

# Apache Nifi

## Automação e Integração de Dados

Pedro Saito e Milton Salgado

Laboratório de Métodos de Suporte à Tomada de Decisão

Universidade Federal do Rio de Janeiro

2025-07-27

# Sumário

1. Introdução .....	2	4.2 <i>Custom Processors</i> .....	18
1.1 História .....	3	5. Comparações .....	19
1.2 Definição .....	4	5.1 Airflow .....	20
1.3 Funcionalidades .....	5	5.2 Desvantagens .....	22
1.4 Automação de Tarefas ..	6	5.3 Vantagens .....	23
2. Componentes .....	7	5.4 Quando Usar NiFi .....	25
2.1 Terminologia .....	8	5.5 Tchau .....	26
3. Aplicação .....	13		
3.1 Exemplo .....	14		
4. Códigos .....	16		
4.1 Scripts .....	17		

# 1. Introdução

---

*Contexto:* O Apache Nifi foi inicialmente desenvolvido pela **NSA** (*National Security Agency*) em 2006, anteriormente conhecido por “*NiagaraFiles*”, para ser uma ferramenta de **automação de fluxo de dados**.

Em 2014, o projeto foi incorporado pelo **Apache Software Foundation** como um projeto *Open Source*.

Baseado no paradigma **FBP** (*Flow-Based Programming*).

★ 5.4k **Github**

510 **Contribuintes**

 **Java 21**

Ferramenta *low-code* de automação de fluxo de dados em tempo real (*streaming data*).

**Plataforma primária:** Aplicação Web **Nativa (Java 21)** ou **Docker**.

Casos de uso principais:

- Workflows moderados de ETL (*real time* ou *event-driven*).
- Processos interconectados e em constante mudança.
- Migração de dados entre sistemas legados e modernos
- Sincronização de bancos de dados em tempo real

- Gerenciamento de múltiplos usuários com diferentes permissões.
- **Data Lineage**: Monitoramento e rastreabilidade dos dados
- Subir uma instância no *Tailscale* para uso interno.
- **Tolerância a falhas** com retry automático e dead-letter queues.
- **Escalabilidade horizontal** com clustering gerenciado pelo ZooKeeper.

<https://localhost:8443/nifi/>

O Apache NiFi automatiza diversas tarefas críticas:

### **Cibersegurança**

- Coleta de logs de segurança
- Correlação de eventos suspeitos
- Análise de tráfego real

### **Observabilidade**

- Métricas agregadas
- Logs em tempo real
- Dashboards automáticos

### **Compliance e Auditoria**

- Rastreamento de dados sensíveis
- Relatórios automáticos
- Políticas de retenção
- Conformidade LGPD/GDPR

### **Detecção de Anomalias**

- Análise de padrões em transações
- Comportamentos anômalos
- Integração com ML models
- Alertas em tempo real

## 2. Componentes

---



- **Dataflow Manager** : Usuário Nifi com permissões para adicionar, remover ou modificar os componentes do *Nifi Flow*.
- **FlowFile**: Unidade básica de dados. Composta pelos *Attributes* (metadados) e o *Content* (conteúdo). Principais metadados:
  - **UUID**: Identificador único universal de 128 bits (RFC 4122 do IETF).
  - **Filename**: Nome do arquivo/dado.
  - **Path**: Diretório estruturado para armazenamento no disco ou serviço externo.

- **Processor:** Componente lógico básico responsável pelas tarefas:
  - Exemplos de *Processors*: GetFile, PutFile, ExecuteSQL, InvokeHTTP, ReplaceText, UpdateAttribute, LogAttribute, RouteOnAttribute.

### Ingestão

### Transformação

### Roteamento

- **Relacionamentos:** Cada processador pode ter zero ou mais relacionamentos definidos. Saída lógica do processador (**Success** ou **Failure**). Após o término do processamento, o FlowFile deve ser roteado/transferido para outro componente escolhido pelo DataFlow Manager.

- **Canvas:** Interface visual primária para o *design*, gerenciamento e monitoramento dos *dataflows*.
- **Conexão:** Componente “físico” que interliga quaisquer dois elementos do *dataflow* (sobretudo *Processors*). Ligação física entre 2 componentes do fluxo de dados.
- **Bulletin:** Os componentes podem reportar *bulletins*, isto é, equivalente a mensagens de *logs* com tempo e nível de severidade (*Debug*, *Info*, *Warning* e *Error*).

