

# PROJETO FINAL TECNOLOGIAS DE BIG DATA E INTELIGÊNCIA ARTIFICIAL

## GRUPO 8 | INDIGESTION

### Nome dos Alunos:

Alan Batista

Manuella Paez

Mario José C M Prado

Wislom Diogo Almeida

### Coordenador:

Prof. Fabio Jardim



# Agenda

- 1. Contextualização do trabalho
- 2. Visão e objetivo do projeto
- 3. Documentação da solução
  - i. Diagrama da arquitetura e descrição dos serviços
  - ii. Detalhamento e configurações técnicas
- 4. Demonstração da solução e entregáveis



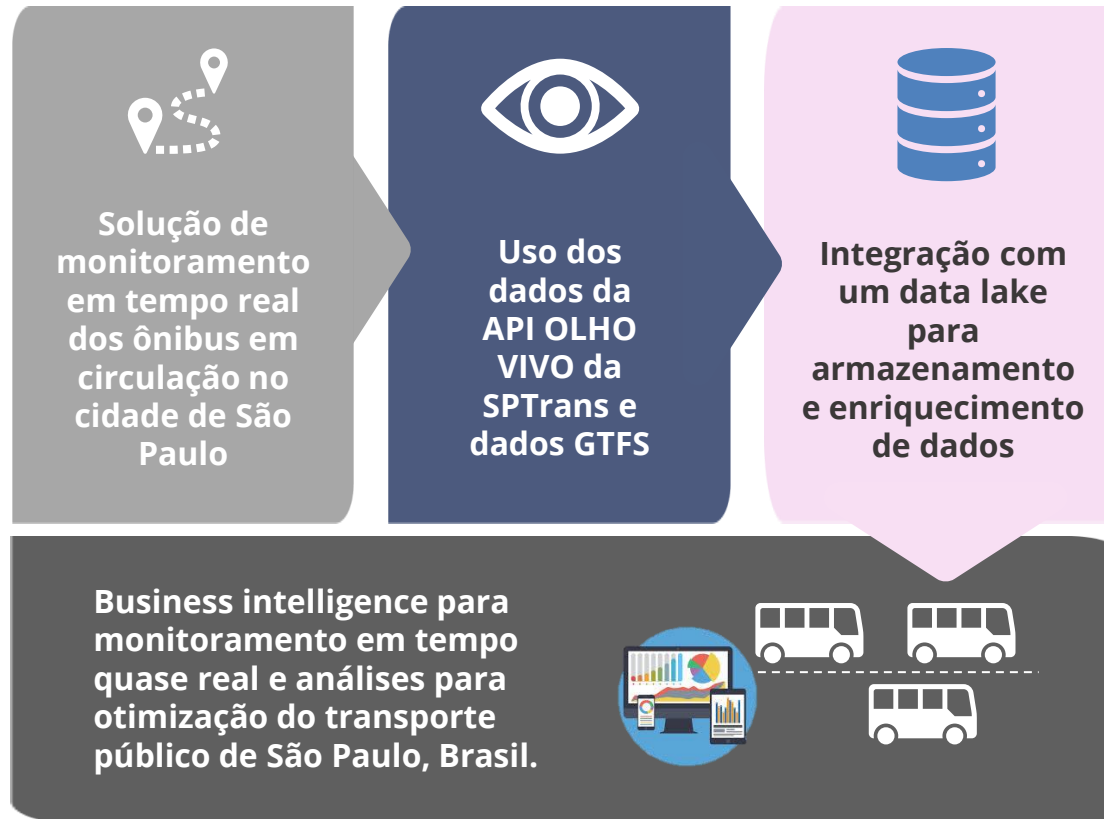
# 1. Contextualização do trabalho

O sistema de transporte público de São Paulo atende **milhões de pessoas** diariamente, e a **eficácia na gestão da frota de ônibus** é crucial para **assegurar uma boa qualidade do serviço público** prestado e, conseqüentemente, a satisfação da população com o tema.

O trabalho possui como objetivo a construção de uma **aplicação** que possibilite o **monitoramento em tempo quase real dos ônibus** em circulação no estado de São Paulo e que **ofereça métricas e KPIs importantes para tomada de decisão**.



