

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

3



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, consequentemente, 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**.

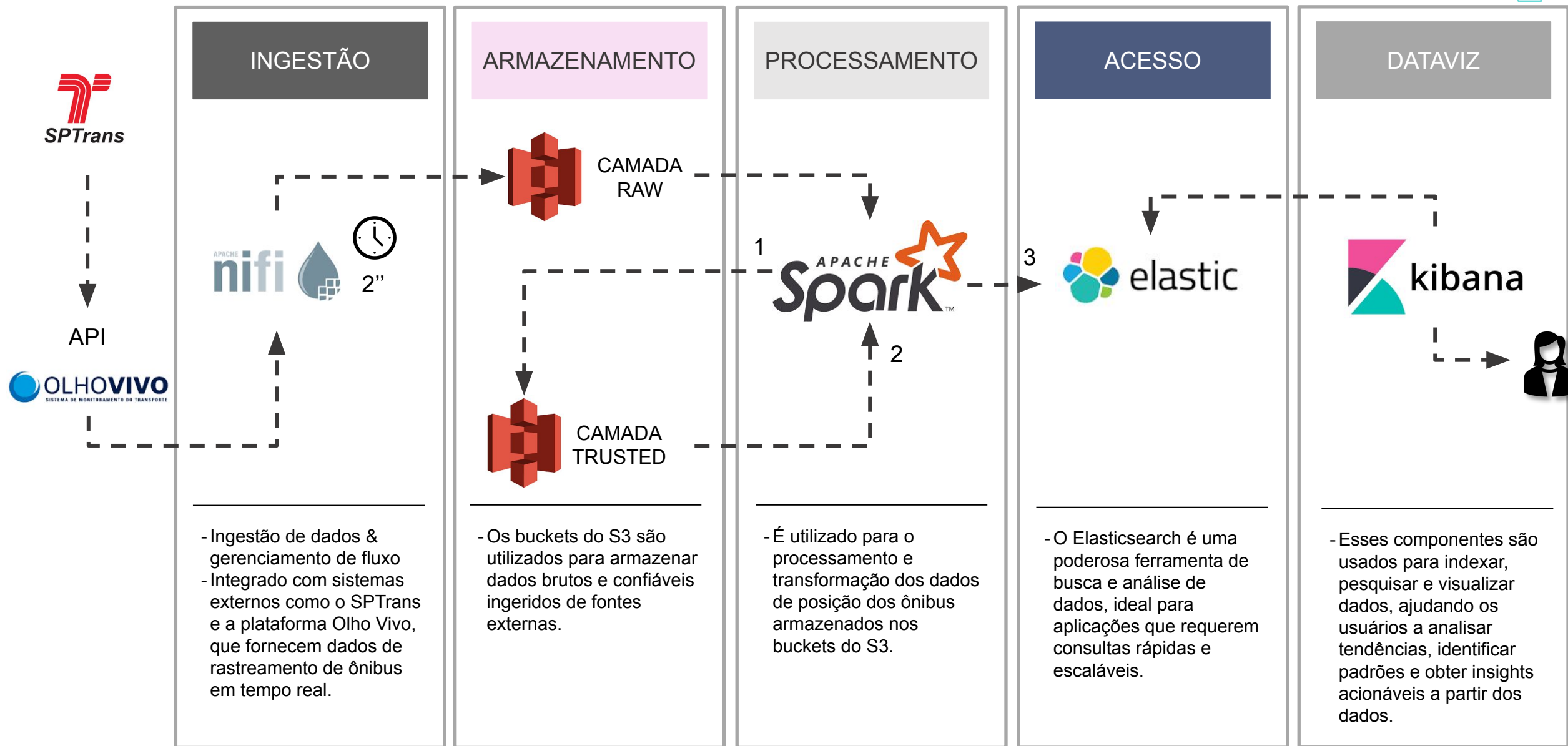


2. Visão e objetivo do Projeto



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

5



3.ii Solução | Detalhamento e configurações técnicas

6

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



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.



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 montados para persistência de dados:

`/minio_data/raw` e `/minio_data/trusted`.

Health Check:

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

3.ii Solução | Detalhamento e configurações técnicas

7

INGESTÃO E
ARMAZENAMENTO

PROCESSAMENTO

ACESSO E DATAVIZ

Credenciais de acesso:

Apache SPARK:

- URL: <http://localhost:8080> (UI do Spark Master)
- **Portas expostas - Master:**
- 7077: Porta do master para gerenciar os workers.
- URL: <http://localhost:8081> (UI do Spark Worker)
- **Portas expostas - Worker:** 8081: Interface web do worker.

