquecidos: Almacenaremos el texto de la web y sus métricas técnicas. Esto incluye velocidad de carga y optimización móvil. Usaremos estos datos para influir en la clasificación de resultados.

1.2. Diseño de Colecciones y Patrones de Modelado

Estructuramos la base de datos principal en dos grandes colecciones. Aplicamos patrones de diseño NoSQL distintos según la necesidad:

A) Colección Principal: web_pages (Patrón de Documento Embebido)

No usamos tablas separadas para URL, SEO y métricas como en SQL. Incrustamos todo en un único documento. El buscador lee un solo bloque del disco duro para obtener la información. El proceso es muy rápido.

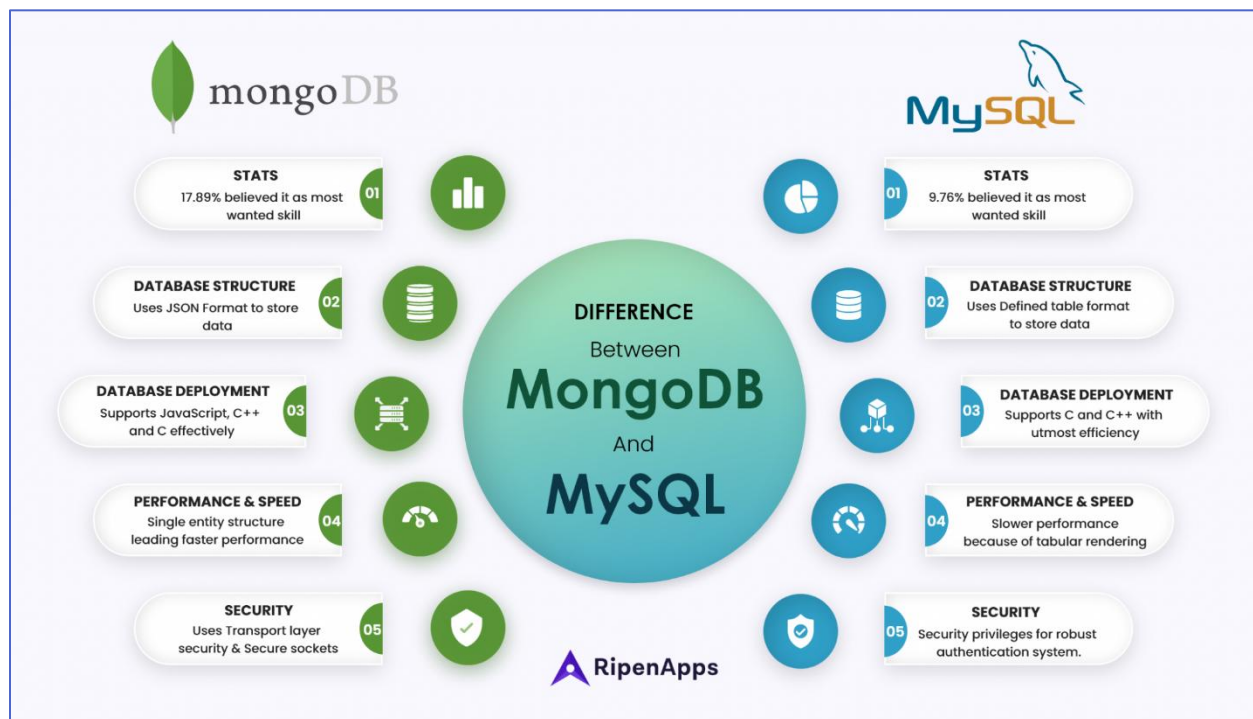
```

C:\> Users > sibko > Documents > GitHub > Cesur-DAM > Asignaturas 2 > HCL (Big Data) > UD5 > Caso_Practico_1 > BIGDATAUD5_1.j
1 // Ejemplo del modelo de documento en la colección 'web pages'
2 {
3   "_id": ObjectId("652b4c9a3f..."),
4   "url": "https://www.ejemplo.es/tecnologia/5g",
5   "dominio": "ejemplo.es",
6   "hash_url": "a1b2c3d4e5f6...", // índice único para evitar indexar dos veces la misma web
7   "contenido": {
8     "titulo": "El avance del 5G en zonas rurales",
9     "cuerpo_texto": "El despliegue de la red móvil avanza a buen ritmo...",
10    "idioma": "es-ES",
11    "tiene_multimedia": true
12  },
13  "metadatos_seo": {
14    "keywords": [
15      "5G",
16      "tecnología",
17      "redes",
18      "rural"
19    ],
20    "fecha_ultimo_rastreo": ISODate("2026-02-21T09:00:00Z"),
21  },
22  "metricas_tecnicas": {
23    "latencia_media_ms": 12,
24    "optimizado_movil_amp": true
25  },
26  "autoridad": {
27    "page_rank_score": 8.7, // Este dato lo pre-calculamos por la noche
28    "trust_score": 9.1
29  },
30  "enlaces_salientes": [
31    "https://www.otraweb.es/info",
32    "https://www.gob.es/telecomunicaciones"
33  ]
34 }
```

