

como o  armazena

# TRILHÓES

de mensagens



# O DISCORD

O Discord é uma plataforma de comunicação largamente utilizada nos dias de hoje. Foi lançado em maio de 2015 e rapidamente conquistou o gosto do público. Em 2016, a plataforma, extraordinariamente, já registrava uma base ativa mensal de 10 milhões de usuários.

O DISCORD USA AS  
MELHORES TECNOLOGIAS DO  
MERCADO, CERTO?

# ERRADO!

# ERRADO!



2015

ITERAÇÃO RÁPIDA

POUCA DOR DE CABEÇA

MIGRAÇÃO RÁPIDA NO FUTURO



# mongoDB

2015-2016

**ESTRUTURA FLEXÍVEL**



**DESENVOLVIMENTO RÁPIDO**

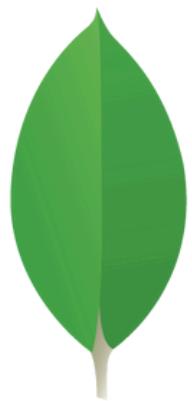
**SEM SHARDING**

**+**

**REPLICA SET UNICO**



**SIMPLICIDADE**



# mongoDB

2015-2016

**ESTRUTURA FLEXÍVEL**



**DESENVOLVIMENTO RÁPIDO**

**SEM SHARDING**

**+**

**REPLICA SET UNICO**

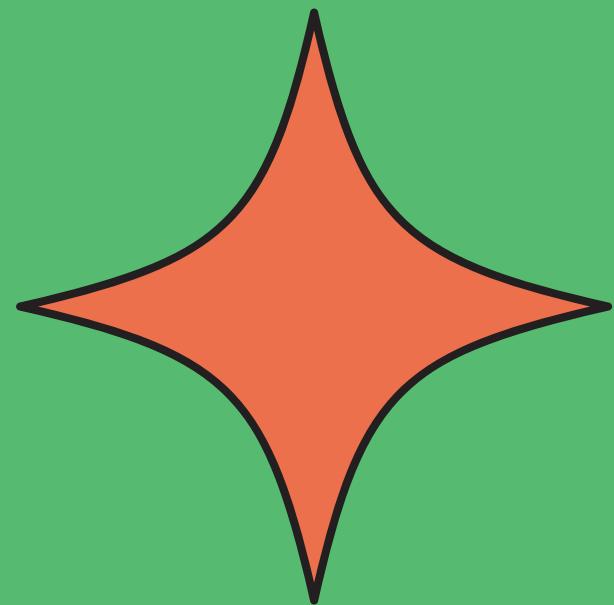


**SIMPLICIDADE**



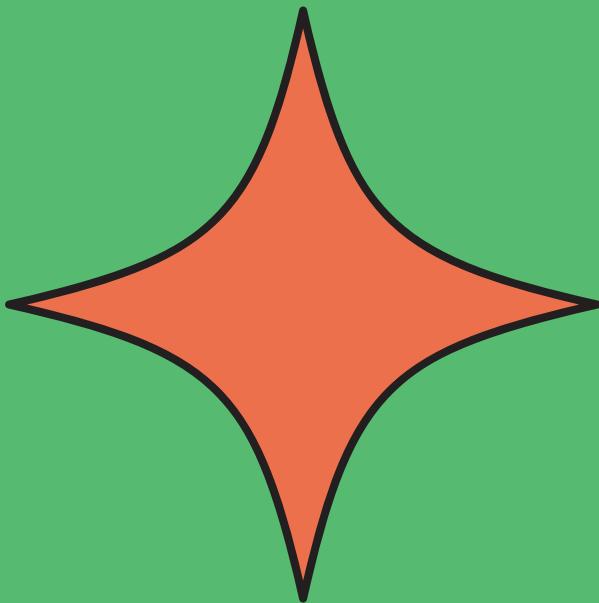
**MIGRAÇÃO FÁCIL NO FUTURO**

# FUNCIONOU!



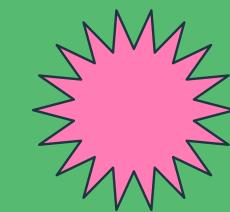
# FUNCIONOU!

*para o que elas precisavam*

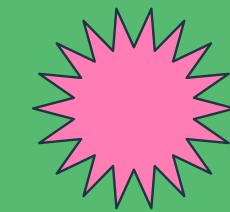


**ATÉ 100 MILHÓES DE MENSAGENS  
CONSEGUIRAM UM SISTEMA BEM  
ROBUSTO SEM ESCALA HORIZONTAL E  
UM ÚNICO ÍNDICE**

# COMO AS MENSAGENS ERAM ARMAZENADAS?



Uma conversa pertence à um canal e elas possuem uma ordem de tempo nas mensagens enviadas.



índice composto nos atributos channel\_id e created\_at

The screenshot shows a Discord direct message window titled '@Bacon - Discord'. The left sidebar lists various servers and users, including 'Amigos' (marked with a red '1'), 'Nitro', 'Loja', and several channels like 'Primata09', 'Delkrusty', and 'Discord OFFICIAL'. The main conversation is between two users: Bacon and skineura. The messages are timestamped by date:

- Bacon (18/03/2025, 00:23): musica do diabo
- skineura (18/03/2025, 14:51): hj to meio rock wins
- Bacon (23/03/2025, 14:33): agora voce ta mandando muito
- skineura (23/03/2025, 15:25): hj to meio sertanejo wins

Each message includes a small Spotify-style preview card showing album art and song titles. On the right side of the window, there is a profile sidebar for Bacon, which includes a large profile picture, the name 'Bacon', the status 'itsdibas\_ ★ #', and sections for 'Membro desde' (10 de mai. de 2016) and 'Servidores em comum' (6). Below the sidebar is a button labeled 'Ver Perfil Completo'.