JUPYTER:

- URL: <http://localhost:8888>
- **Portas expostas:**
- 8888: Porta padrão para acessar a interface web do Jupyter Notebook.
- 4040 a 4043: Portas usadas para monitoramento 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:

2. 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

3. 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

INGESTÃO E
ARMAZENAMENTO

PROCESSAMENTO

ACESSO E DATAVIZ

Credenciais de acesso:

Elastic Search:

- URL: <http://localhost:49200> (API)

Kibana:

- URL: <http://localhost:45601>



GITHUB

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 single node

.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&timeout=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.

PROJETO FINAL DATA ENGINEERING

GRUPO 8 | INDIGESTION

Nome dos Alunos:

Alan Batista

Manuella Paez

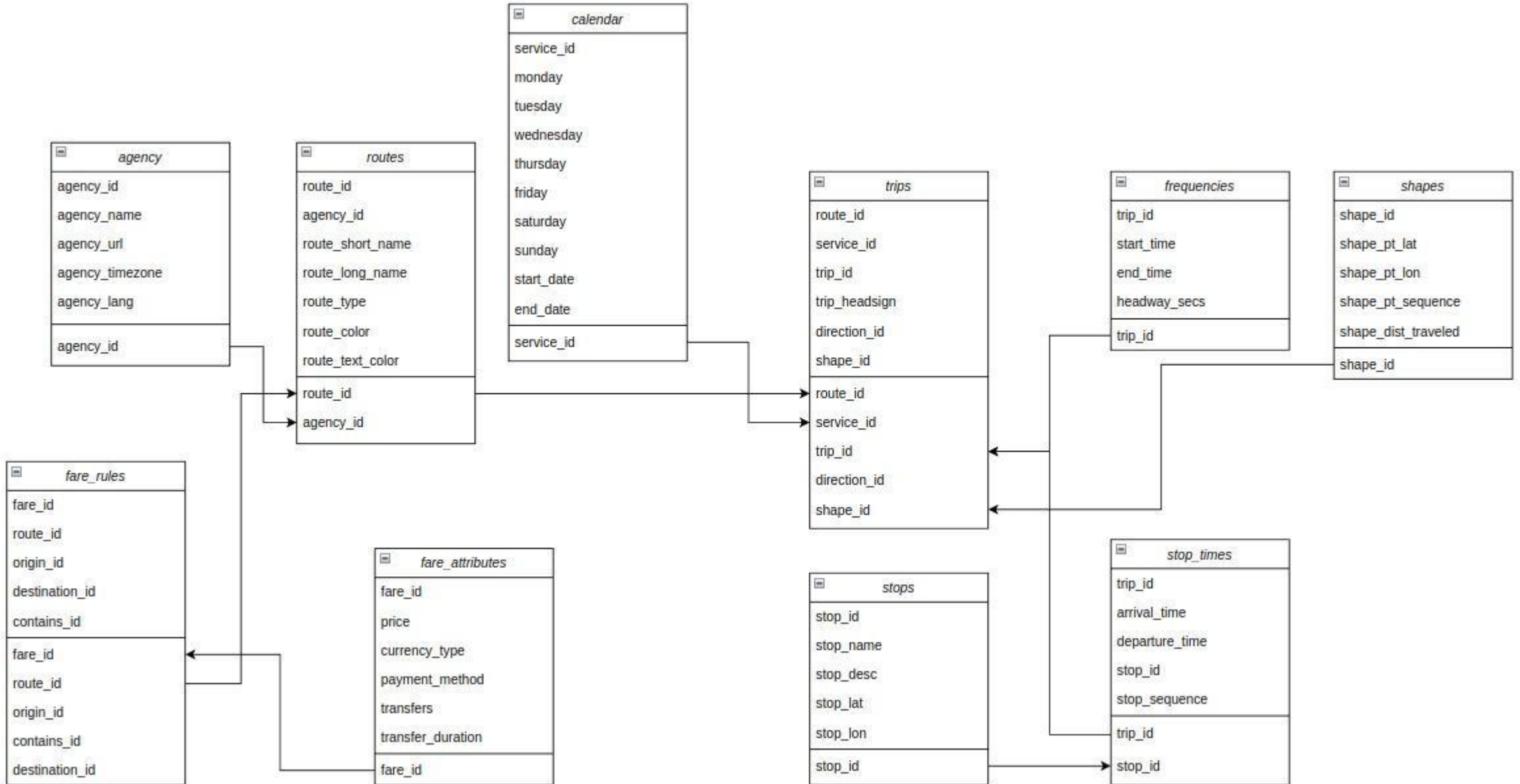
Mario José C M Prado

Wislom Diogo Almeida

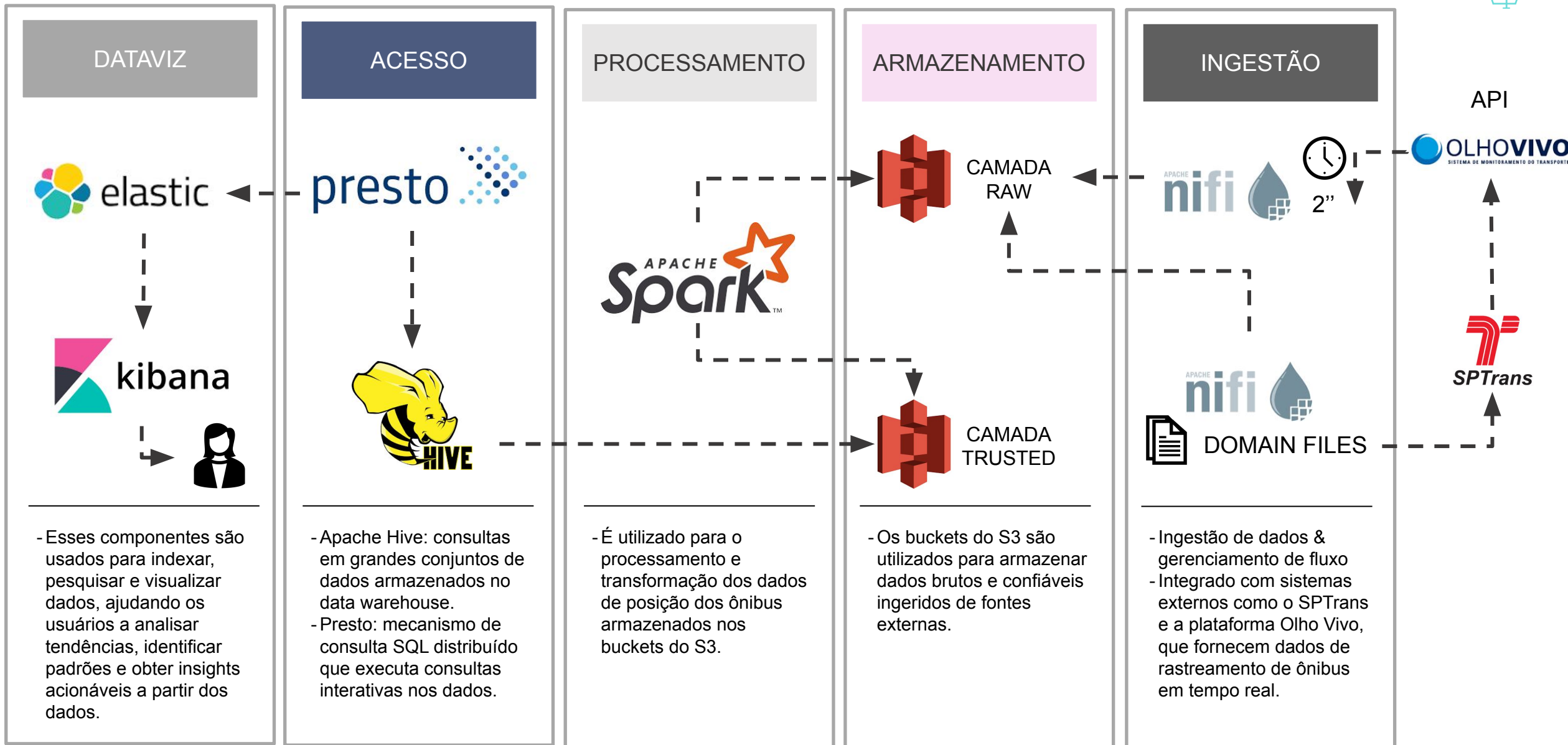
Coordenador:

Prof. Fabio Jardim





3.i Solução | Arquitetura da solução



3.ii Solução | Documentação da aplicação



INGESTÃO E ARMAZENAMENTO

PROCESSAMENTO E ACESSO

DATAVIZ

Credenciais de acesso:

NiFi:

- **URL:** <http://localhost:49090>
- **Porta:** 9090 (mapeada para a porta 49090 no host)

