

Uma Arquitetura de Stream Processing e ETL Serverless na AWS

Maycon Viana Bordin
Data Engineer @ Sicredi





Sicredi



*Uma instituição
financeira diferente.
Uma instituição
financeira cooperativa.*

TRANSFORMAÇÃO DIGITAL



Novo Core Bancário



Plataforma para
inovação

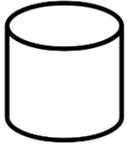


Experiência digital

DATA LAKE do **WOOOP**



Files



RDMS



NoSQL



Events

GOVERNANCE



Credentials / Roles / Permissions

INTERNAL CONSUMERS




Metabase

INGESTION



Kinesis



Microservices



SQS



S3



PROCESSING



EMR



Lambda



presto

STORAGE



S3



RDMS



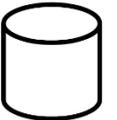
MongoDB



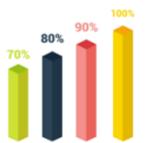
Tableau



CSV Files



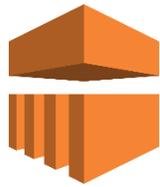
RDMS



Reports

Serverless

Serverless na



EMR



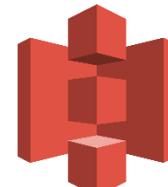
Lambda



SQS



Kinesis



S3



RDS



DynamoDB



SNS



ElasticSearch
Service



Athena



Glue



Redshift



QuickSight

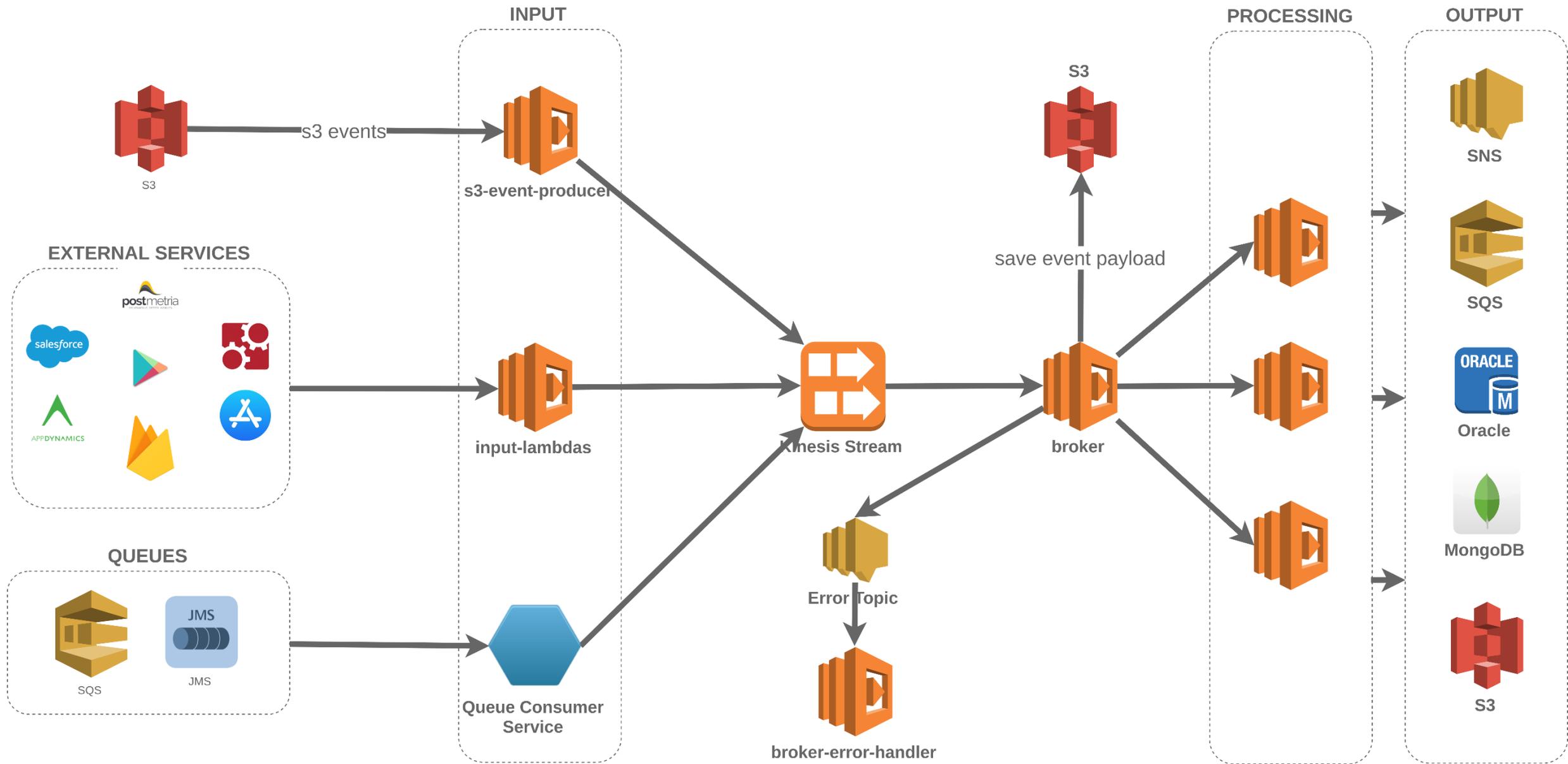


Step Functions

Por que escolhemos Serverless?