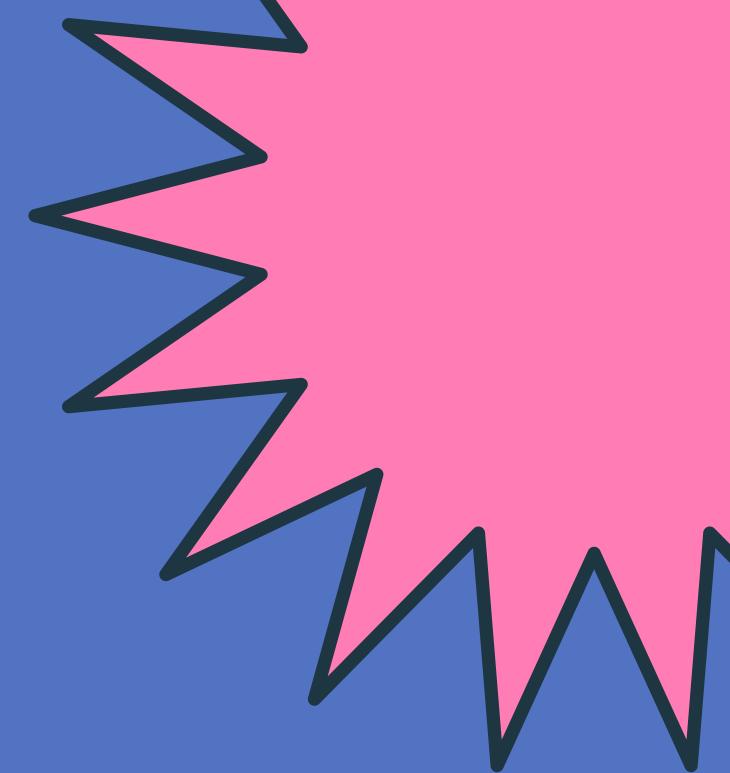
# O COMEÇO DO FIM

AS MENSAGENS  
ARMAZENADAS  
PASSARAM DE  
100 MILHÓES

DADOS NÃO CABEM NA  
MEMÓRIA

+ SEEKS NO DISCO

+ LATÊNCIA



# MIGRAÇÃO PARA O CASSANDRA (2016/2017)

## **REQUISITOS DEFINIDOS**

**ESCALABILIDADE LINEAR**

**PERFORMANCE PREVISIVEL**

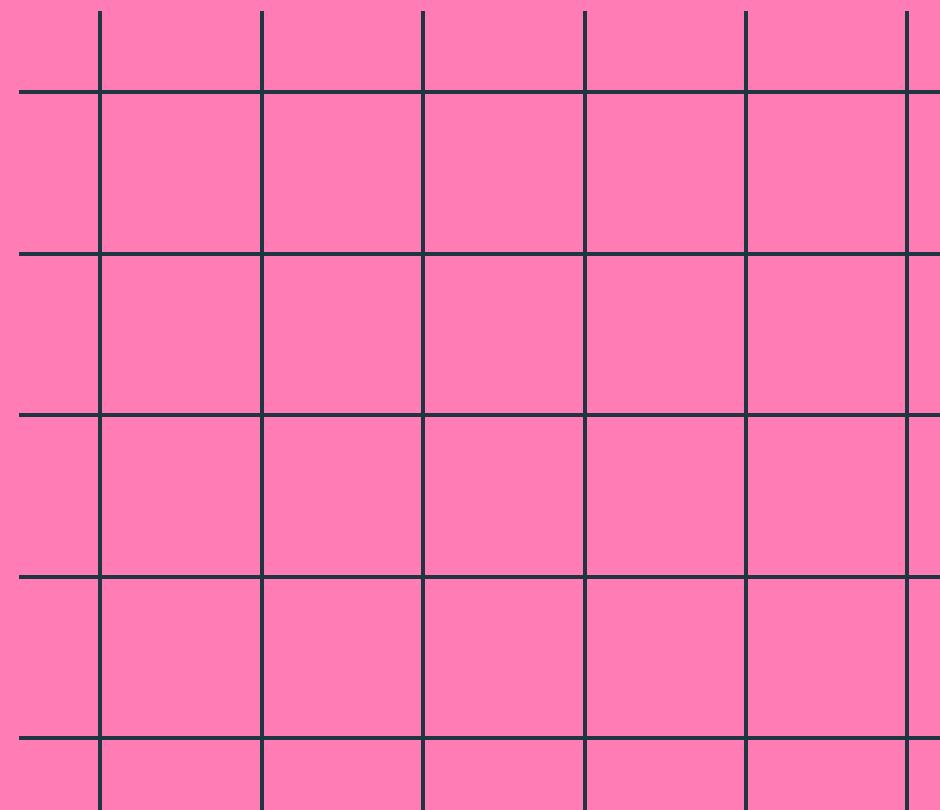
**RECUPERAÇÃO AUTOMÁTICA**

**NÃO É UMA BLOB-STORE**

**BAIXA MANUTENÇÃO**

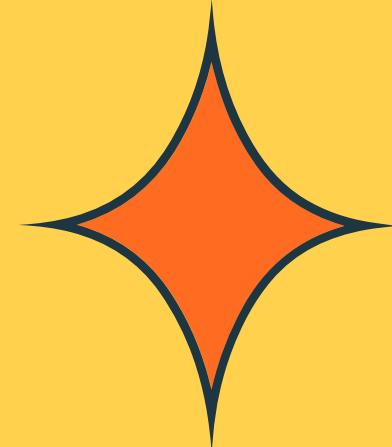
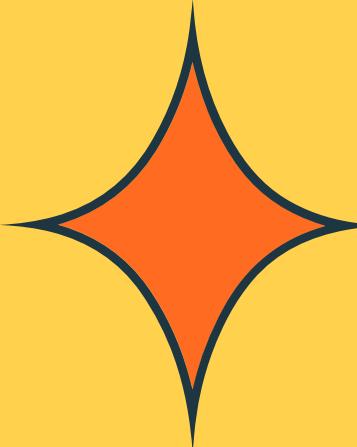
**OPEN-SOURCE**

**TECNOLOGIA CONSOLIDADA**





**CASSANDRA**



# KKV-STORE

## KEY KEY VALUE



Chave de partição: determina o nó e localização física do dado no servidor



Chave de agrupamento: identificador único da tupla e atributo de ordenamento

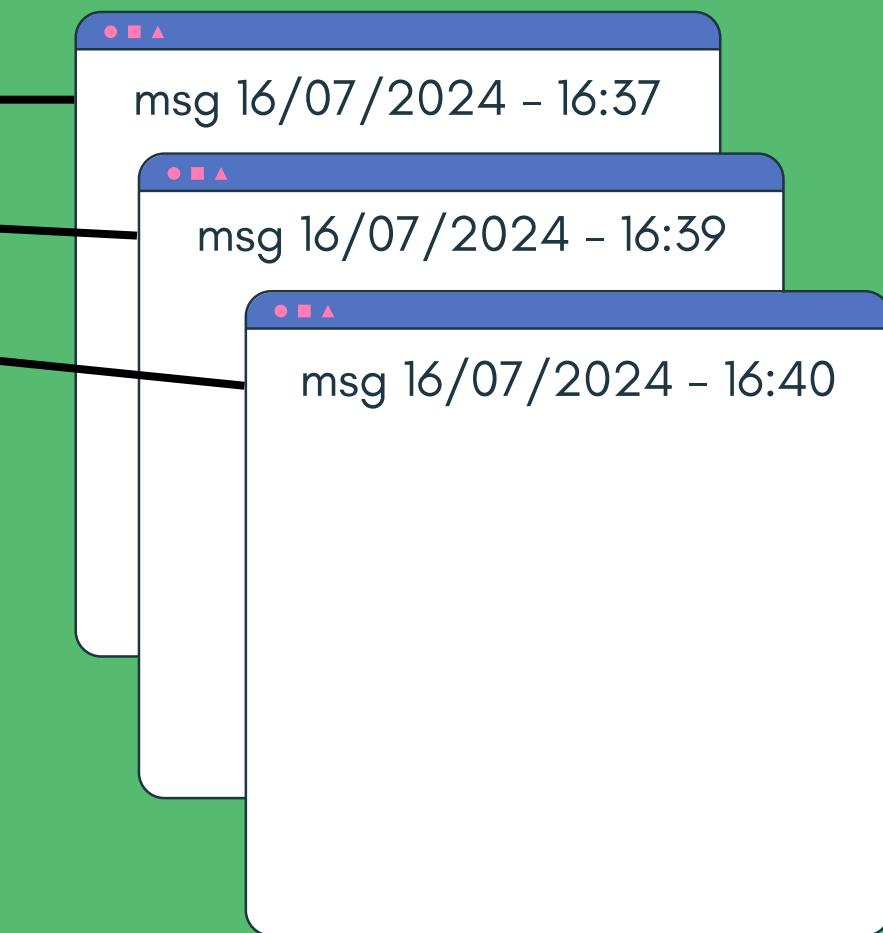


# KKV-STORE

## KEY KEY VALUE



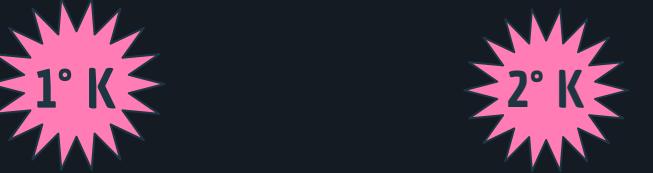
1° K



2° K

# NO DISCORD

```
CREATE TABLE messages (
    channel_id bigint,
    message_id bigint,
    author_id bigint,
    content text,
    PRIMARY KEY (channel_id, message_id)
) WITH CLUSTERING ORDER BY (message_id DESC);
```