## 2.1 Terminologia

- **Controller Service:** Componente que provê funcionalidades e configurações para outros elementos (sobretudo *Processors*).
  - **Exemplo:** DBCPConnectionPool fornece serviço de pool de conexões para um ou mais processadores. Propriedades obrigatórias:
    - **Database Connection URL:** jdbc:sqlserver://ip;port;db
    - **Database Driver Class Name:** com.Microsoft.sqlserver.jdbc.SQL
    - **Driver's Path:** ~/mssql-drivers/ptb/jars
    - **Database User & Password:** sa \*\*\*\*

## 2.1 Terminologia

- **Funnel**: Componente que permite combinar dados de diferentes conexões em uma só.
- **Process Group**: Agrupamento lógico de *Processors* para facilitar o gerenciamento do *dataflow*.
- **Remote Process Group**: Transferência de dados (*FlowFiles*) para uma instância remota do *Nifi*.
- **Reporting Tasks**: Tasks executadas no *background* que providenciam relatórios de utilização e monitoramento da instância *Nifi*.
  - Exemplos: `MonitorDiskUsage`, `MonitorMemory`.

### 3. Aplicação

---

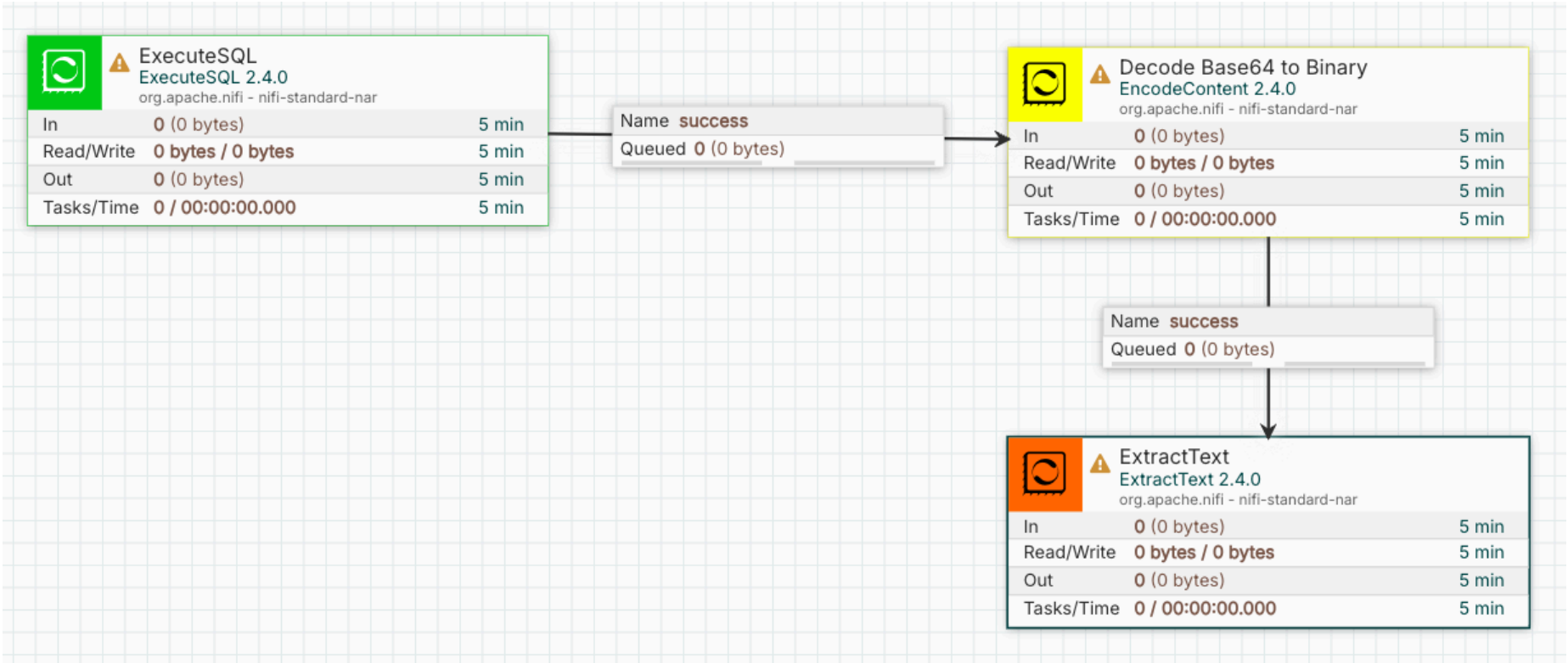


Figura 1: Fluxo de consulta do base64, decodificação e extração de texto.