A Arquitetura



e tudo escrito em



A photograph of a Parisian street scene. In the foreground, the back of a man's head and shoulder is visible as he looks towards a building. The building has a light-colored facade with rectangular panels. The street is narrow and lined with tall, multi-story buildings with many windows and balconies. The sky is visible in the distance, showing a mix of blue and white clouds. The overall atmosphere is that of a busy, historic city street.

O Framework dentro do Framework

Processamento de Eventos

Reprocessamento

Controle de Fluxo (Backpressure)

Broker de Eventos

Controle de Erros

Garantia de Processamento

Controle de Estado

Dados

Formato Padronizado

Tipos de Dados

Schema

Transporte

Kinesis + KPL

Codec ORC

Cópia de Eventos no S3

Metadados

100% Automatizado

Dashboard

Métricas

Troubleshooting

Aplicações

SQS

SNS

JSON

CSV

XML

XLS

Oracle

MongoDB

Teradata

S3

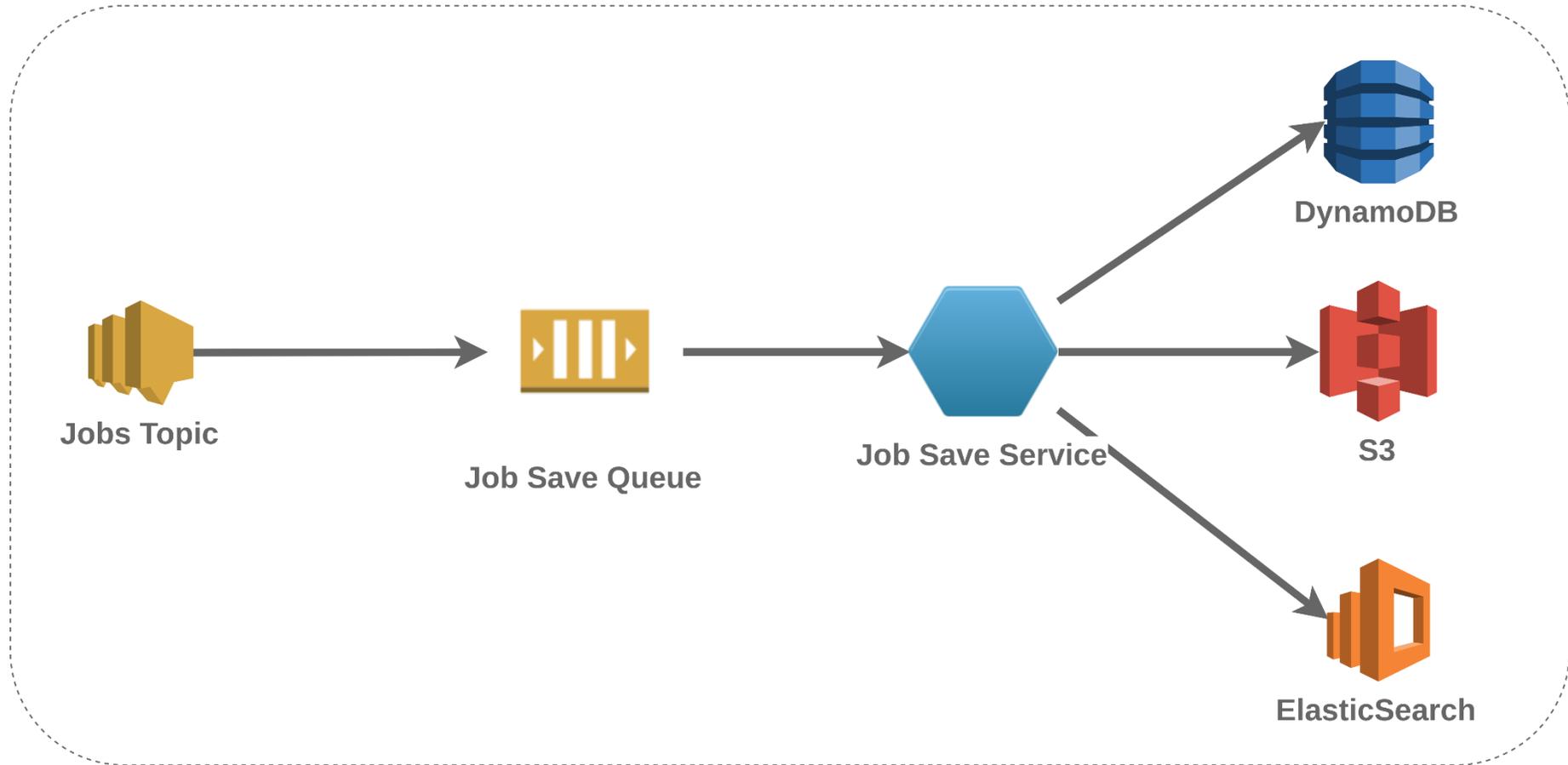
 `JsonMapProcessor.scala`

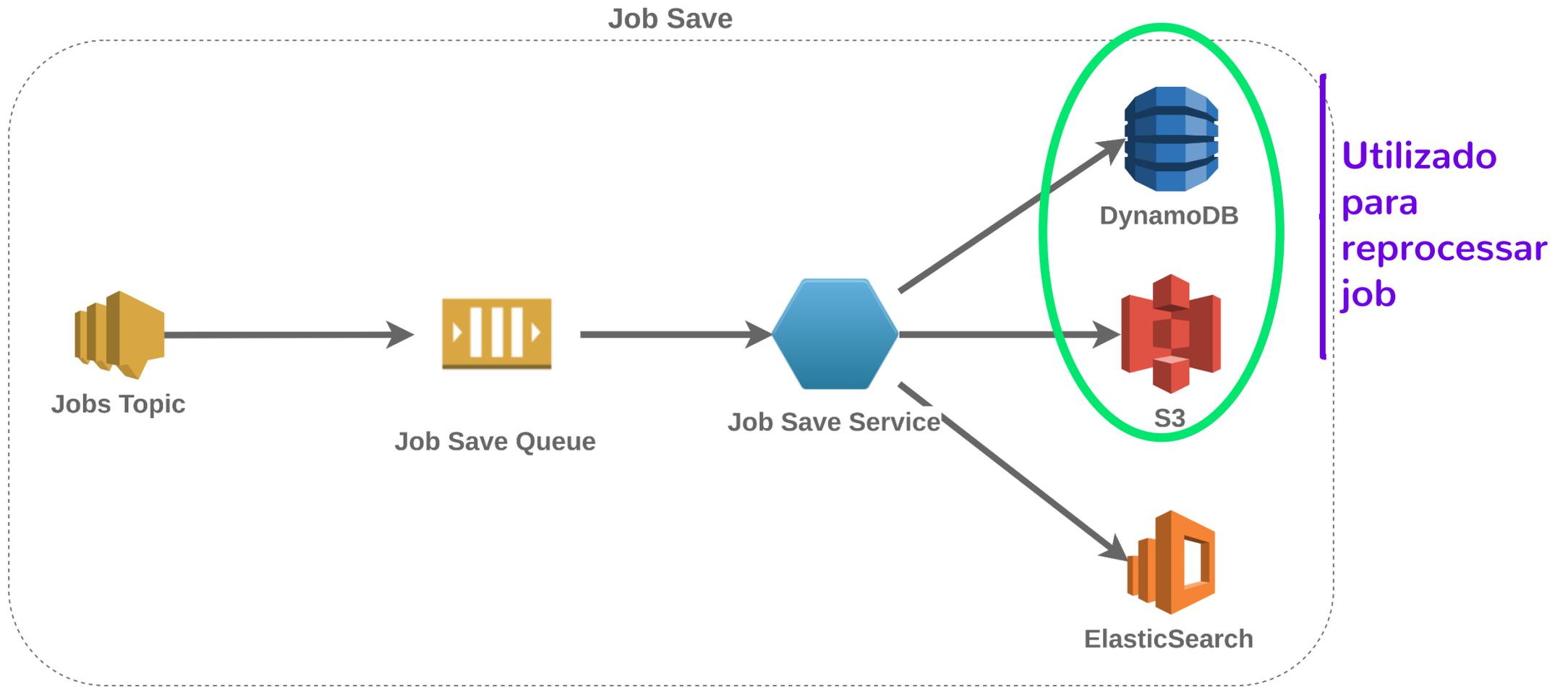
Raw

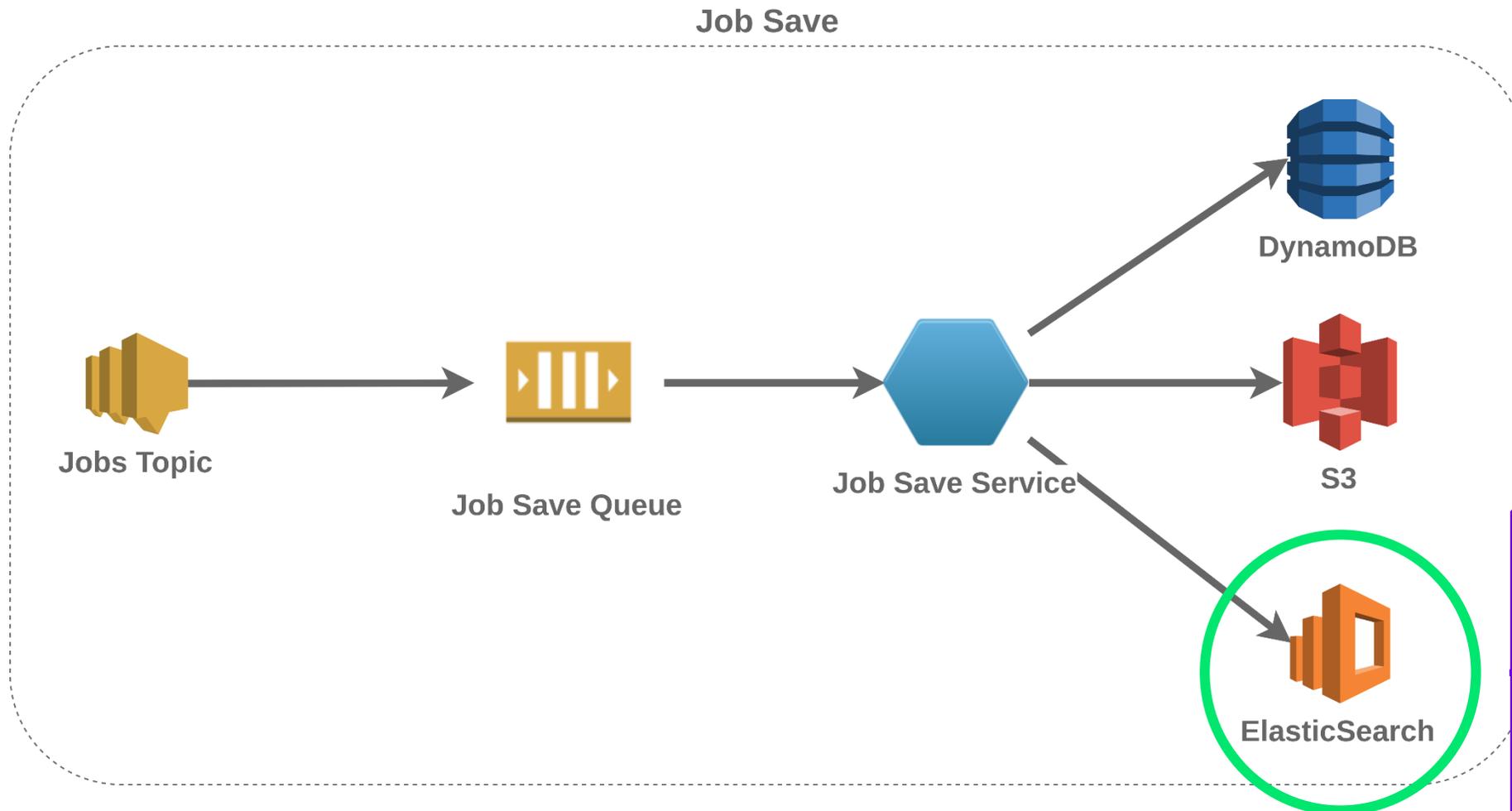
```
1  class JsonMapProcessor extends BaseProcessor {
2    lazy val fields: Map[String, List[String]] = config.getStringListMap("fields")
3    lazy val jsonField: String                = config.getString("json_field").getOrElse("data")
4    lazy val outputStream: String            = config.getOutputStream
5
6    override def process(input: EventRecord): Unit = {
7      val rawData = input.getString(jsonField).get
8
9      if (rawData == null || rawData.trim.isEmpty) {
10         log.info("Record is empty, ignoring it.")
11         return
12       }
13
14       val record = fields.map(field => {
15         field._1 -> JsonPathUtil.readToStringMultiple(field._2, rawData)
16       })
17
18       emit(outputStream, record)
19     }
20 }
```