PROBLEMA COM  
ESSA  
MODELAGEM

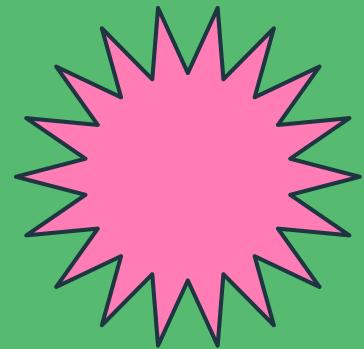
PARTIÇÕES  
GRANDES DEMAIS

# SOLUÇÃO: BUCKETING

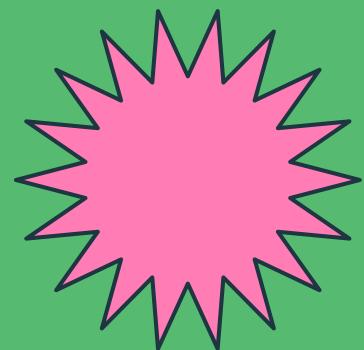
```
CREATE TABLE messages (
    channel_id bigint,
    bucket int,
    message_id bigint,
    author_id bigint,
    content text,
    PRIMARY KEY ((channel_id, bucket), message_id)
) WITH CLUSTERING ORDER BY (message_id DESC);
```



# SUBINDO EM PROD



Modificaram o código para fazer os reads/writes tanto do Mongo quanto do Cassandra



author\_id == null ? Como?