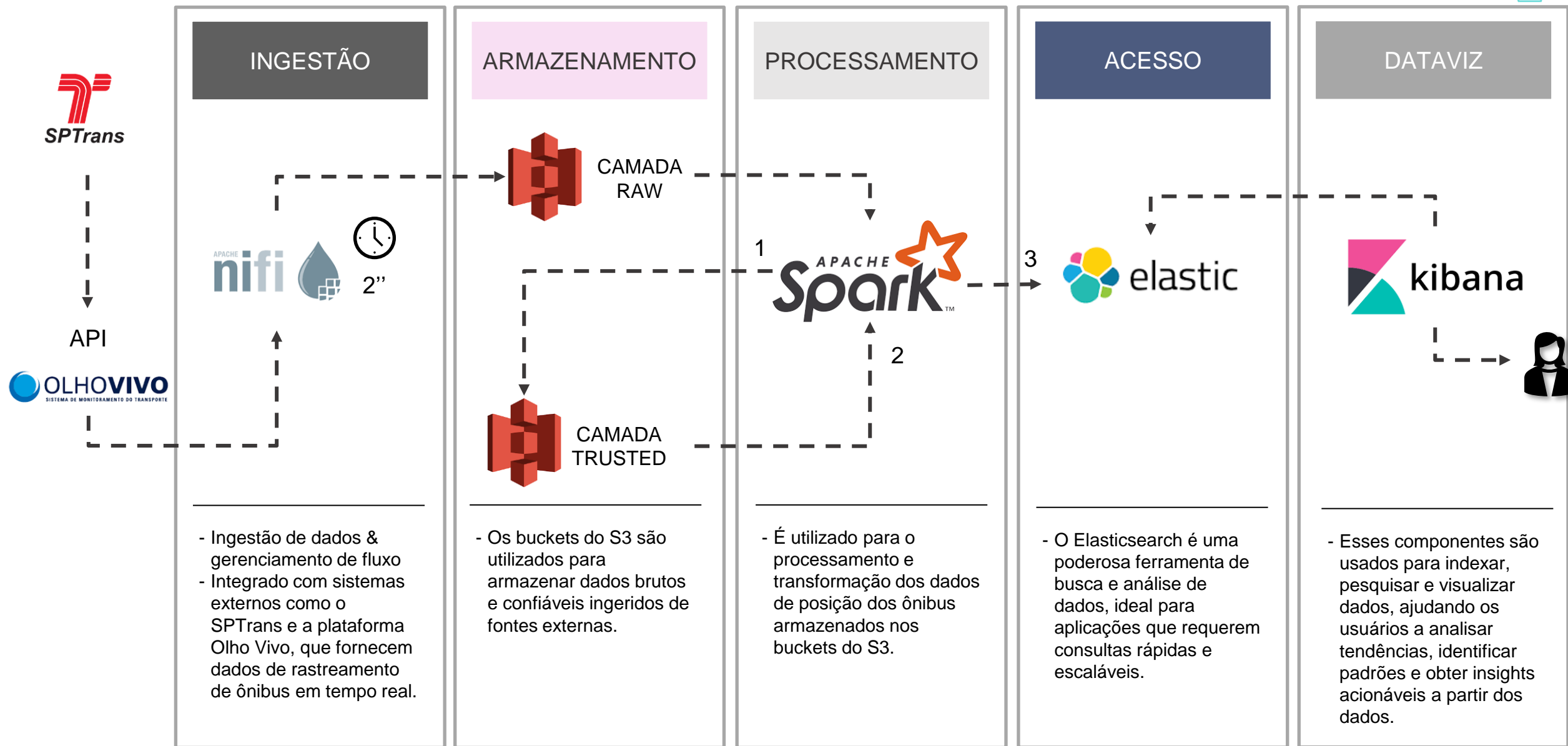
## 2. Visão e objetivo do Projeto





### 3.i Solução | Arquitetura da solução e descrição dos serviços

5



### INGESTÃO E ARMAZENAMENTO

### PROCESSAMENTO

### ACESSO E DATAVIZ

Credenciais de acesso:

#### NiFi:

- URL: <http://localhost:49090>
- Porta padrão para acesso à interface web do NiFi : 9090

#### MinIO:

Console: <http://localhost:49001>

#### Serviço de

Armazenamento: <http://localhost:49000>

#### Portas expostas:

- 9000: Porta padrão para acesso ao serviço MinIO
- 9001: Porta do console de gerenciamento

#### Credenciais de Acesso:

Usuário: admin / Senha: minioadmin



GITHUB

Detalhamentos  
adicionais

1



#### Configuração:

**Imagem Utilizada:** apache/nifi:\${NIFI\_VERSION}

**Container Name:** nifi-otmzsp

**Hostname:** nifi-otmzsp

#### Ambiente:

**NIFI\_WEB\_HTTP\_PORT:** 9090

**NIFI\_WEB\_HTTPS\_HOST:** nifi

**Timezone (TZ):** America/Sao\_Paulo

#### Volumes:

- `./volumes/nifi/util:/util` | Diretório para utilitários.
- `./volumes/nifi/util/jar:/util/jar` | Diretório para arquivos JAR.
- `./volumes/nifi/conf:/opt/nifi/nifi-current/conf` | Diretório de configuração do NiFi.

#### Comando:

```
sh -c "ln -snf /usr/share/zoneinfo/$(TZ) /etc/localtime && echo $(TZ) > /etc/timezone" | Configuração do timezone.
```

#### Recursos:

**Limite de Memória:** 2 GB para garantir desempenho adequado.

2



#### Configuração:

**Imagem Utilizada:** minio/minio:\${MINIO\_VERSION}

**Container Name:** minio-otmzsp

#### Estrutura de Pastas:

- Raw: Diretório para armazenamento de dados brutos.
- Trusted: Diretório para armazenamento de dados confiáveis.