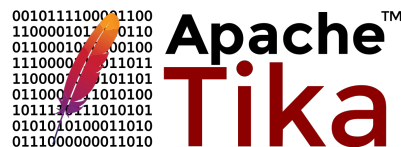
## 3.1 Exemplo

## 3. Aplicações

```
SELECT TOP 10 base64encode  
FROM [pgmconsultaprocesso].[dbo].[documentos]  
WHERE  
    sucesso = 1  
AND ext = '.pdf'
```



Base64 → Binary





## 4. Códigos

---

## 4.1 Scripts

Dois *Processors* primordiais para execução de comandos externos:

- **ExecuteScript**: Executa scripts em linguagens como Python, Groovy, JavaScript diretamente no processador.
- **ExecuteStreamCommand**: Passa o conteúdo do *FlowFile* como input, executa uma série de comandos arbitrários e retorna o resultado da transformação.

```
import sys, json
from datetime import datetime
data = json.load(sys.stdin)
data['modificada_em'] = datetime.utcnow().isoformat()
json.dump(data, sys.stdout)
```

python

Listagem 1: *ExecuteStreamCommand* para adicionar metadado *datetime*.

## 4.2 Custom Processors

Na ausência de *Processors* que cumprem a tarefa desejada, pode-se definir seus próprios *Processors* em **Java** e importá-los no *dataflow*.

```
public class ExampleProcessor extends AbstractProcessor {  
    public static final Relationship REL_FAILURE = new Relationship  
        .builder()  
        .description("Failed processing")  
        .name("failure")  
        .build();  
  
    public static final Relationship REL_SUCCESS = new Relationship  
        .builder()  
        .description("Succeed processing")  
        .name("success")  
        .build();  
  
    @Override  
    public Set<Relationship> getRelationships() {  
        return new HashSet<>() {{  
            add(REL_FAILURE);  
            add(REL_SUCCESS);  
        }};  
    }  
}
```



.Nar (Nifi Archive)

## 5. Comparações

---

## 5.1 Airflow

Atributos	Apache Nifi	Apache Airflow
Interface	GUI <i>Drag-and-Drop</i>	Código Python (DAGs)
Escalabilidade	Cluster horizontal	Executors
Periodicidade	<i>Event-Driven</i>	Time-based schedules
Utilização	Interface Gráfica	Python <i>script</i>
Monitoramento	<i>Provenance events</i>	Logging + interface
Concorrência	Back-pressure, priorização, load- balancing entre nós	Execução paralela via Executors (Celery, Kubernetes, Local)

## 5.1 Airflow

Atributos	Apache Nifi	Apache Airflow
Data Handling	<i>FlowFiles</i> com conteúdo + atributos; streaming ou micro-batch	Task outputs persistidos (XComs, arquivos, databases); principalmente batch
Failure Handling	Retry automático, dead-letter queues, rastreamento de proveniência	Políticas de retry por task, alertas SLA, intervenção manual

## 5.2 Desvantagens

- **Reprodutibilidade:** Dificuldade para copiar um *workflow* de uma instância do *Nifi* para outra. É preciso recriar cada componente manualmente.
- **Dinamidade:** Dificuldade na modelagem de *dataflows* que se ajustam constantemente (mais complexos).
  - Exemplo: SQL Dinâmico.
- **Depuração:** Herda características do paradigma **FBP** que concebem os *Processadores* como caixas pretas. Em algum momento será preciso olhar por de trás dos panos.

## 5.3 Vantagens

- **Facilidade no Uso:** Interface simples para construção de fluxos “comuns”, desde que não exigem funcionalidades externas.
- **Escalabilidade Horizontal:** *Deploy* de um Cluster de servidores do *Nifi* geridos pelo *ZooKeeper*.
- **Variedade das Fontes de Dados:** Variedade razoável de *Processors* para ingestão de dados.
  - Exemplos: Sistemas de Arquivo, Banco de Dados, APIs e Protocolos de Transferência, E-mail, Serviços de Streaming como Kafka.
- **Processamento em Tempo Real:** Event-driven com baixa latência e alta performance.



## 5.3 Vantagens

- **Data Lineage Completo:** Rastreabilidade total dos dados para compliance e auditoria.
- **Tolerância a Falhas:** Recuperação automática e gerenciamento de erros robusto.
- **Governança de Dados:** Controle granular de acesso, segurança e políticas de dados.

## 5.4 Quando Usar NiFi

### ✓ Ideal para:

Fluxos complexos e dinâmicos

Rastreabilidade completa

Processamento em tempo real

Integração de múltiplas fontes

Equipes não-técnicas

Workflows *event-driven*

### ✗ Não recomendado para:

Transformações complexas

Processamento batch pesado

Computação intensiva

Modelagem com DAGs é preferível





Figura 2: *Niagara falls.*