# CONSISTÊNCIA EVENTUAL

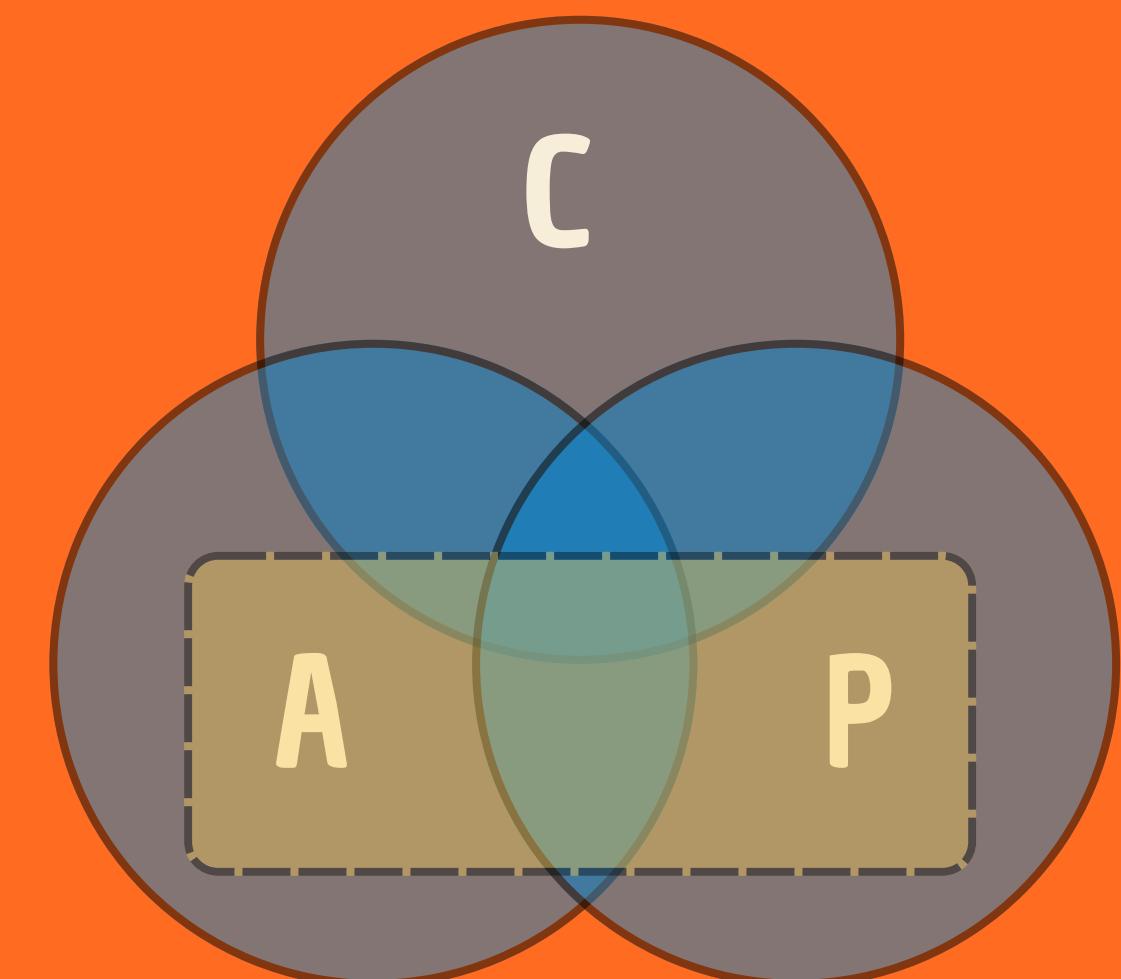
ALTA DISPONIBILIDADE

RESISTÊNCIA A PARTIÇÃO

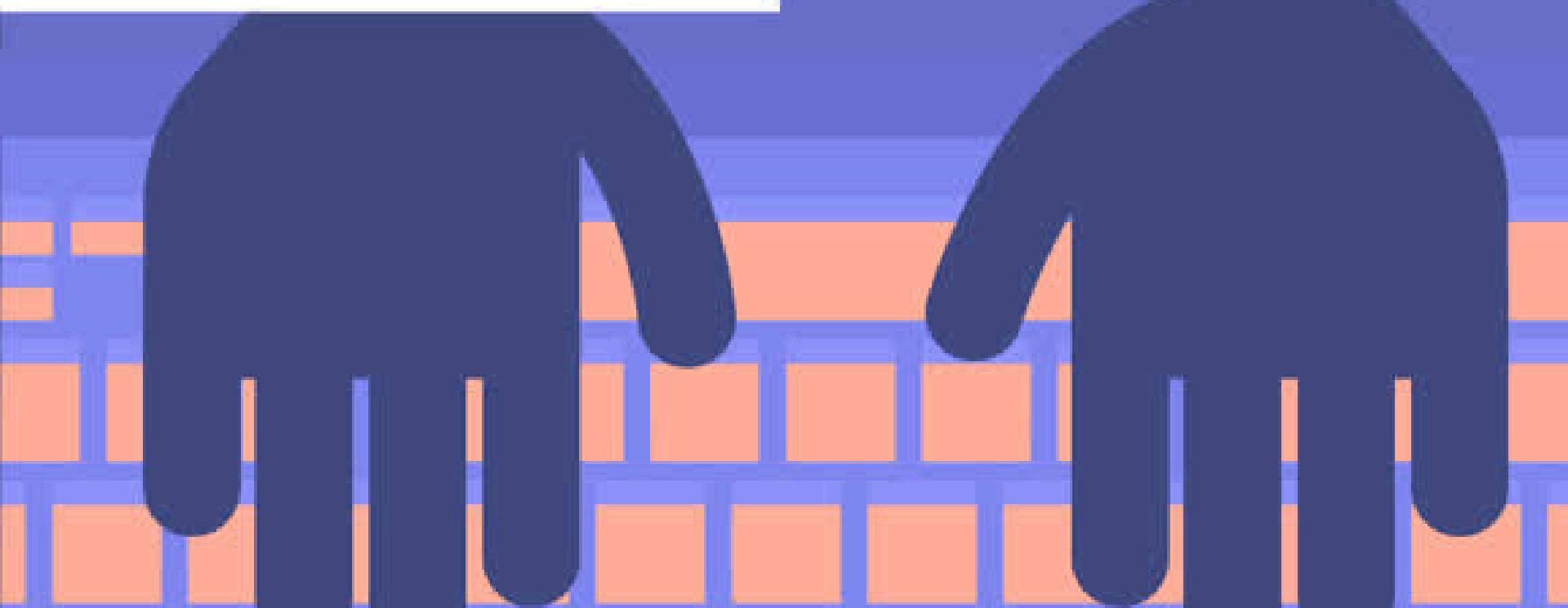
CASSANDRA É AP

UPSERT

LAST-WRITE-WINS P/ COLUNA

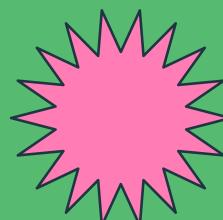


I STAN HE

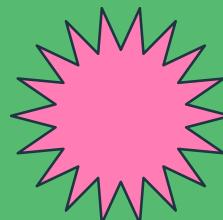


# SOLUÇÃO: MENSAGEM COM AUTHOR\_ID == NULL → CORROMPIDA

## WRITES INEFICIENTES



No Cassandra, null → tombstone



De 16 campos da mensagem, estavam escrevendo 12 campos nulls.

## TOMBSTONES: DADOS DELETADOS SÃO MARCADOS COMO TOMBSTONES AFIM DE PRESERVAR A CONSISTÊNCIA EVENTUAL.

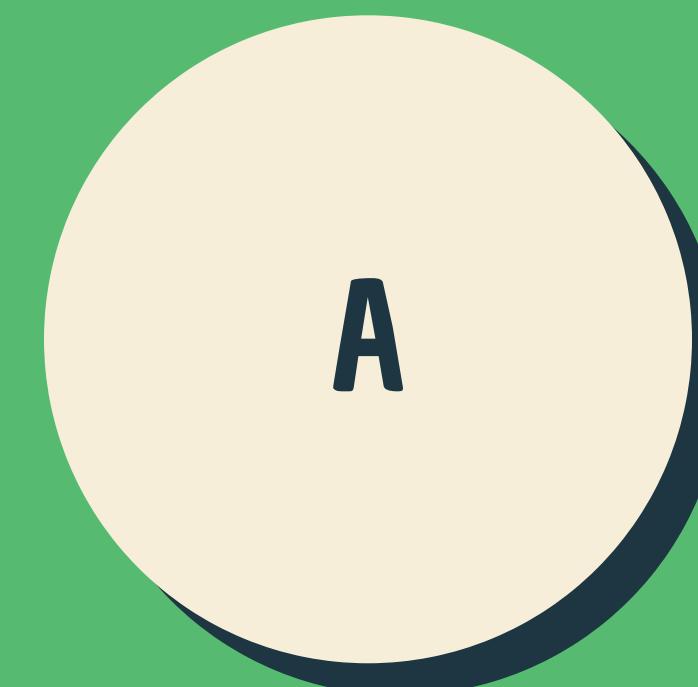
# TOMBSTONES:

DADOS DELETADOS SÃO MARCADOS COMO  
TOMBSTONES AFIM DE PRESERVAR A  
CONSISTÊNCIA EVENTUAL.