#### Volumes:

Mapeamento do diretório local `./volumes/minio/data` para `/data` no container, garantindo persistência de dados.

#### Health Check:

- **Comando:** ["CMD", "mc", "ready", "local"]
- **Intervalo:** 5 segundos
- **Retries:** 5 tentativas
- **Timeout:** 5 segundos

### INGESTÃO E ARMAZENAMENTO

### PROCESSAMENTO

### ACESSO E DATAVIZ

Credenciais de acesso:

#### Apache SPARK:

- URL: <spark://spark-master-otmzsp:7077>

#### Portas expostas - Master:

- 7077: Porta do master para gerenciar os workers.
- 8080: Interface web do master.
- **Portas expostas - Worker:** 8081: Interface web do worker.

#### JUPYTER:

##### Portas expostas:

- 8888: Porta padrão para acessar a interface web do Jupyter Notebook.
- 4040 a 4043: Portas usadas para monitoramento do Spark (UI do Spark).



GITHUB

Detalhamentos  
adicionais



#### Configuração:

**Imagem Utilizada:** apache/spark:\${SPARK\_VERSION}

#### Container Name:

- Master: *spark-master-otmzsp*
- Worker: *spark-worker-otmzsp*

#### Modos de Operação:

##### 1. Master:

##### Variáveis de Ambiente:

SPARK\_MODE: master

SPARK\_MASTER\_HOST: spark-master-otmzsp

TZ: America/Sao\_Paulo

**Limite de Memória:** 2 GB.

**Comando:** /opt/spark/bin/spark-class  
org.apache.spark.deploy.master.Master

##### 2. Worker:

##### Variáveis de Ambiente:

SPARK\_MODE: worker

SPARK\_MASTER\_URL: spark://spark-master-otmzsp:7077

SPARK\_WORKER\_MEMORY: 1g

TZ: America/Sao\_Paulo

**Limite de Memória:** 1 GB.

**Comando:** /opt/spark/bin/spark-class  
org.apache.spark.deploy.worker.Worker spark://spark-master-otmzsp:7077



#### Configuração:

**Imagem Utilizada:** jupyter/pyspark-notebook:latest

**Container Name:** *jupyter-otmzsp*

#### Ambiente:

**JUPYTER\_TOKEN:** "" | Desabilita o token de segurança para acesso.

#### Volumes:

*../notebooks:/home/jovyan/work* | Diretório para armazenar notebooks, permitindo persistência de dados e fácil acesso.

#### Comando:

*start-notebook.sh --NotebookApp.token="" --  
NotebookApp.password=""* | Inicia o servidor Jupyter sem token ou senha, facilitando o acesso.

## 3.ii Solução | Detalhamento e configurações técnicas

8

INGESTÃO E  
ARMAZENAMENTO

PROCESSAMENTO

ACESSO E DATAVIZ

Credenciais de acesso:

### Elastic Search:

#### Portas expostas:

- 9200: Porta padrão para acesso à API REST do Elasticsearch.
- 9300: Porta para comunicação entre nós (cluster).

### Kibana:

#### Portas expostas:

- 5601: Porta padrão para acessar a interface web do Kibana.

Detalhamentos  
adicionais

5



### Configuração:

**Imagem Utilizada:** elasticsearch:7.17.20

**Container Name:** elasticsearch-otmzsp

**Hostname:** elasticsearch-otmzsp

### Ambiente:

discovery.type: single-node | configuração para executar em modo de nó único)

.ES\_JAVA\_OPTS: "-Xms2g -Xmx2g" | configurações de memória do Java

.xpack.security.enabled: "false" | desabilita a segurança para simplificar a configuração.

### Volumes:

./volumes/elasticsearch/esdata:/usr/share/elasticsearch/data | Mapeamento para persistência de dados.

### Health Check:

**Comando:** curl -sS --fail http://elasticsearch-otmzsp:9200/\_cluster/health?wait\_for\_status=yellow&time out=0s

**Intervalo:** 1 segundo

**Retries:** 3 tentativas

**Start Period:** 20 segundos

**Timeout:** 5 segundos

6



### Configuração:

**Imagem Utilizada:** kibana:7.17.20

**Container Name:** kibana-otmzsp

**Hostname:** kibana-otmzsp

### Ambiente:

ELASTICSEARCH\_HOSTS: "http://elasticsearch-otmzsp:9200" | configuração para conectar ao Elasticsearch.

### Volumes:

./volumes/kibana/data:/usr/share/kibana/data | Mapeamento para persistência de dados, garantindo que as configurações e dashboards sejam mantidos.

### Dependências:

**Elasticsearch:** O Kibana depende do Elasticsearch, sendo necessário garantir que o serviço do Elasticsearch esteja saudável antes do início do Kibana.



GITHUB



# PROJETO FINAL DATA ENGINEERING

## GRUPO 8 | INDIGESTION

### Nome dos Alunos:

Alan Batista

Manuella Paez

Mario José C M Prado

Wislom Diogo Almeida

### Coordenador:

Prof. Fabio Jardim