Metadados

Job Save

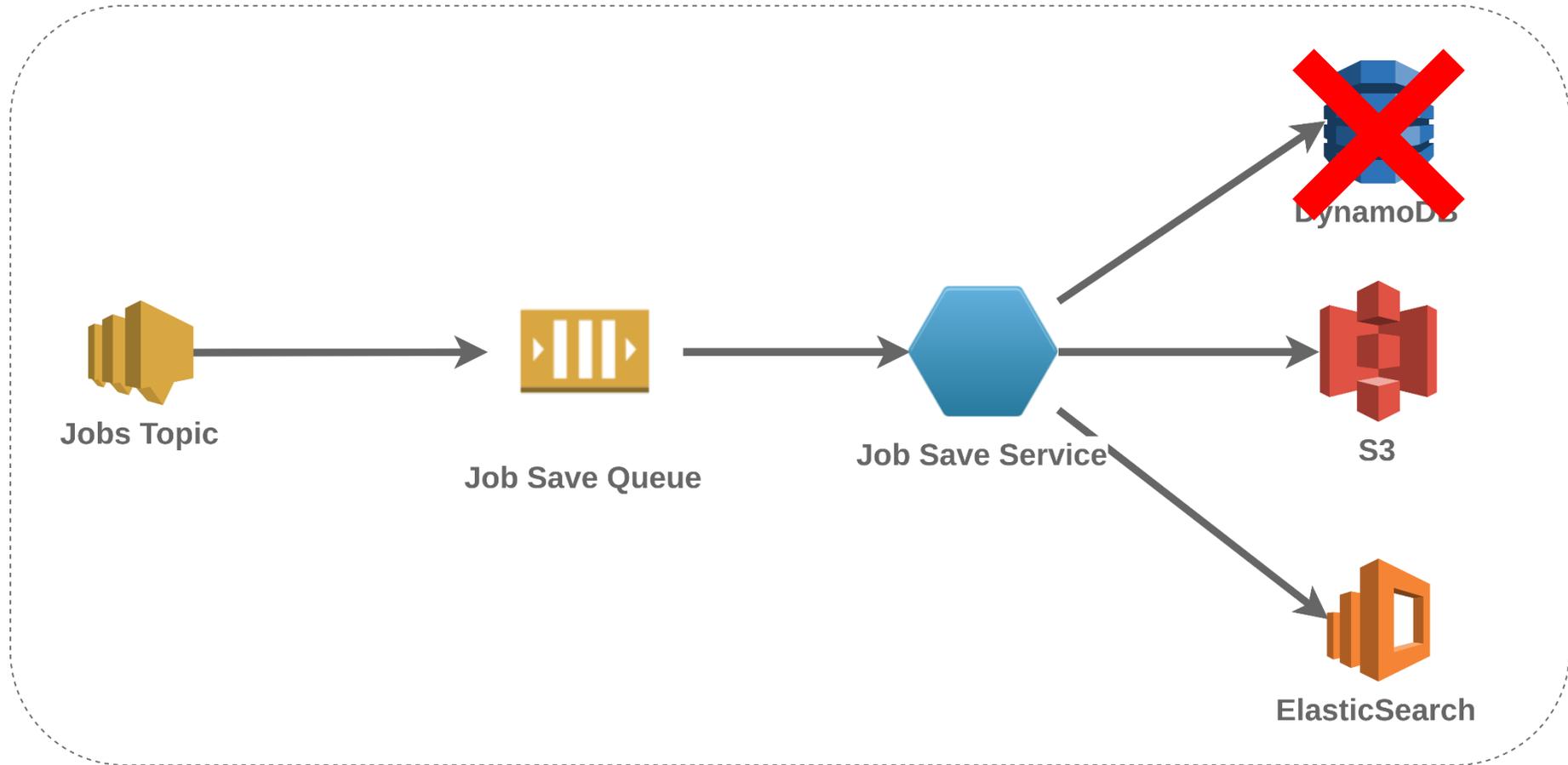






Utilizado para buscas e controle de paralelismo e de falhas

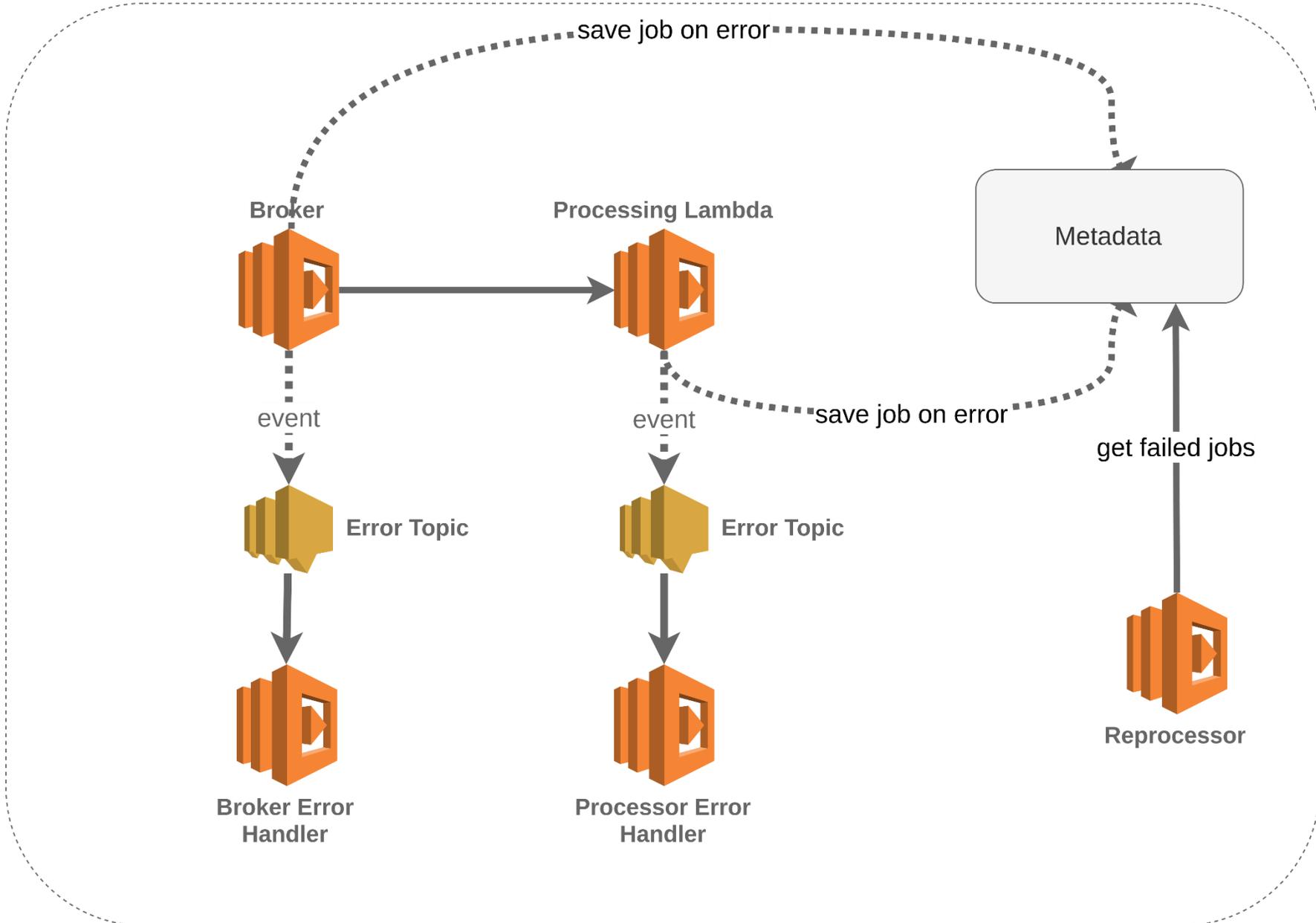
Job Save



A black and white photograph of a dog sitting on a laptop keyboard. The dog is positioned in the center of the keyboard, with its front paws on the keys. The laptop is open, and the screen is visible above the keyboard. The background is a textured, light-colored surface, possibly sand or dirt. A semi-transparent dark horizontal band is overlaid across the middle of the image, containing the text "Lidando com Erros" in a bright green, sans-serif font.