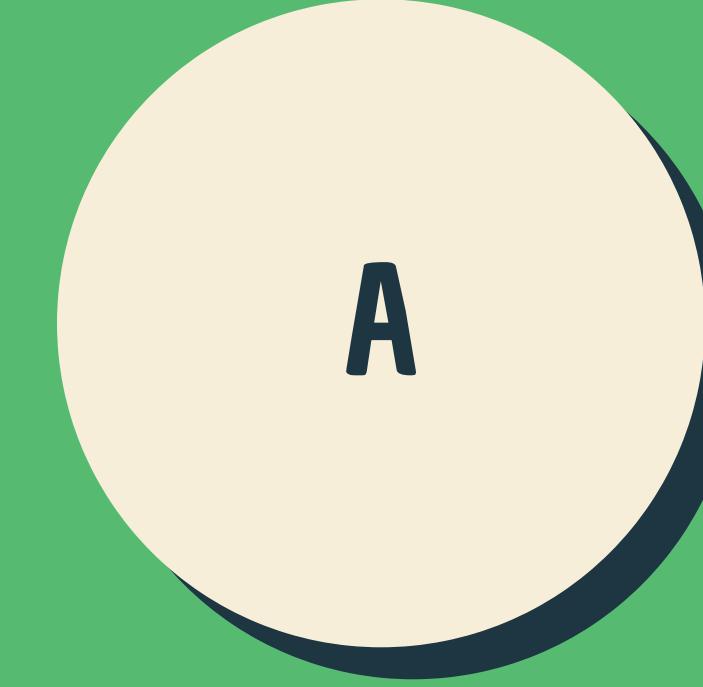
NÓ-1



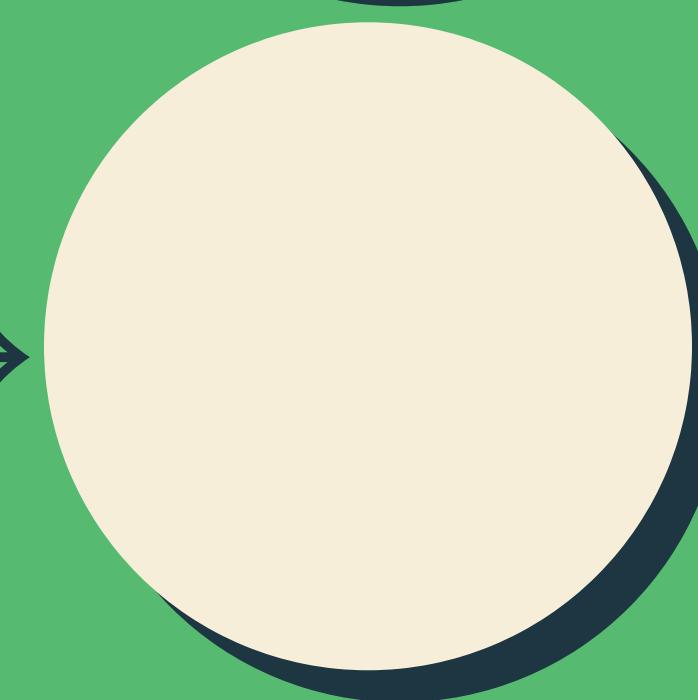
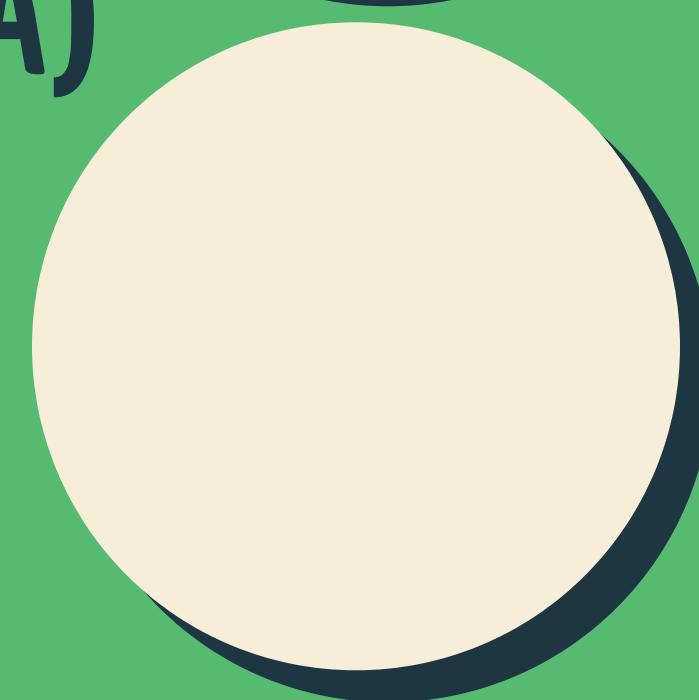
NÓ-2



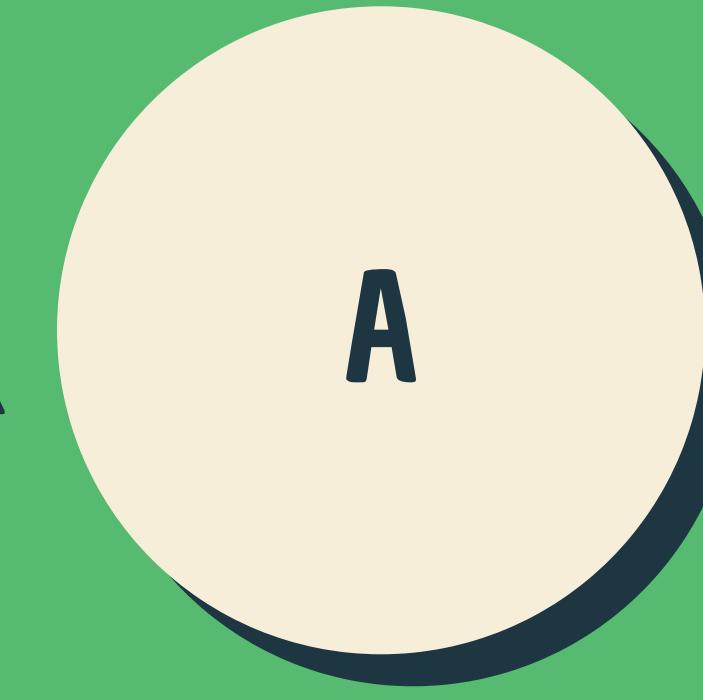
NÓ-3



DELETE(A)



X



UPDATE(A)



# TOMBSTONES:

DADOS DELETADOS SÃO MARCADOS COMO  
TOMBSTONES AFIM DE PRESERVAR A  
CONSISTÊNCIA EVENTUAL.

NÓ-1

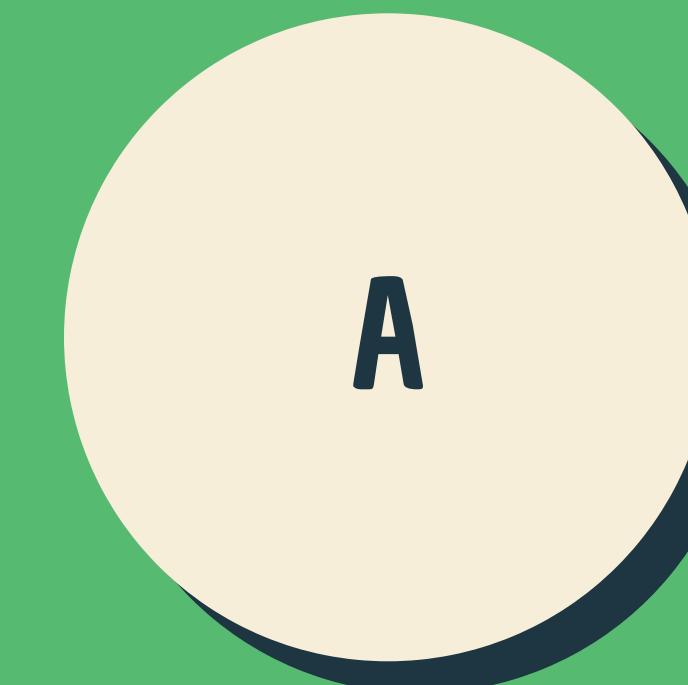


DELETE(A)



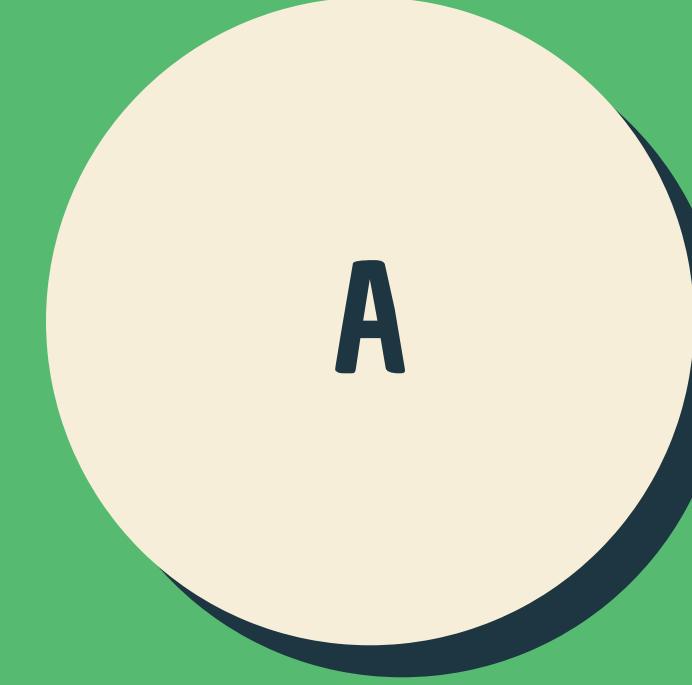
TOMBSTONE(A)