MinIO:

Console: <http://localhost:49001>

Serviço de

Armazenamento: <http://localhost:49000>

Portas:

- 9000 (serviço de armazenamento) mapeada para a porta 49000 no host
- 9001 (console) mapeada para a porta 49001 no host

Credenciais de Acesso:

Usuário: admin / **Senha:** minioadmin



README/GITHUB

1



Macro-etapas

Definir serviços principais (NiFi e MinIO) e criar rede personalizada otmzsp-network

Detalhamento

1. Arquivo docker-compose.yml
2. Comando iniciar docker-compose

Pré-configurações

Volume montado para persistência de dados e logs (ex.: repositórios de banco de dados, arquivos de fluxo, conteúdo, etc.).
Definido fuso horário de SP.
Versões corretas do NiFi e MinIO configuradas nas variáveis de ambiente `#{NIFI_VERSION}` e `#{MINIO_VERSION}` no arquivo `.env`.
Diretório de dados:
`/mini_data/raw` (brutos)
`/mini_data/trusted` (confiáveis)

2



Gerar certificado da API Olho Vivo

1. Obter certificado
2. Importar o certificado para um keystore Java

Configuração de diretório `util/jks` exista antes de executar o comando.

Documentação detalhada:
[Aqui](#)

Dicionário de dados:
[Aqui](#)

3



Configurar o `StandardSSLContextService` no NiFi

1. Acessar a interface web do NiFi
2. Adicionar e configurar o `StandardSSLContextService`
3. Aplicar a config. Reiniciar o serviço se necessário

NiFi configurado para utilizar o `StandardSSLContextService` em seus processadores que se comunicam via HTTPS.

3.ii Solução | Documentação da aplicação



INGESTÃO E
ARMAZENAMENTO

PROCESSAMENTO E
ACESSO

DATAVIZ

Credenciais de acesso:

Apache SPARK:

• **URL:** <http://localhost:49090>

Portas Spark Master:

• 7077

• 8080

Portas Spark Worker: 8081

JUPYTER:

• **Console:** <http://localhost:49001>

• **Serviço de Armazenamento:** <http://localhost:49000>

Portas:

• 4040

• 4041

• 4042

• 4043

 README/GITHUB



Macro-etapas

Gerenciar os recursos e distribuir as tarefas entre os nós do cluster e Executar as tarefas distribuídas

Detalhamento

1. Executar o Spark Master com o comando `deploy.master`
2. Executar o Spark Worker com o comando `deploy.worker`

Pré-configurações



Tratamento dos dados

1. Executar o comando `start notebook`



Tratamento dos dados

1. Executar o comando `start notebook`

3.ii Solução | Documentação da aplicação



INGESTÃO E
ARMAZENAMENTO

PROCESSAMENTO E
ACESSO

DATAVIZ

Credenciais de acesso:

Elastic Search:

• **URL:** <http://localhost:49090>

Portas Spark Master:

• 7077

• 8080

Password: "12345"

Portas Spark Worker: 8081

JUPYTER:

Console: <http://localhost:49001>

Serviço de Armazenamento: <http://localhost:49000>

Portas:

• 4040

• 4041

• 4042

• 4043

 README/GITHUB

1  elastic

2  kibana

3

Macro-etapas	XXXXXX	XXXXXX	XXXXXX
Detalhamento	1. XXXXXX	1. XXXXXX	1. XXXXXX
Pré-configurações			

3.ii Solução | Análise de dados



Camada Raw: dados brutos, sem transformações



Camada Trusted: Dados tratados



Camada Bussiness: Dados enriquecidos

Dados extraídos diretamente das camadas de origem (API Olho Vivo e do GTFS);

Dados referentes à solução de monitoramento near real time dos ônibus em circulação.

Ex.: Latitude e longitude dos veículos em tempo real, horários de chegada previstos,...

Integração dos dados a partir das múltiplas fontes (API e GTFS) e análise de possíveis cruzamentos de dados;

Tratamento dos dados e aplicação de regras de limpeza.

Dados enriquecidos;

Manipulação dos dados para consolidação em métricas e KPIs, definidos para atender às necessidades de negócio;

Base definida para consumo por sistemas de visualização e relatórios analíticos.



Exemplos de KPIs

- Velocidade média do ônibus
- Tempo total das viagens
- Tempo médio em paradas
- Horário estimado de chegada



Melhorias: avaliação upgrade plano pago do mapa para visualização linhas tracejados

Limitação contronada pelo campo keyword