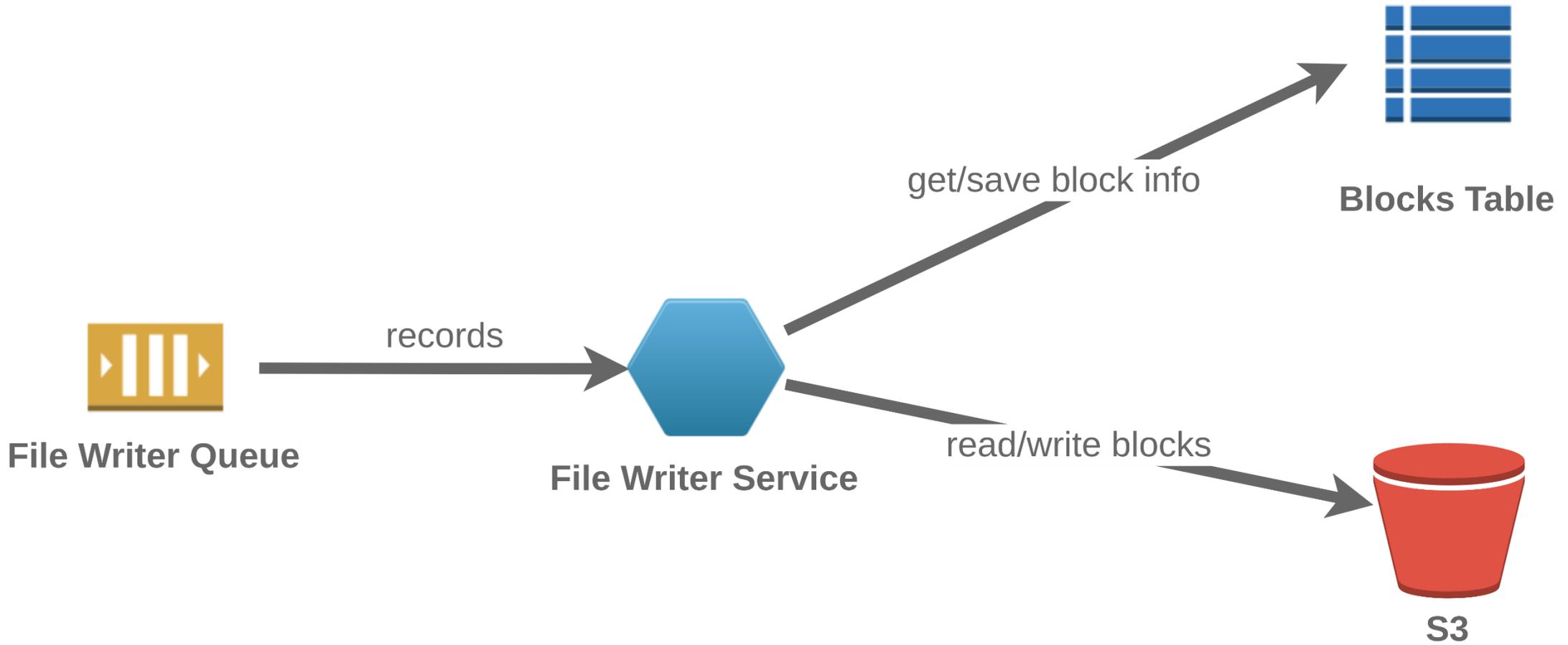
Lidando com Erros

Error Handling



Do Oracle ao (S3 + Presto)
em
Near Real Time

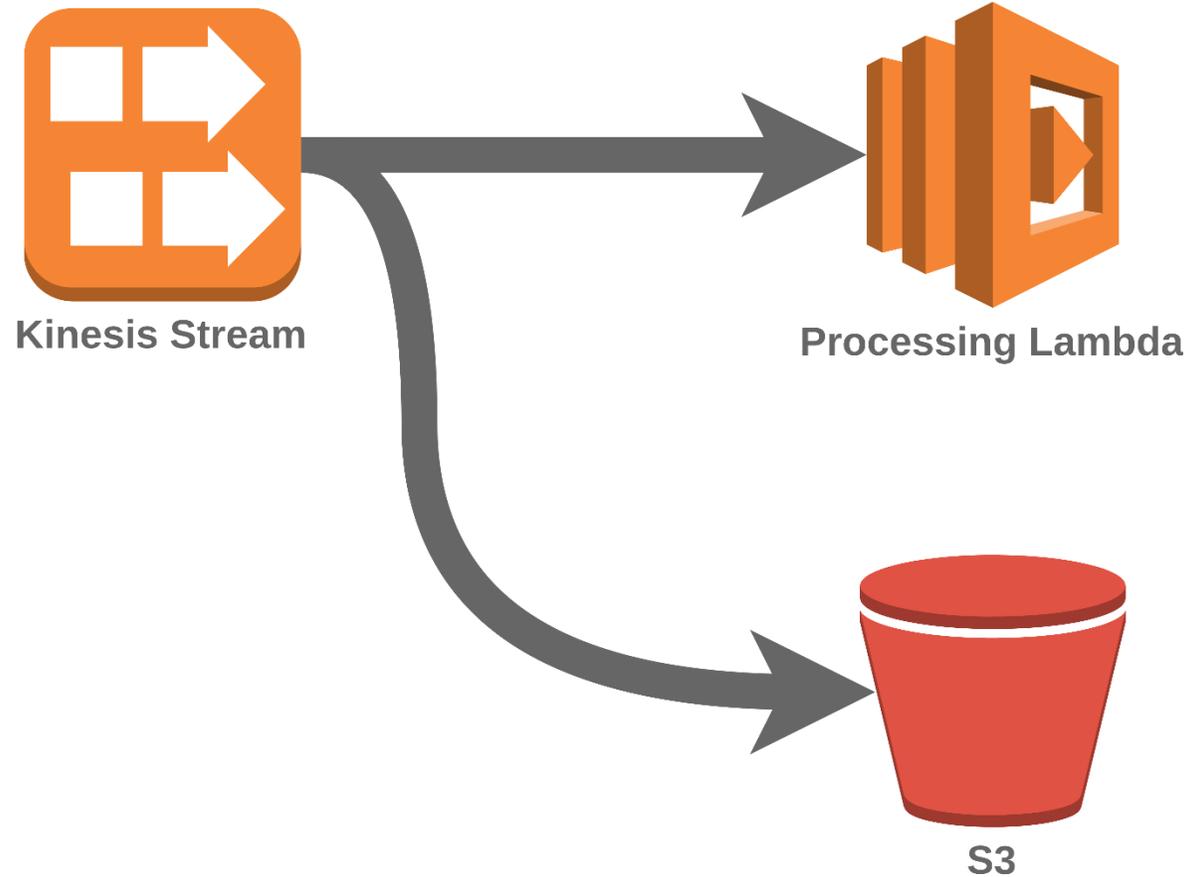
File Writer



Amazon Kinesis



Faça Cópia das suas Streams!





Quanto custa uma stream?

1 shard = 1 MB/s ou 1000
registros/s

Você paga pela hora/shard

E por cada carga de **PUT** de
25KB

A grayscale photograph of a hand turning a faucet handle. Water is flowing from the faucet into a sink. The background is a plain wall.

Reduzindo os custos do Kinesis

- Faça agregação dos dados com o **Kinesis Producer Library (KPL)**



Reduzindo os custos do Kinesis

- Faça agregação dos dados com o **Kinesis Producer Library** (KPL)
- Dimensione corretamente suas streams

A grayscale background image showing a hand turning a faucet on, with water flowing into a sink. The hand is positioned as if just finished turning the handle, and a thin stream of water is falling from the spout. The overall scene is dimly lit, emphasizing the contrast between the hand and the water.

Reduzindo os custos do Kinesis

- Faça agregação dos dados com o **Kinesis Producer Library** (KPL)
- Dimensione corretamente suas streams
- Uma stream **X** N Streams

Como Utilizamos o Kinesis

- Utilizamos apenas uma stream no Kinesis com streams “lógicas” identificadas pela partitioning key.