NÓ-2



TOMBSTONE(A)

NÓ-3



X

UPDATE(A)

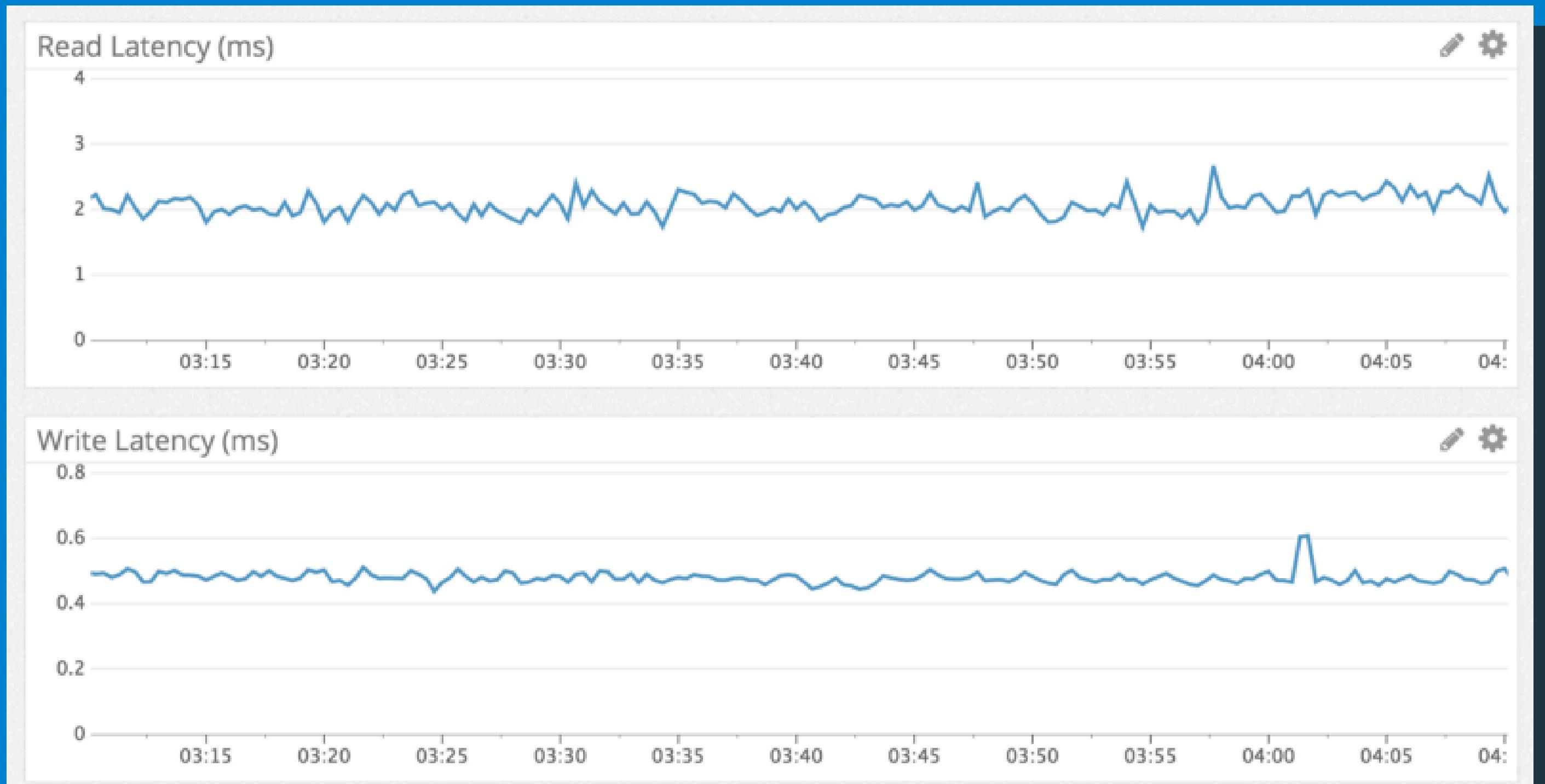


# DISCORD X TOMBSTONE:

NO CASO ANTERIOR, O DELETE/UPDATE ACONTECE EM NÓS DIFERENTES.

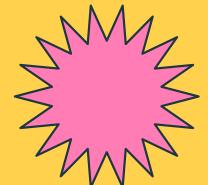
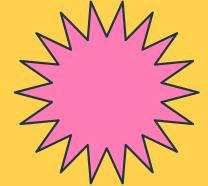
NÓ CASO DO DISCORD, ACONTECEM NO MESMO NÓ, RESULTANDO EM UMA TUPLA ZUMBI: UM DADO QUE RESSUCITOU DE FORMA INCONSISTENTE, ERA PARA SER UM TOMBSTONE, MAS FOI RESSUCITADO POR CONTA DO UPSERT.

# RESULTADOS

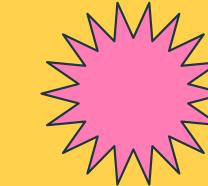
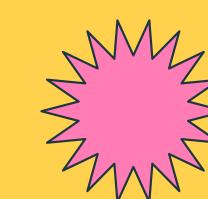


# A GRANDE SURPRESA

EM UM BELO DIA, O CASSANDRA PAROU DE FUNCIONAR DIREITO.

-  Pausas de até 10 segundos para Garbage Collection
-  Servidor público do subreddit “Puzzles & Dragons” teve todas suas milhões mensagens, exceto uma, apagadas via API do Discord
-  Resultado: milhões de tombstones que o Cassandra precisava checar

