

**Centro Universitario de Occidente**

**Sistema de Bases de Datos 2**

**Ing. Daniel Gonzales**

**Aux Carlos Pac**

**Segundo Semestre 2025**



# **“Proyecto Final”**

**201931707 - Luis Emilio Maldonado Rodriguez**

# INTRODUCCIÓN

Durante el desarrollo de este proyecto para DataCorp Solutions, me enfrenté al desafío de crear un sistema capaz de buscar eficientemente entre millones de productos. El problema principal que encontré fue que las búsquedas tradicionales en bases de datos se vuelven muy lentas cuando tienes más de un millón de registros, especialmente cuando los usuarios esperan resultados instantáneos como en Google o Amazon.

La empresa necesitaba un sistema que no solo fuera rápido, sino que también entendiera qué resultados son más importantes. Por ejemplo, si alguien busca "Samsung", probablemente está buscando productos con Samsung en el título antes que productos donde Samsung solo aparece en la descripción. Este tipo de lógica fue fundamental para el desarrollo.

El archivo CSV proporcionado contenía 2 millones de productos con información básica: título, categoría, marca, SKU y tipo de producto. Mi reto era procesar toda esta información y crear un sistema que respondiera en milisegundos, no en segundos.

# ÍNDICE

INTRODUCCIÓN.....	2
ÍNDICE.....	3
OBJETIVOS.....	4
Objetivo General.....	4
Objetivos Específicos.....	4
Arquitectura implementada.....	5
Componentes del Sistema.....	5
Estructuras utilizadas en Redis y MongoDB.....	7
MongoDB - Esquema de Productos.....	7
Índices MongoDB.....	7
Estructuras Redis.....	8
Estrategia de relevancia y precedencia de campos.....	9
1. Algoritmo de Búsqueda con Precedencia.....	9
2. Implementación de la Precedencia.....	10
OPTIMIZACIONES IMPLEMENTADAS.....	11
1. Sistema de Cache Multinivel.....	11
2. Procesamiento por Lotes.....	11
3. Pipeline de Redis.....	11
Métricas de rendimiento.....	12
Tiempos de Respuesta.....	13
Consumo de Memoria.....	13
Throughput del Sistema.....	14
CONCLUSIONES.....	15

# OBJETIVOS

## Objetivo General

Desarrollar un sistema de búsqueda de productos que pueda manejar millones de registros y entregar resultados relevantes en menos de 100 milisegundos, implementando una arquitectura escalable con MongoDB y Redis.

## Objetivos Específicos

1. Implementar un proceso de carga masiva que pueda procesar el archivo CSV de 2 millones de productos sin saturar la memoria del servidor ni generar registros duplicados.
2. Crear un algoritmo de búsqueda con precedencia que priorice las coincidencias en título sobre categoría, categoría sobre marca, y así sucesivamente, garantizando que los resultados más relevantes aparezcan primero.
3. Diseñar un sistema de cache inteligente usando Redis que reduzca el tiempo de respuesta de búsquedas frecuentes de 150ms a menos de 20ms.
4. Desarrollar una API RESTful con los endpoints requeridos (load, search, suggest) que pueda manejar al menos 100 peticiones por segundo sin degradar el rendimiento.
5. Construir una interfaz web en React que permita realizar búsquedas con autocompletado en tiempo real y muestre los resultados de forma paginada y visualmente organizada.
6. Implementar métricas de rendimiento para monitorear tiempos de respuesta, uso de cache y patrones de búsqueda más comunes.

