

# Relatório Técnico

## 1. Objetivo do Case

O objetivo do case é desenvolver uma solução de engenharia de dados para consumir informações de partidas do jogo League of Legends em tempo real utilizando a API oficial da Riot Games. O fluxo de dados deve ser processado e armazenado em um Data Lake, estruturado em camadas (bronze, silver, gold). Além disso, a solução deve permitir monitoramento e observabilidade, garantindo a rastreabilidade do fluxo de dados, detecção de falhas e identificação de gargalos de desempenho.

Github: <https://github.com/diegomeyer/datamaster>

### 1.1. Conta Developer Riot

- Crie uma conta no portal Riot Games
- Faça a verificação do email cadastrado
- Aceite os termos
- Crie de fato a conta de developer
- Acesse o link Developer Portal
- Gere uma nova chave da API

### 1.2. Configurando o Ambiente

- Execute o comando  
`./start_services.sh API_KEY`

## 2. Arquitetura de Solução e Arquitetura Técnica

### 2.1. Arquitetura de Solução

A solução foi projetada em um pipeline de dados com as seguintes etapas principais:

- Extração:
  - i. Dados de partidas são consumidos da API da Riot Games.
  - ii. As informações são enviadas para o tópico summoners no Kafka.
- Processamento em Tempo Real:
  - i. Dados do Kafka são consumidos pelo Spark.
  - ii. Os dados são escritos na camada bronze do Data Lake.
- Transformação e Limpeza:
  - i. Dados da camada bronze são processados e normalizados para a camada silver.
- Agregação:
  - i. A camada gold é criada com agregações e métricas prontas para análise, como KPIs.
- Monitoramento:

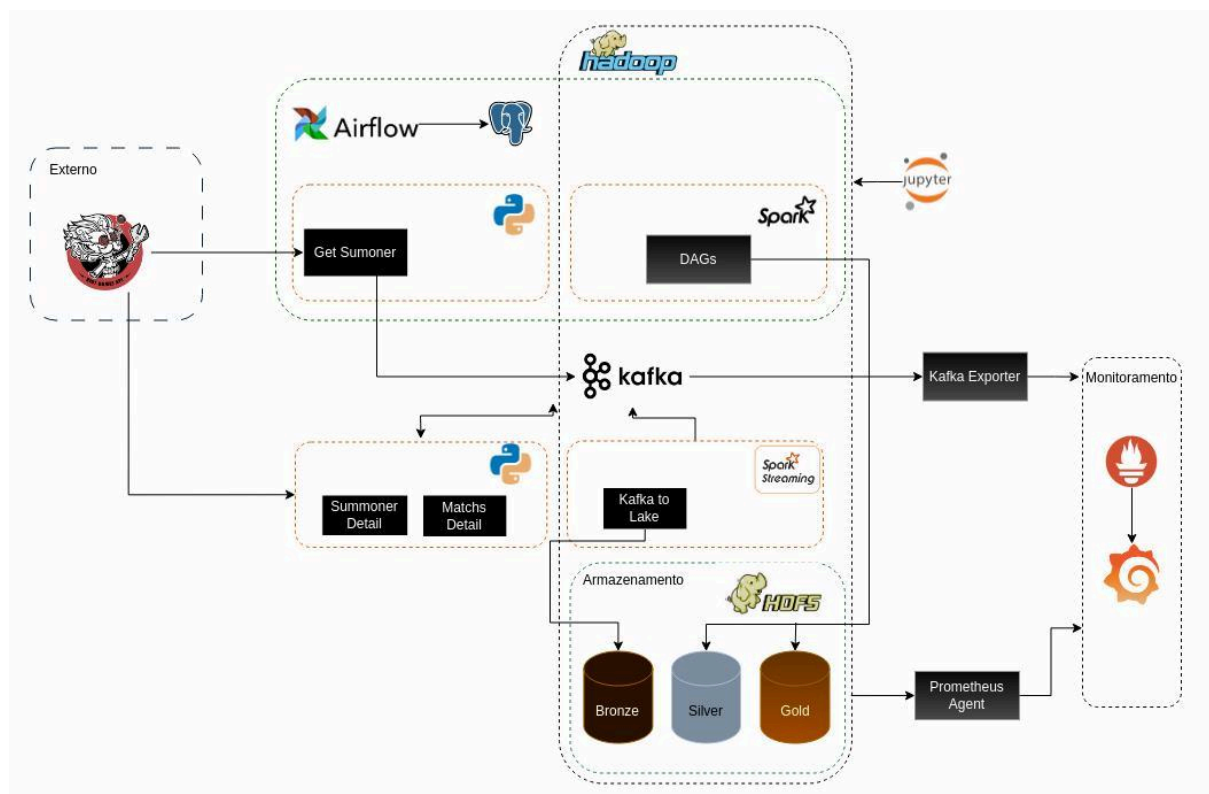
- i. Métricas de todo o pipeline são monitorados usando Prometheus e Grafana.