Como Utilizamos o Kinesis

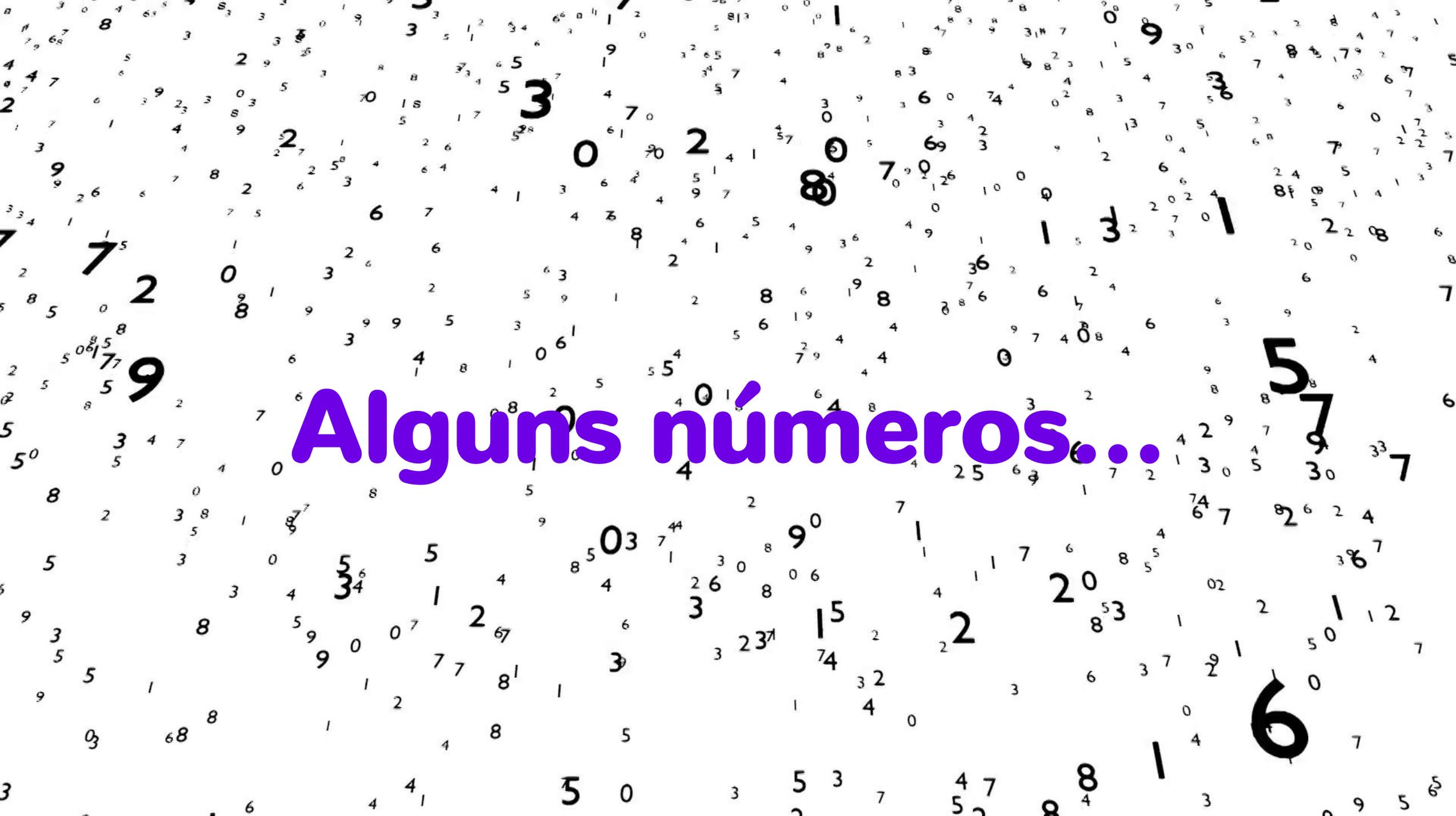
- Utilizamos apenas uma stream no Kinesis com streams “lógicas” identificadas pela partitioning key.
- Agregamos os dados em micro-batches, e ainda utilizamos o KPL.

Como Utilizamos o Kinesis

- Utilizamos apenas uma stream no Kinesis com streams “lógicas” identificadas pela partitioning key.
- Agregamos os dados em micro-batches, e ainda utilizamos o KPL.
- Codificamos os dados em **ORC** com strings em binário.

Como Utilizamos o Kinesis

- Utilizamos apenas uma stream no Kinesis com streams “lógicas” identificadas pela partitioning key.
- Agregamos os dados em micro-batches, e ainda utilizamos o KPL.
- Codificamos os dados em **ORC** com strings em binário.
- Todas streams possuem schema, mantemos o caos na RAW.



Alguns números...

12
Milhões
jobs executados

192
Milhões
registros processados

+120 GB
dados processados



PONTOS FORTES
PONTOS FRACOS



PONTOS FORTES

- Pouco gerenciamento
- Integração simples entre serviços da AWS
- Estabilidade dos serviços
- Facilidade em escalar serviços

PONTOS FRACOS

- Custos podem crescer rapidamente
- Subutilização de recursos (ex. **Kinesis**)
- Explosão de chamadas Lambda
- Saiba quando sua Lambda deve virar um microserviço ou EC2



Nossa arquitetura está em constante **evolução**. Começamos pequeno e crescemos de forma **orgânica**, na medida em que os cenários mudam.

Obrigado!



mayconbordin



mayconbordin@gmail.com