B) Colección Secundaria: search_analytics (Patrón de Referencia)

Guardamos los **registros estadísticos** de búsqueda en otra colección. Esto evita que los documentos de las webs crezcan demasiado. Aquí usamos referencias y guardamos el `_id` de la web.

Un usuario busca "5G" y hace clic en la web de ejemplo. Guardamos ese evento en `search_analytics` para alimentar nuestro algoritmo. Esta colección acepta escrituras masivas constantes. No bloquea las lecturas de las búsquedas.



IMPLEMENTACIÓN TÉCNICA DEL BUSCADOR

Un buscador es más que un almacén digital. Es un motor de recuperación de información. Explicaré la resolución de las operaciones críticas en MongoDB.

2.1. El Motor de Búsqueda Textual (Full-Text Search)

El mayor reto era buscar palabras concretas dentro de millones de párrafos. Un "find" normal con expresión regular recorre la base de datos documento por documento. Esto se llama **Collection Scan**. En Big Data, este proceso falla y tumba el servidor.

Implementaremos **Índices de Texto (\$text)** como solución. MongoDB crea una estructura de **árbol (B-Tree)** interna. Esta mapea cada palabra (**token**) a los documentos que la contienen.

Aplicamos pesos (**weights**). La palabra "5G" tiene más valor en el título principal de la web. Vale menos en la letra pequeña del pie de página.

```
.js BIGDATAUD5_1.js U X
C: > Users > sibko > Documents > GitHub > Cesur-DAM > Asignaturas 2 > HCL (Big Data) > UD5 > Caso_Practico_1 > .js BIGDATAUD5_1.js > ...
1 // 1. Script de creación del índice (Setup inicial)
2 db.web_pages.createIndex([
3   {
4     "contenido.titulo": "text",
5     "contenido.cuerpo_texto": "text",
6     "metadatos_seo.keywords": "text",
7     "dominio": "text"
8   },
9   {
10    weights: {
11      "contenido.titulo": 15,      // Prioridad máxima (si está en el título, puntúa más)
12      "metadatos_seo.keywords": 10, // Prioridad alta
13      "dominio": 5,              // Prioridad media
14      "contenido.cuerpo_texto": 1  // Prioridad normal (texto genérico)
15    },
16    default_language: "spanish" // Vital: aplica stemming (raíces) y elimina stopwords (el, la, los)
17  }
18 ])
```

Un usuario busca "**despliegue 5G rural**". El backend ejecuta esta consulta eficiente:

```

BIGDATAUD5_2.js U X
C: > Users > sibko > Documents > GitHub > Cesur-DAM > Asignaturas 2 > HCL (Big Data) > UD5 > Caso_Practico_1 > BIGDATAUD5_2.js >
1 // 2. Query de recuperación de resultados del usuario
2 db.web_pages.find(
3   { $text: { $search: "despliegue 5G rural" } },
4   { score: { $meta: "textScore" } } // Extrae la puntuación de relevancia que calcula MongoDB
5 )
6   .sort({ score: { $meta: "textScore" } }) // Ordena los resultados de mayor a menor relevancia
7   .limit(20) // Paginación básica: Solo traemos 20 resultados para no saturar la red