## 2.2. Arquitetura Técnica

### Tecnologias Utilizadas:

Tecnologia	Função
Kafka	Sistema de mensageria para ingestão de dados em tempo real.
Kafka Exporter	Export as métricas do Kafka para o Prometheus
Spark	Processamento distribuído de dados.
HDFS	Armazenamento em Data Lake com suporte a grandes volumes.
Jupyter	Ambiente para exploração dos dados
Airflow	Scheduler Jobs
Prometheus	Coleta de métricas para monitoramento.
Grafana	Visualização de métricas em dashboards.
Docker Compose	Orquestração dos serviços.

### Arquitetura:



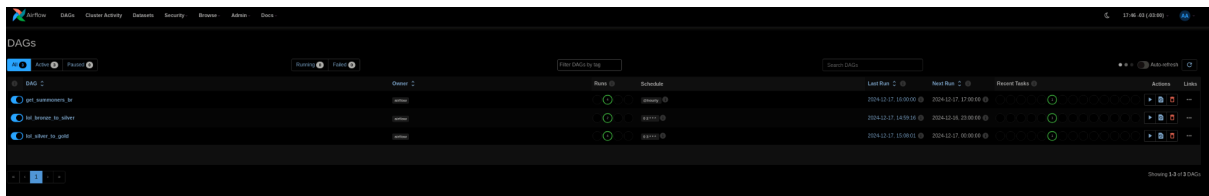
### 3. Explicação sobre o Case Desenvolvido

#### 3.1. Extração de Dados

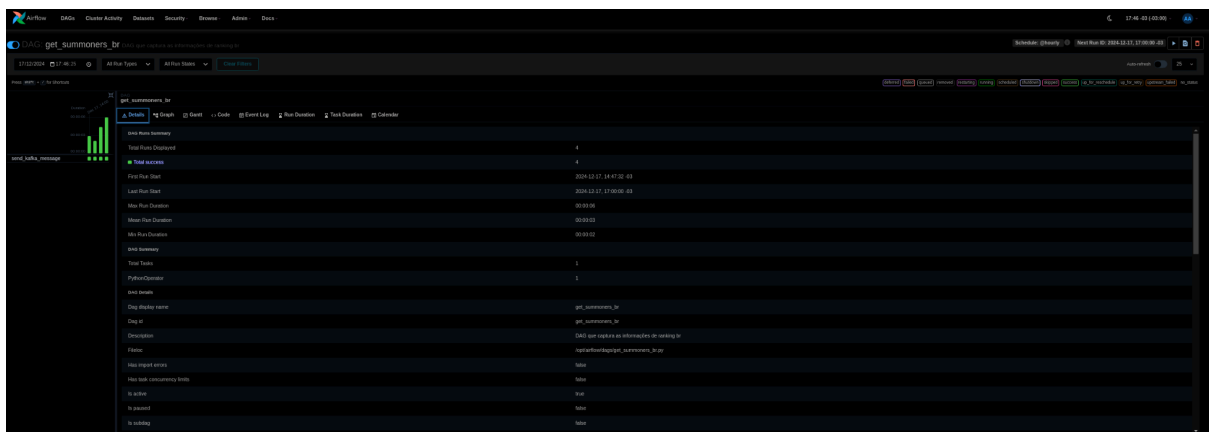
O script `get_summoners_br.py` consulta a API da Riot Games e está agendado para executar de hora em hora, onde coleta os IDs dos jogadores da região do brasil e envia informações sobre partidas para o tópico Kafka `summoners`.

O script streaming `kafka_summoner_details.py` recebe o ID do jogador e busca os IDs de das últimas 3 partidas e envia esses IDs para o topico Kafka `summoners_details`

O script streaming `kafka_matches.py` recebe o ID da partida e busca as informações da partida e envia o json para o tópico Kafka `matches`.