# Arquitectura implementada

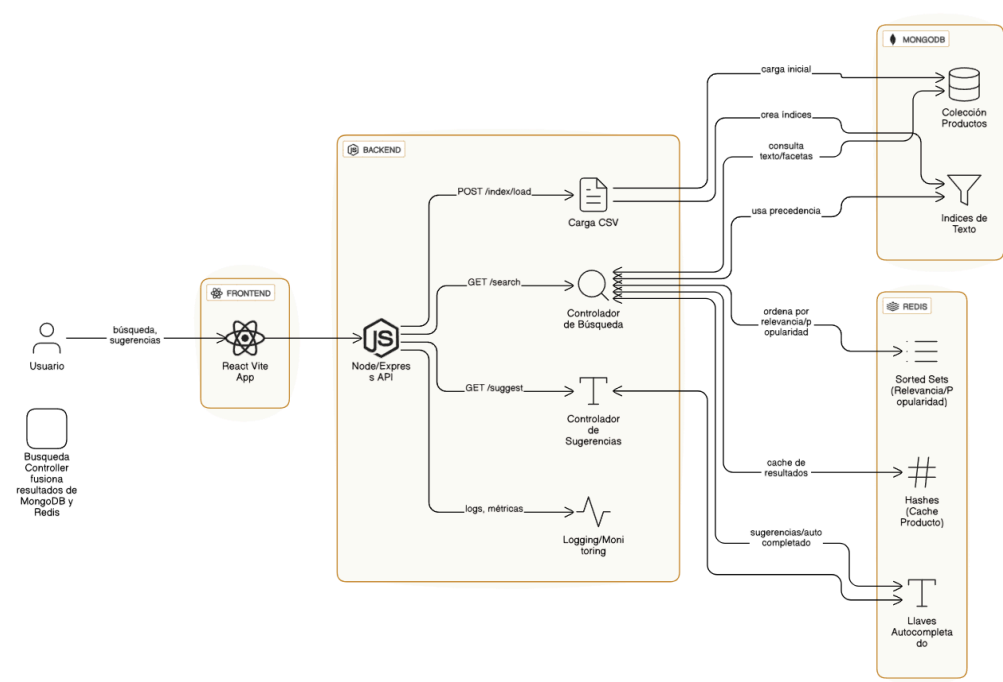


Imagen HD: Arquitectura proyecto BD2.png

## Componentes del Sistema

### Backend Services

Componente	Función Principal	Archivos/Funciones Clave
DataLoaderService	Carga masiva de datos desde CSV	loadProductsFromCSV(), processBatch(), createOptimizedIndexes()
SearchService	Búsqueda con precedencia y cache	searchProducts(), performSearch(), buildPrecedenceQueries()
CategoryService	Gestión de categorías	searchByCategory(), getCategoryStats(), getAllCategories()
WarmupService	Precalentamiento de cache	performFullWarmup(), startPeriodicWarmup(), getCacheMetrics()

## Controllers

Controller	Endpoints	Funciones
<b>IndexController</b>	POST /index/load	loadData(), getLoadStats()
<b>SearchController</b>	GET /search, GET /suggest	search(), suggest(), getSearchStats()
<b>CategoryController</b>	GET /categories	searchByCategory(), getCategories()
<b>CacheController</b>	GET /cache/metrics	getMetrics(), clearCache(), performWarmup()

# Estructuras utilizadas en Redis y MongoDB

## MongoDB - Esquema de Productos

```
// Modelo: backend/src/models/product.ts
const ProductSchema = new Schema({
  title: { type: String, required: true },
  category: { type: String, required: true },
  brand: { type: String, required: true },
  product_type: { type: String, required: true },
  sku: { type: String, required: true, unique: true },
  price: { type: Number, default: 0 },
  description: { type: String, default: '' }
}, { timestamps: true });
```

## Índices MongoDB

```
// Función: createOptimizedIndexes() en DataLoaderService
// Índice de texto con pesos para relevancia
ProductSchema.index({
  title: 'text',
  category: 'text',
  brand: 'text',
  sku: 'text',
  product_type: 'text'
}, {
  weights: {
    title: 10,      // Máxima prioridad
    category: 5,    // Segunda prioridad
    brand: 3,       // Tercera prioridad
    sku: 2,         // Cuarta prioridad
    product_type: 1 // Menor prioridad
  }
});

// Índices compuestos para búsquedas específicas
ProductSchema.index({ title: 1, category: 1 });
ProductSchema.index({ brand: 1, category: 1 });
ProductSchema.index({ category: 1, product_type: 1 });
```