```

2.2. Algoritmo de Ranking y el Aggregation Framework

La búsqueda simple de palabras favorece a las webs con spam. Estas repiten el mismo término muchas veces. Debemos calcular la "autoridad" de cada web según los enlaces recibidos. Esto sigue el concepto clásico del PageRank.

MongoDB no es una base de datos orientada a grafos. Pero tiene una herramienta potente para **cálculos masivos**: el Aggregation Framework. Es la evolución de MapReduce.

Diseñamos un proceso por lotes (**Batch Processing**). MongoDB ejecuta un Pipeline de agregación cada madrugada con poco tráfico. El proceso realiza lo siguiente:

\$unwind: Desglosa los arrays de enlaces_salientes.

\$group: Agrupa las webs por dominio y cuenta los enlaces entrantes de sitios de alta autoridad.

\$merge: Actualiza la colección original y sobrescribe el campo autoridad.page_rank_score.

Realizamos los cálculos complejos de noche. Las lecturas de los usuarios son instantáneas durante el día.

2.3. Hiper-Localización (Búsqueda Geoespacial)

Queremos aportar más valor. El usuario busca "restaurante". No tiene sentido mostrarle uno que está a 500 km. **MongoDB** permite **añadir índices geoespaciales (2dsphere)** a los documentos. Conocemos las coordenadas aproximadas del usuario. El backend puede añadir un filtro espacial a la búsqueda:

```
1 db.web_pages.find({
2   "contenido.cuerpo_texto": /restaurante/i,
3   "localizacion_negocio": {
4     $near: {
5       $geometry: { type: "Point", coordinates: [-3.703790, 40.416775] }, // Coordenadas del usuario
6       $maxDistance: 2000 // Busca solo en un radio de 2 kilómetros
7     }
8   }
9 })
```

MongoDB index types

- Single field index
- Compound index
- Multikey index
- Geospatial index
- Test index
- Hashed index



ARQUITECTURA TÉCNICA Y ESCALABILIDAD

FARO funcionará a gran escala. La base de datos manejará **Terabytes o Petabytes** de información. Un solo servidor se saturará aunque tenga mucha RAM (escalabilidad vertical). Debemos usar la **Escalabilidad Horizontal** y mantener la Alta Disponibilidad.

3.1. Sharding (Fragmentación de la Base de Datos)

Usaremos un **clúster de Sharding**. Tomaremos el total de las páginas web guardadas. Las repartiremos físicamente entre varios servidores distintos. La arquitectura de Sharding en MongoDB planteada tiene tres piezas:

Mongos (Routers): Son los servidores de "puerta de entrada". Reciben la búsqueda del usuario. No guardan datos y solo dirigen el tráfico hacia el servidor correcto.

Config Servers: Servidores pequeños que guardan el "mapa" de ubicación de cada fragmento de datos.

Shards (Nodos de Datos): Los servidores donde residen los documentos (ej. Servidor 1, Servidor 2...).

La decisión crítica: La Clave de Fragmentación (**Shard Key**)

Una mala elección de la Shard Key arruina el rendimiento del clúster. Si fragmentamos por "fecha", todas las webs rastreadas hoy irían al mismo servidor físico. Esto lo saturaría y crearía Jumbo Chunks.

Usaremos **Hashed Sharding** sobre el campo dominio. Aplicaremos una función hash criptográfica al nombre del dominio. Los documentos se repartirán de forma uniforme y aleatoria entre todos los servidores. Así equilibramos la carga de trabajo.