DAGs



DAG - Extração de Dados

#### 3.2. Processamento de Dados

O script `consumer_kafka_to_lake.py`:

- Consome mensagens do Kafka em tempo real no tópico match.
- Escreve os dados brutos e as métricas de cada etapa na camada bronze do Data Lake em formato Parquet.
- Assegura tolerância a falhas com checkpoints.

```
try:
    spark.read.parquet("hdfs://hadoop-namenode:8020/datalake/bronze/matches/*/*.parquet").printSchema()
except AnalysisException as e:
    print(f"Erro do Spark ao acessar os dados: {e}")

root
|-- timestamp_summoner: string (nullable = true)
|-- timestamp_summoner_details: string (nullable = true)
|-- timestamp_match: string (nullable = true)
|-- match_details: string (nullable = true)
```

Schema - Tabela Bronze

O script bronze to silver.py

- DAG que consome a camada bronze e extrai e estrutura as informações mais relevantes.

The screenshot shows the Databricks Jobs console for a job named 'ice\_bronze\_to\_silver'. The job is in a 'Completed' state. The 'Job Summary' section shows that the job was created on 2024-12-17 at 14:58:26.03. The 'Job Details' section shows that the job was created by 'ice\_bronze\_to\_silver' and is associated with the 'ice\_bronze\_to\_silver' cluster. The 'Job History' section shows that the job was completed on 2024-12-17 at 15:00:26.03. The 'Job Details' section also shows the job's configuration, including the job name, job ID, and the job's description.

Job Name	Job ID	Job Status	Job Creation Time	Job Completion Time	Job Duration
ice_bronze_to_silver	12345678901234567890	Completed	2024-12-17 14:58:26.03	2024-12-17 15:00:26.03	00:02:00

## DAG - Bronze to Silver

```
try:
    spark.read.parquet("hdfs://hadoop-namenode:8020/datalake/silver/games/**/*.parquet").show(2,False)
except AnalysisException as e:
    print("Erro do Spark ao acessar os dados: {}".format(e))

[match_id | game_creation|game_duration|game_mode|game_version |map_id|queue_id|platform_id|game_start_timestamp|game_end_timestamp
3083995782 |1734202378631|1093 |CLASSIC |14.24.64.5128|11 |420 |BR1 |1734202791905 |1734122172672
3083995782 |1734381028131|1063 |CLASSIC |14.24.64.5128|11 |420 |BR1 |1734381855205 |1734382918667

only showing top 2 rows
```

```
try: spark.read.parquet("hdfs://hadoop-namenode:8020/datalake/silver/participants/*-parquet").orderBy("win").show(2,False)
except AnalysisException as e:
    print(f'Erro do Spark ao acessar os dados: {e}')

[match_id | participant_id | champion_id | champion_name | team_id | individual_position | lane | kills | deaths | assists | item1 | item2 | item3 | item4 | item5 | item6 | allInPings | itemsPurchased | wardsKilled | wardsPlaced | total_damage_dealt | total_damage_taken | gold_earned | win | ]
303331483414 | 61 | 10rn | 100 | MIDLE | MIDLE | 8 | 14 | 3158 | 16655 | 2031 | 1056 | 1645 | 1056 | 1393 | 0 | 23 | 3 | 11 | 65993 | 21842 | 13036 | 13036 | false |
303331483419 | 516 | 10rn | 100 | TOP | TOP | 2 | 5 | 3047 | 3076 | 3068 | 1031 | 1031 | 1057 | 3340 | 0 | 17 | 3 | 8 | 98327 | 39697 | 7836 | 7836 | false |
```

```
try:
    spark.read.parquet("hdfs://hadoop-namenode:8020/datalake/silver/teams_bans/*.*.parquet").show(2,False)
except AnalysisException as e:
    print(f'Erro do Spark ao acessar os dados: {e}')

+-----+-----+-----+-----+
|match_id|team_id|ban_turn|champion_id|
+-----+-----+-----+-----+
|3040166917|100|15|55|1|
|3040295332|200|9|20|1|
+-----+-----+-----+-----+

only showing top 2 rows
```

```
try:
    spark.read.parquet('hdfs://hadoop-namenode:8020/datalake/silver/teams_stats/**/*.parquet').show(2,False)
except AnalysisException as e:
    print('Erro do Spark ao acessar os dados: {}'.format(e))
```

```
-----+-----+-----+-----+
|match_id|team_id|win|baron_kills|dragon_kills|tower_kills|
-----+-----+-----+-----+
|3040250383|200|true|10|12|18|1|
|3040202159|200|true|1|4|11|1|
-----+-----+-----+-----+
only showing top 2 rows
```