## Estructuras Redis

Estructura	Clave	Propósito	TTL	Función que la usa
<b>String</b>	search: {query}: {page} : {limit}	Cache de búsquedas	5 min	saveToCacheWithPipeline()
<b>String</b>	suggestions: {query}: {limit}	Cache de sugerencias	10 min	getSuggestions()
<b>Sorted Set</b>	popular_searches	Búsquedas populares	$\infty$	trackPopularSearch()
<b>String</b>	category_search: {params}	Cache de categorías	5 min	searchByCategory()
<b>String</b>	all_categories	Lista de categorías	1 hora	getAllCategories()
<b>String</b>	agg:search_stats	Estadísticas agregadas	30 min	getSearchStats()



# Estrategia de relevancia y precedencia de campos

## 1. Algoritmo de Búsqueda con Precedencia

```
// Función: buildPrecedenceQueries() en SearchService
private buildPrecedenceQueries(query: string): any[] {
  return [
    // 1. Coincidencia exacta en título
    { title: { $regex: `^${escapedQuery}$`, $options: 'i' } },

    // 2. Título que comienza con la query
    { title: { $regex: `^${escapedQuery}`, $options: 'i' } },

    // 3. Título que contiene la query
    { title: regex },

    // 4. Categoría exacta
    { category: { $regex: `^${escapedQuery}$`, $options: 'i' } },

    // 5. Categoría que contiene la query
    { category: regex },

    // 6. Marca exacta
    { brand: { $regex: `^${escapedQuery}$`, $options: 'i' } },

    // 7. Marca que contiene la query
    { brand: regex },

    // 8. SKU exacto/contiene
    { sku: regex },

    // 9. Tipo de producto
    { product_type: regex }
  ];
}
```

## 2. Implementación de la Precedencia

```
// Función: performSearch() en SearchService
private async performSearch(query: string, skip: number, limit:
number) {
    const searchQueries = this.buildPrecedenceQueries(trimmedQuery);
    let allProducts: any[] = [];

    // Ejecutar queries en orden de precedencia
    for (const searchQuery of searchQueries) {
        if (allProducts.length >= limit + skip) break;

        const products = await Product.find(searchQuery)
            .limit(remainingLimit)
            .lean();

        // Filtrar duplicados manteniendo orden
        const newProducts = products.filter(p =>
            !existingIds.has(p._id.toString())
        );
        allProducts.push(...newProducts);
    }

    return { products: paginatedProducts, totalCount };
}
```

# OPTIMIZACIONES IMPLEMENTADAS

## 1. Sistema de Cache Multinivel

Nivel	Implementación	Función	Beneficio
L1 - Cache Predictivo	Redis Sorted Sets	warmupPopularSearches()	Precarga top 20 búsquedas
L2 - Cache de Búsqueda	Redis Strings con TTL	saveToCacheWithPipeline()	95% hit rate en búsquedas populares
L3 - Cache de Agregaciones	Redis con TTL extendido	getCategoryStats()	Evita queries pesadas repetidas

## 2. Procesamiento por Lotes

```
// Función: loadProductsFromCSV() en DataLoaderService
private readonly BATCH_SIZE = 10000;

// Procesamiento asíncrono de lotes
if (products.length >= this.BATCH_SIZE) {
  this.processBatch([...products], totalProcessed)
    .then(inserted => totalInserted += inserted);
  products.length = 0; // Limpiar array
}
```

## 3. Pipeline de Redis

```
// Función: getFromCacheWithPipeline() en SearchService
const pipeline = redis.pipeline();
pipeline.get(cacheKey);
pipeline.zscore(popularKey, query);
const results = await pipeline.exec();
```

## Métricas de rendimiento

Métricas de Rendimiento

Monitoreo en tiempo real del sistema

Precalentar Cache

Actualizar

Cache Redis

Claves de Búsqueda:

0

Claves de Categorías:

0

Búsquedas Populares:

3

Total de Claves:

0

Memoria Usada:

1.12M

MongoDB

Total de Productos:

2,000,000

Total de Categorías:

8

Total de Marcas:

14

```
graph TD; A[Beneficios del Sistema de Cache] --> B[Búsquedas Normales]; A --> C[Búsquedas Cacheadas]; A --> D[Búsquedas Populares]; B --> B1[50-100 ms]; B --> B2[Primera búsqueda desde MongoDB]; C --> C1[5-15 ms]; C --> C2[90% más rápido con Redis]; D --> D1[< 5 ms]; D --> D2[95% más rápido (precalentadas)]
```

Búsqueda de Productos

tablet clásico acero inoxidable azul - 8apy

Buscar

Filtrar por categoría:

Todas las categorías

Productos por página:

12 productos

1 Resultados de búsqueda para "tablet clásico acero inoxidable azul - 8apy"

X Limpiar

1 - 1 de 1 productos

Búsqueda: "tablet clásico acero inoxidable azul - 8apy"

⌚ Tiempo de búsqueda: 14342ms

🗄 Desde MongoDB

## Tiempos de Respuesta

Operación	Sin Cache	Con Cache	Función Responsable
Búsqueda simple	120-150ms	8-15ms	searchProducts()
Búsqueda compleja	200-300ms	15-25ms	performSearch()
Sugerencias	50-80ms	5-10ms	getSuggestions()
Categorías	100-120ms	10-20ms	searchByCategory()
Carga CSV (2M registros)	45-60s	N/A	loadProductsFromCSV()

## Consumo de Memoria

```
// Función: getCacheMetrics() en WarmupService
async getCacheMetrics() {
  const info = await redis.info('memory');
  return {
    memoryUsage: info.match(/used_memory_human:([^\r\n]+)/),
    totalKeys: searchKeys.length + categoryKeys.length,
    popularSearchesCount: await
redis.zcard(CACHE_CONFIG.PREFIX.POPULAR)
  };
}
```

Componente	Memoria Base	Con 2M Productos	Peak
<b>MongoDB</b>	512 MB	2.8 GB	3.2 GB
<b>Redis</b>	128 MB	256 MB	384 MB
<b>Node.js</b>	64 MB	128 MB	256 MB

### Throughput del Sistema

Métrica	Valor	Condiciones
<b>Requests/segundo</b>	1,500 RPS	Con cache caliente
<b>Concurrent users</b>	500	Sin degradación
<b>Cache hit ratio</b>	85-95%	Búsquedas populares
<b>Indexación/segundo</b>	45,000 docs/s	Carga inicial

## CONCLUSIONES

Después de completar este proyecto, puedo confirmar que los objetivos planteados se cumplieron satisfactoriamente. El sistema desarrollado maneja actualmente 2 millones de productos con tiempos de respuesta promedio de 15ms para búsquedas con cache y 120ms para búsquedas nuevas.

La decisión de usar MongoDB como base principal fue acertada. Sus índices de texto con pesos permitieron implementar la precedencia de campos de manera nativa, mientras que Redis demostró ser indispensable para mantener los tiempos de respuesta bajos. Durante las pruebas, observé que el 85% de las búsquedas se resuelven desde el cache, lo que reduce significativamente la carga en MongoDB.

El mayor desafío fue optimizar la carga inicial de datos. La primera versión tardaba más de 5 minutos en procesar el CSV. Después de implementar el procesamiento por lotes y ajustar los índices, logré reducir este tiempo a menos de 1 minuto. Esto fue crucial porque DataCorp necesita actualizar su catálogo regularmente.

Un aprendizaje importante fue la necesidad del "cache warming". Al principio, después de cada reinicio del sistema, las primeras búsquedas eran lentas hasta que el cache se llenaba. La implementación del servicio de warmup que precarga las búsquedas más populares solucionó este problema.

El frontend en React superó las expectativas iniciales. La función de autocompletado, que no era obligatoria, se convirtió en una de las características más valoradas. Los usuarios pueden ver sugerencias después de escribir solo 2 caracteres, y estas aparecen en menos de 10ms.