# SOLUÇÕES

-  Redução do tempo de vida dos tombstones
-  Otimização da query

# MIGRAÇÃO PARA O SCYLLADB (2023)

**AO ATINGIR A MARCA DE TRILHÓES DE  
MENSAGENS, O CASSANDRA PASSOU A  
APRESENTAR LIMITAÇÕES**

## LIMITAÇÕES

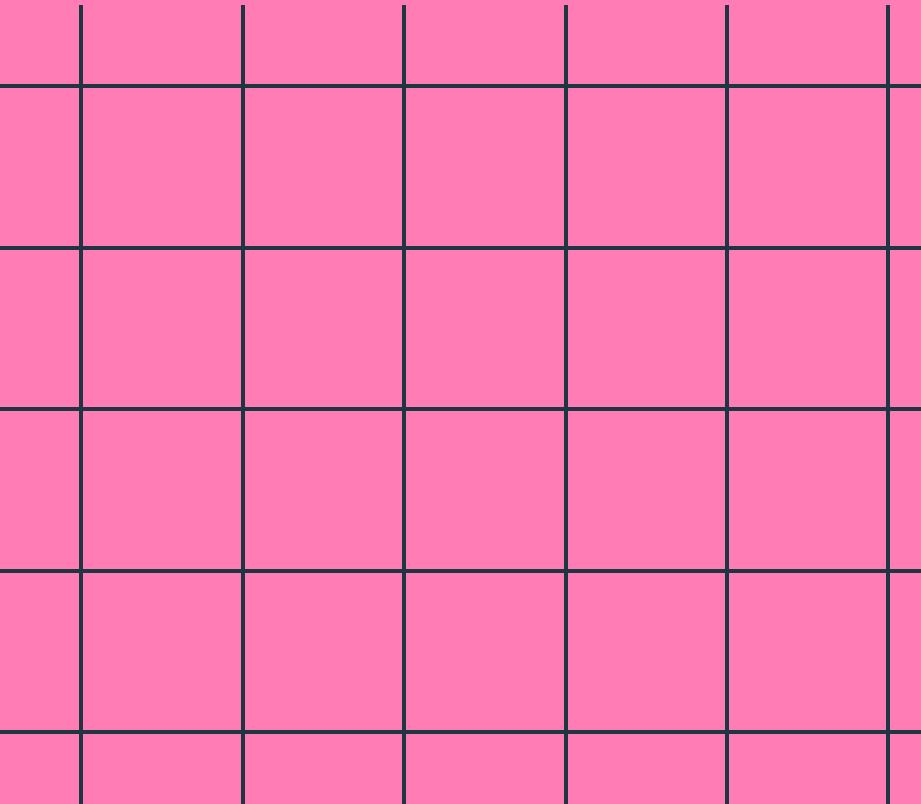
LATÊNCIA IMPREVISTAS

PAUSAS DE GC

HOT PARTITIONS

MANUTENÇÃO CUSTOSA

COMPACTAÇÃO INEFICIENTE



# SCYLLADB

Compatível com Cassandra

Escrito em C++

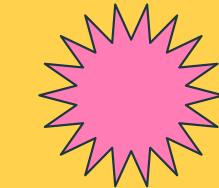
Arquitetura shard-per-core

## ARQUITETURA SHARD-PER-CORE

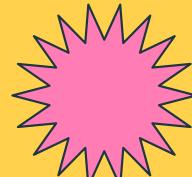
Cada núcleo físico do processador (core) é responsável por um shard (partição de dados) independente dentro do mesmo nó, processando requisições de forma isolada e paralela.



Remove locks globais

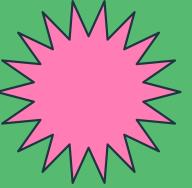


Maior aproveitamento de hardware



Processo mais previsível e eficiente

# MOTIVAÇÕES

-  **Eliminação de pausas de GC:** não depende da JVM
-  **Isolamento de carga:** graças à shard-per-core
-  **Desempenho superior:** testes e benchmarks
-  **Facilidade de manutenção:** compactação, reparo e monitoramento

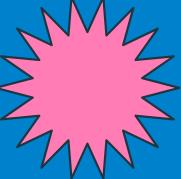
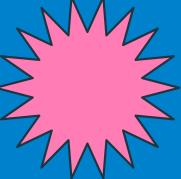
# BENCHMARKS

Métrica	ScyllaDB (shard-per-core)	Cassandra (tradicional)
Throughput máximo (ops/s)	até 300k	até 80k
Latência p99 de leitura	até 4.5 ms	até 265 ms
Latência p99 de escrita	até 7.9ms	até 494 ms
Tempo para adicionar nó	~37 minutos	~107 minutos
Tempo para substituir nó	~54 minutos	~209 minutos
Tempo de compaction	~36 minutos (32x mais rápido)	~22 a 38 horas

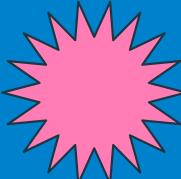
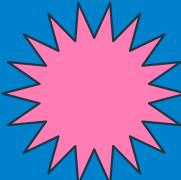
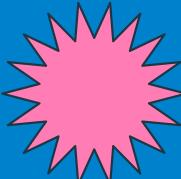
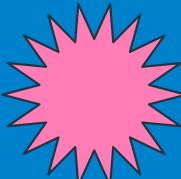
*Fontes: [ScyllaDB vs Cassandra Performance Comparison](#), [ScyllaDB Benchmarks](#), [ScyllaDB vs Apache Cassandra](#)*

# RUST DATA SERVICES

## Problema

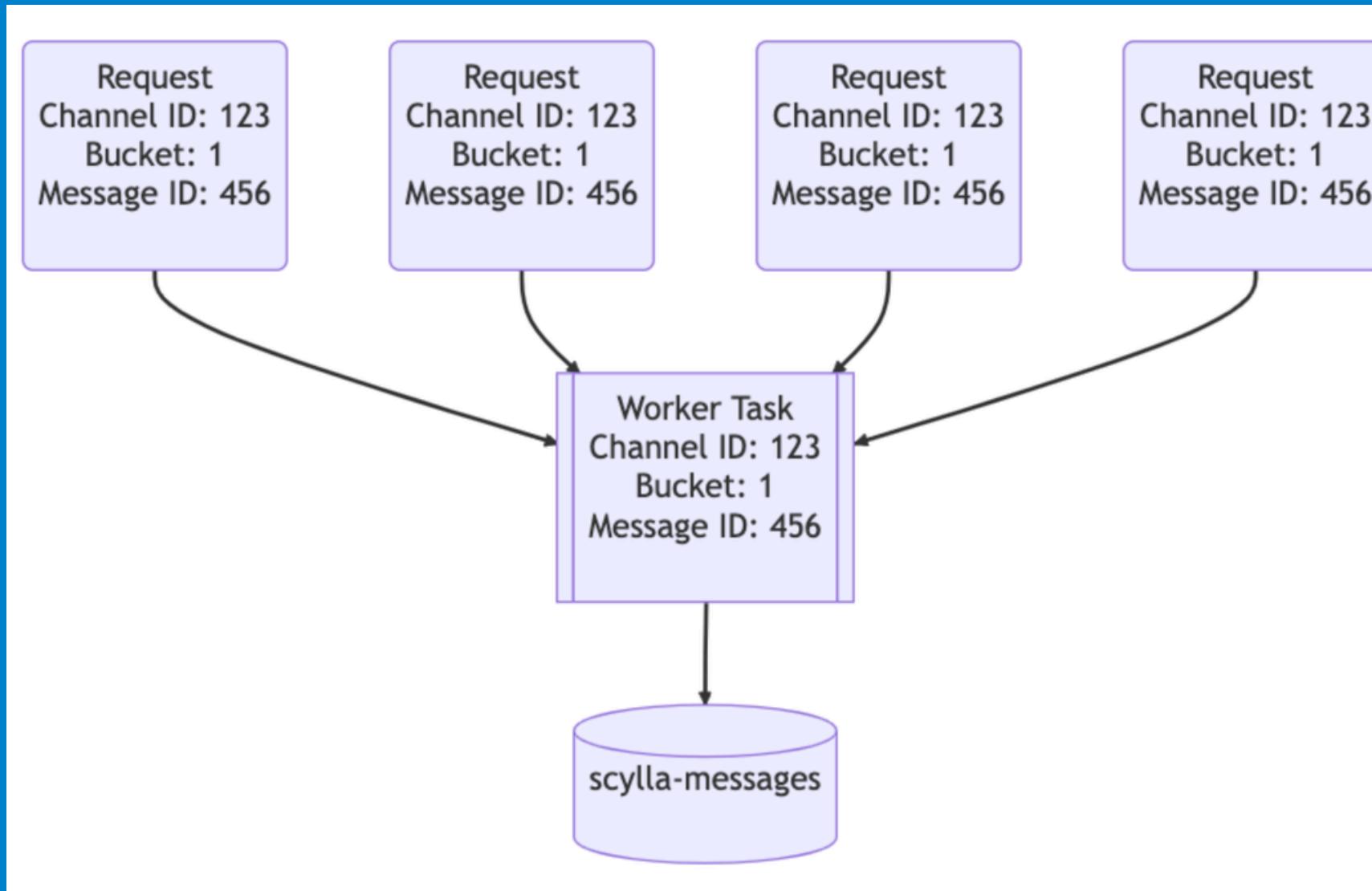
-  Picos de tráfego em determinados canais
-  Menções @everyone ou grandes anúncios/eventos