3.2. Alta Disponibilidad (Replica Sets) y Teorema CAP

Un buscador no puede "caerse". Por cada Shard físico, montaremos un **Replica Set (un grupo de replicación)**. Este tendrá 3 nodos distribuidos en distintos Centros de Datos:

Nodo Primario: Es el jefe y recibe todas las operaciones de escritura (cuando guardamos webs nuevas).

Nodos Secundarios (x2): Mantienen una copia exacta de los datos del Primario en tiempo real.

El servidor Primario puede quemarse o la red puede caer. Los nodos Secundarios lo notan enseguida mediante **pulsos internos (heartbeats)**. Celebran una "elección" automática en un par de segundos. Uno de ellos se convierte en el nuevo Primario. El sistema se **repara solo** y sin nuestra intervención.

Consideración del **Teorema CAP**:

El Teorema CAP aplica a bases de datos distribuidas. Solo puedes elegir dos opciones entre Consistencia (C), Disponibilidad (A) y Tolerancia a Particiones (P). MongoDB prioriza 'CP' por defecto (Consistencia estricta). Pero para un buscador configuraremos la prioridad en AP (Disponibilidad). Configuraremos las llamadas (Read Preferences) para leer de los nodos secundarios. Preferimos mostrar resultados con unos minutos de desfase (consistencia eventual). Esto es mejor que mostrar una pantalla de error por servicio no disponible.

SEGURIDAD Y GOBERNANZA DEL DATO

Seguimos las normativas de protección de datos (RGPD). La seguridad es obligatoria. La implementación en MongoDB incluirá:

Cifrado en Reposo (Encryption at Rest): Usaremos el motor de almacenamiento WiredTiger de MongoDB. Lo configuraremos con cifrado AES-256 a nivel de disco. Alguien puede robar un disco duro físico del servidor. Los datos serán ilegibles sin la llave. Cifrado en Tránsito: Usaremos certificados TLS/SSL para proteger todas las comunicaciones internas entre los routers mongos, los shards y las aplicaciones cliente.

Control de Acceso (RBAC): Crearemos usuarios de base de datos con roles estrictos. Los bots de rastreo solo tendrán permisos de escritura. La API del buscador mantendrá permisos exclusivos de lectura. Bloquearemos comandos destructivos como dropDatabase. Así, ni siquiera un administrador podrá borrar el clúster por error.



CONCLUSIONES

Diseñé la arquitectura del buscador FARO sobre MongoDB. Este trabajo confirmó que las tecnologías NoSQL transformaron la ingeniería de datos masiva.

Este informe explica cómo la flexibilidad del formato BSON evita los fallos típicos de las bases de datos relacionales ante la variedad de la web. Trazamos una guía técnica clara. Usamos índices \$text junto con búsquedas geoespaciales. También diseñamos pipelines de agregación para calcular la relevancia de los sitios por la noche.

El diseño distribuido usa Sharding y Replica Sets. El sistema soportará la información actual. Podrá escalar horizontalmente en el futuro con más servidores económicos. Esta estructura mantiene la alta disponibilidad. La arquitectura está lista para el despliegue.

REFERENCIAS

<https://www.mongodb.com/es/docs/manual/installation/>

<https://www.mongodb.com/es/docs/manual/text-search/>

<https://studio3t.com/knowledge-base/articles/mongodb-aggregation-framework/>

<https://www.mongodb.com/es/docs/manual/geospatial-queries/>

<https://www.mongodb.com/es/docs/manual/sharding/>

<https://www.mongodb.com/es/docs/manual/core/replica-set-elections/>

MongoDB Curso, Introducción Practica a NoSQL

<https://www.youtube.com/watch?v=lWMemPN9t6Q>

MongoDB in 5 Minutes with Eliot Horowitz

<https://www.youtube.com/watch?v=EE8ZTQxa0AM>

How Google searches one document among Billions of documents quickly?

<https://www.youtube.com/watch?v=CeGtgouT8eA>

Aggregation | MongoDB | Tutorial 10

<https://www.youtube.com/watch?v=Kk6Er0c7srU>