### Saídas das tabelas Silvers

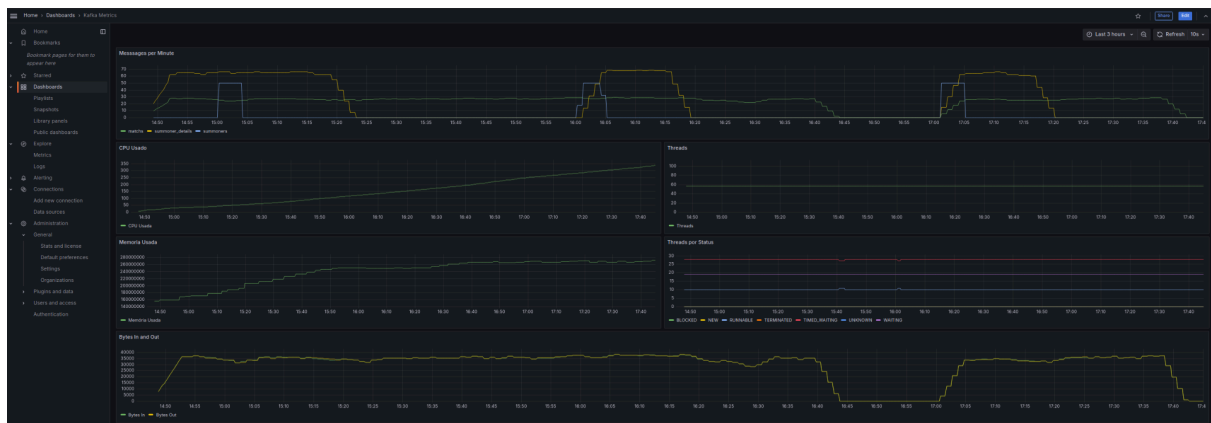
O script silver to gold.py

- DAG que consome a camada silver e realiza agregações.

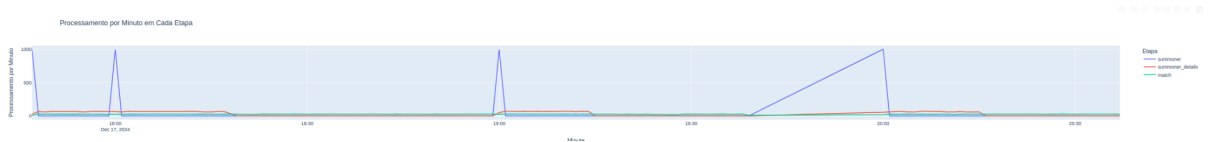




Hadoop Métricas



Kafka Métricas



Throughput - Etapas

### 3.5. LGPD

Para os dados que trabalhamos na API, não temos dados sensíveis. Caso existisse dados sensíveis poderíamos utilizar os métodos:

- Generalização
- Supressão
- K-Anonimidade
- Tokenização

## 4. Melhorias

### 4.1. Escalabilidade:

- Implementar particionamento no Kafka e paralelismo no Spark para suportar maior volume de dados.
- Configurar múltiplos nós no cluster HDFS para maior capacidade de armazenamento.

#### 4.2. Observabilidade Avançada:

- Adicionar rastreamento distribuído com OpenTelemetry para monitorar o tempo de processamento em cada componente do pipeline.

#### 4.3. Governança de Dados:

- Adicionar políticas de retenção em cada camada.
- Adicionar expurgo dos dados

#### 4.4. Segurança:

- Configurar autenticação e autorização no Kafka.
- Criptografar os dados sensíveis armazenados no Data Lake.
- Chaves de API devem ser armazenadas em um cofre de senhas.

## 5. Considerações Finais

A solução desenvolvida apresenta um pipeline robusto para ingestão, processamento e armazenamento de dados em tempo real, com suporte a monitoramento e observabilidade. As tecnologias utilizadas garantem escalabilidade e flexibilidade, atendendo às demandas do case.

Com as melhorias sugeridas, a solução pode ser ainda mais eficiente e confiável, garantindo maior governança e capacidade de lidar com volumes crescentes de dados.