

# Modelagem e Armazenamento de Mensagens de uma Rede Social Banco de Dados Cassandra e DataStax Astra

Guilherme César Athayde

Novembro de 2025

## 1 Introdução

Este relatório descreve a implementação de um sistema de armazenamento de mensagens de uma rede social para a comunidade universitária da UFSCar, utilizando o DataStax Astra DB (Cassandra as a Service) e a biblioteca `astrappy` em Python.

O sistema foi projetado para atender aos requisitos de uma rede social moderna, com foco em escalabilidade, performance e flexibilidade na modelagem de dados.

## 2 Scripts CQL e Implementação

### 2.1 Tarefa 1: Criar a Coleção (Tabela)

No DataStax Astra DB, utilizamos a abordagem de *collections* que é equivalente a tabelas no Cassandra tradicional, mas com maior flexibilidade por ser schemaless. A estrutura dos documentos foi definida da seguinte forma:

```
1 collection_name = "social_messages"
2
3 # Dropar colecao se ja existir (para testes)
4 db.drop_collection(collection_name)
5
6 # Criar a colecao com estrutura:
7 #
8 #   "_id": "message_id",
9 #   "user_id": "user_xxx",
10 #   "user_age": 25,
11 #   "topic": "tecnologia",
12 #   "message_text": "Texto da mensagem",
13 #   "timestamp": "2025-11-20T10:30:00"
14 #
15
16 collection = db.create_collection(collection_name)
```

Listing 1: Criação da coleção

### Equivalente CQL Cassandra:

```
1 CREATE TABLE IF NOT EXISTS social_messages (
2     message_id TEXT,
3     user_id TEXT,
4     user_age INT,
5     topic TEXT,
6     message_text TEXT,
7     timestamp TIMESTAMP,
8     PRIMARY KEY (user_id, timestamp, message_id)
9 ) WITH CLUSTERING ORDER BY (timestamp DESC);
10
11 CREATE INDEX IF NOT EXISTS idx_topic
12 ON social_messages (topic);
```

Listing 2: CQL equivalente para Cassandra tradicional

## 2.2 Tarefa 2: Inserir Dados (20+ mensagens)

Foram inseridas 25 mensagens com distribuição variada entre 10 usuários diferentes e 6 tópicos distintos:

```
1 users = [f"user_{str(i).zfill(3)}"
2         for i in range(1, 11)] # user_001 a user_010
3 user_ages = {user: random.randint(18, 65)
4             for user in users}
5 topics = ["politica", "saude", "tecnologia",
6            "educacao", "esportes", "entretenimento"]
7
8 # Gerar 25 mensagens com diferentes usuarios,
9 # tópicos e datas
10 messages = []
11 base_date = datetime.now() - timedelta(days=30)
12
13 for i in range(25):
14     user = random.choice(users)
15     topic = random.choice(topics)
16     timestamp = base_date + timedelta(
17         days=random.randint(0, 30),
18         hours=random.randint(0, 23),
19         minutes=random.randint(0, 59))
20
21     message = {
22         "_id": f"msg_{str(i+1).zfill(3)}",
23         "user_id": user,
24         "user_age": user_ages[user],
25         "topic": topic,
26         "message_text": f"{message_text}"
27                         (Tema: {topic})",
28         "timestamp": timestamp.isoformat()
29     }
30     messages.append(message)
```

```

31
32 result = collection.insert_many(messages)

```

Listing 3: Inserção de mensagens

**Equivalente CQL:**

```

1 INSERT INTO social_messages
2   (message_id, user_id, user_age, topic,
3    message_text, timestamp)
4 VALUES
5   ('msg_001', 'user_001', 25, 'tecnologia',
6    'Mensagem sobre tecnologia',
7    '2025-11-01T10:30:00');
8 -- Repetir para todas as 25 mensagens...

```

Listing 4: CQL para inserção de dados

### 2.3 Tarefa 3: Recuperar Mensagem Específica

Implementação da consulta para recuperar mensagens de um usuário específico:

```

1 # Buscar uma mensagem do user_001
2 target_user = "user_001"
3 user_message = collection.find_one(
4     {"user_id": target_user}
5 )
6
7 # Buscar todas as mensagens de um usuario especifico
8 all_user_messages = list(collection.find(
9     {"user_id": target_user}
10 ))

```

Listing 5: Consulta de mensagens por usuário

**Equivalente CQL:**

```

1 -- Buscar uma mensagem especifica
2 SELECT * FROM social_messages
3 WHERE user_id = 'user_001'
4 LIMIT 1;
5
6 -- Buscar todas as mensagens de um usuario
7 SELECT * FROM social_messages
8 WHERE user_id = 'user_001';

```

Listing 6: CQL para consulta por usuário

### 2.4 Tarefa 4: Índice e Frequência de Tópicos

Para consultar a frequência de tipos de mensagens por usuário:

```

1 # Recuperar todas as mensagens
2 all_messages = list(collection.find({}))