## Solução Proposta

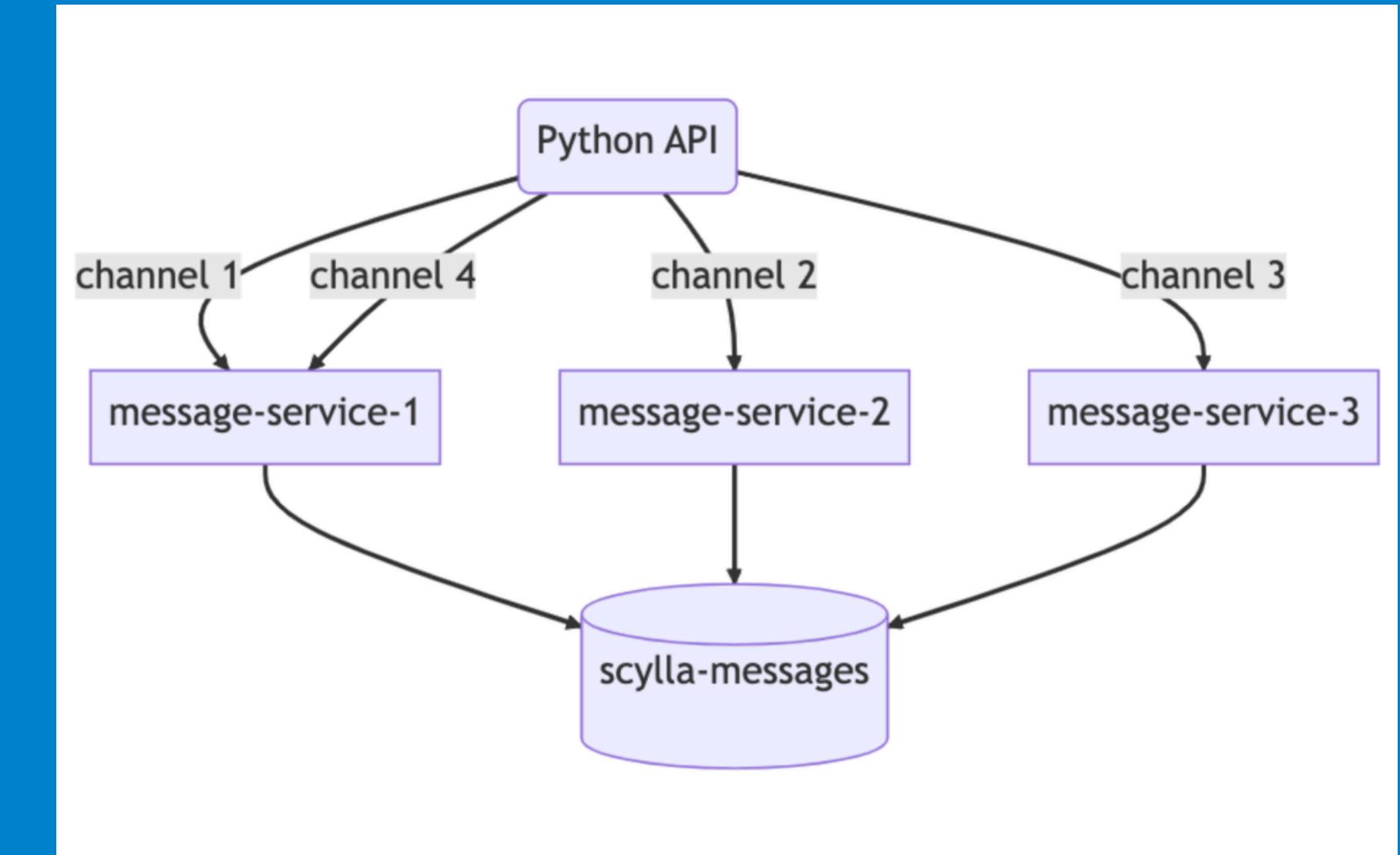
-  Criação de intermediário entre API e BD
-  Coalescência de requisições
-  Roteamento consistente
-  Concorrência segura

# RUST DATA SERVICES

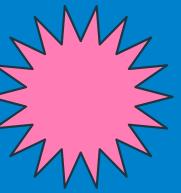
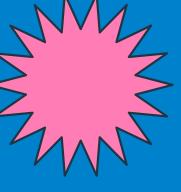
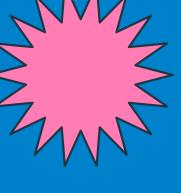
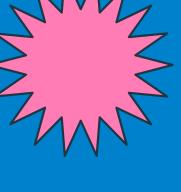
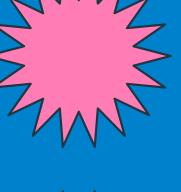
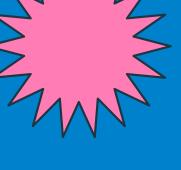
## Coalescência de Requisições



## Roteamento Consistente

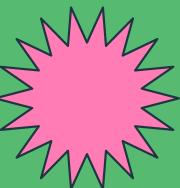
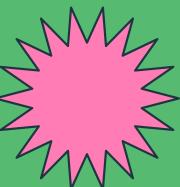


# A MIGRAÇÃO

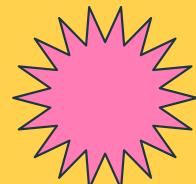
-  **Leitura reversa:** impeditivo, essencial para a aplicação
-  **Provisionamento do novo cluster:** SSDs com RAID
-  **Desenvolvimento de ferramenta de migração**
-  **Gravação duplicada:** novas mensagens salvas no Scylla e Cassandra
-  **Validação automática:** comparações entre leituras nos dois bancos
-  **Resolução de tombstones:** travamento nos últimos 0,0001%

# MODELAGEM DE DADOS

NÃO MUDOU!

-  Chave primária (channel\_id, bucket)
-  Ordenação por message\_id DESC
-  Otimizações nativas ScyllaDB como a distribuição otimizada de shards

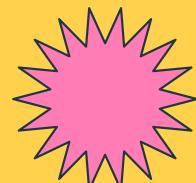
# RESULTADOS



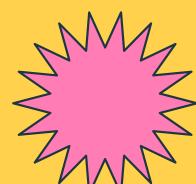
**Redução de nós:** 177 para 72



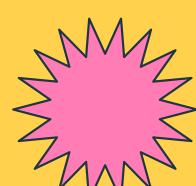
**Aumento de capacidade dos nós:** 9TB vs 4TB



**Latência reduzida:** p99 de 40-125ms para 15ms em leitura.  
Estabilizou em 5ms na escrita.



**Redução de custo operacional:** menos nós garantem menos manutenções manuais e maior previsibilidade



**Resiliência comprovada:** suportou picos globais sem degradação significativa de performance