```

```

3
4 # Agrupar por usuario e contar frequencia de topics
5 user_topic_frequency = {}
6
7 for msg in all_messages:
8     user = msg['user_id']
9     topic = msg['topic']
10
11    if user not in user_topic_frequency:
12        user_topic_frequency[user] = {}
13
14    if topic not in user_topic_frequency[user]:
15        user_topic_frequency[user][topic] = 0
16
17    user_topic_frequency[user][topic] += 1
18
19 # Exibir resultados
20 for user, topics in sorted(
21     user_topic_frequency.items()
22):
23     print(f"\n{user}:")
24     for topic, count in sorted(
25         topics.items(),
26         key=lambda x: x[1],
27         reverse=True
28     ):
29         print(f"    - {topic}: {count} mensagem(ns)")

```

Listing 7: Consulta de frequência de tópicos por usuário

#### Equivalente CQL (agregação):

```

1 -- No Cassandra, a agregacao deve ser feita
2 -- em nivel de aplicacao
3 -- ou usando Spark/DataStax Analytics
4
5 SELECT user_id, topic, COUNT(*) as frequency
6 FROM social_messages
7 GROUP BY user_id, topic;
8
9 -- Alternativa com indice:
10 CREATE MATERIALIZED VIEW messages_by_user_topic AS
11     SELECT user_id, topic, message_id, timestamp
12     FROM social_messages
13     WHERE user_id IS NOT NULL
14         AND topic IS NOT NULL
15     PRIMARY KEY ((user_id, topic), timestamp,
16                 message_id)
17     WITH CLUSTERING ORDER BY (timestamp DESC);

```

Listing 8: CQL para agregação de frequências

### 3 Discussão sobre a Modelagem

#### 3.1 Escolhas de Modelagem

A modelagem escolhida seguiu os princípios do Cassandra/DataStax Astra:

1. **Desnormalização:** Ao contrário de bancos relacionais, optamos por armazenar todos os dados da mensagem em um único documento, incluindo idade do usuário, que poderia estar em uma tabela separada em um modelo relacional.
2. **Query-First Design:** A estrutura foi projetada considerando as consultas mais frequentes:
  - Busca por usuário
  - Busca por data/timestamp
  - Busca por tópico
3. **Primary Key:** Em uma implementação Cassandra tradicional, a chave primária seria (`user_id`, `timestamp`, `message_id`), onde:
  - `user_id`: Partition key (distribui dados entre nós)
  - `timestamp`: Clustering key (ordena dados dentro da partição)
  - `message_id`: Garante unicidade
4. **Schemaless com Astra:** O uso de collections no Astra DB oferece flexibilidade schemaless, permitindo evolução do schema sem migração de dados.

#### 3.2 Estratégia de Indexação

##### 3.2.1 Índices Automáticos do Astra DB

O DataStax Astra DB gerencia índices automaticamente para otimizar consultas. Para campos frequentemente consultados como `user_id` e `topic`, o Astra cria índices secundários dinamicamente.

##### 3.2.2 Índices Secundários (Cassandra Tradicional)

Em um ambiente Cassandra tradicional, implementaríamos:

- **Índice em topic:** Para permitir consultas por tipo de mensagem
- **Materialized View:** Para agregações pré-computadas de frequência por usuário e tópico

##### 3.2.3 Trade-offs

Vantagens:

- Consultas por usuário são extremamente rápidas (partition key lookup)
- Mensagens ordenadas por timestamp dentro de cada partição

- Escalabilidade horizontal natural do Cassandra
- Índices secundários permitem consultas flexíveis por tópico

#### **Desvantagens:**

- Índices secundários podem impactar performance de escrita
- Agregações (contagem de frequências) são mais eficientes quando processadas na aplicação
- Consultas cross-partition (todos os usuários) são menos eficientes

### **3.3 Consultas e Análise de Resultados**

#### **3.3.1 Performance de Consultas**

1. **Consulta por usuário específico:** Operação O(1) usando partition key, extremamente eficiente mesmo com milhões de mensagens.
2. **Consulta de frequências:** Requer full table scan, mas com processamento eficiente em memória. Para produção, recomenda-se:
  - Usar Materialized Views para pré-agregação
  - Implementar cache (Redis) para resultados frequentes
  - Usar Apache Spark para análises complexas
3. **Consulta por tópico:** Índice secundário permite busca eficiente, mas menos otimizada que partition key lookup.

#### **3.3.2 Resultados Obtidos**

A implementação demonstrou:

- Inserção bem-sucedida de 25 mensagens distribuídas entre 10 usuários
- Recuperação instantânea de mensagens por usuário específico
- Cálculo correto de frequências de tópicos por usuário
- Distribuição equilibrada entre os 6 tópicos disponíveis

## **4 Conclusões**

A implementação atendeu todos os requisitos funcionais propostos:

1. ✓ Cada mensagem possui ID único
2. ✓ Armazenamento de user\_id e idade
3. ✓ Classificação por tema/tópico
4. ✓ Registro de timestamp
5. ✓ Consultas rápidas por usuário implementadas