

Oficina HUE

- | | | | |
|-------|------------------------|-------|----------------|
| I. | Conceitos e Definições | X. | Exercícios 1 |
| II. | Componentes | XI. | Workflow Oozie |
| III. | Navegação | XII. | Hands On 2 |
| IV. | Job Browser | XIII. | Exercícios 2 |
| V. | Database Browser | | |
| VI. | Query Editor | | |
| VII. | Impala x Hive | | |
| VIII. | Boas Práticas Impala | | |
| IX. | Hands On 1 | | |

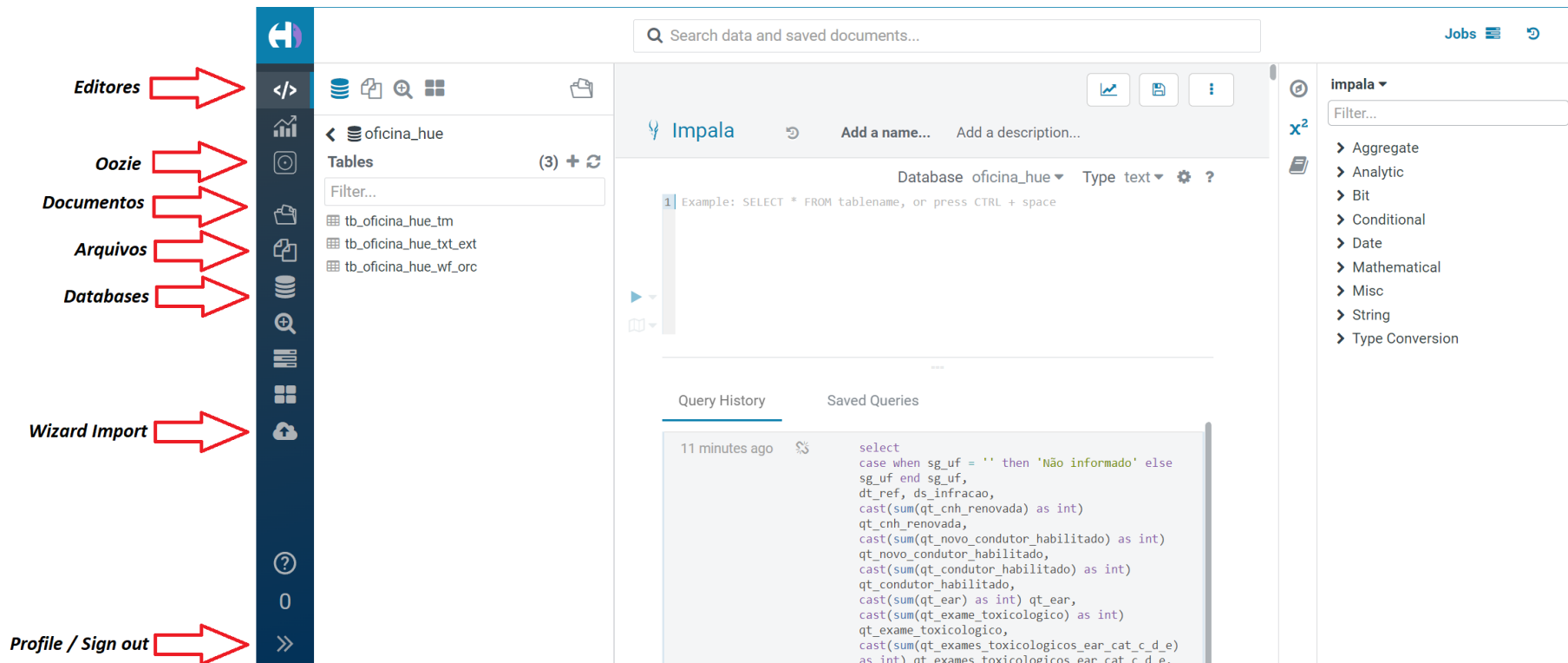
Hue é um assistente de SQL para bancos de dados e DW, desenvolvido com o objetivo de acessar via browser os dados / arquivos armazenados em um cluster HDFS.

- HDFS Browser
- Job Browser
- Database Browser
- Query Editor
- Oozie Workflows



- **HDFS Browser** - Interface web cujo principal objetivo é permitir à interação do usuário com os arquivos armazenados no cluster Hadoop, como um sistema de arquivos.
- **Job Browser** - Ferramenta que nos permite acompanhar tarefas executadas no cluster Hadoop.
- **Database Browser** - Interface que nos permite visualizar de forma gráfica os “databases” existentes no HDFS, e seus componentes.
- **Query Editor** - Editor de consultas HQL que nos permite utilizar diversas ferramentas de consulta (Hive, Impala, etc) aos dados do cluster Hadoop de forma mais amigável.
- **Oozie Workflows** - Aplicação incorporada ao HUE que nos permite a criação de fluxos de trabalho e agendamento dos mesmos para execução no cluster Hadoop.

Barra de atalhos da nova versão do HUE

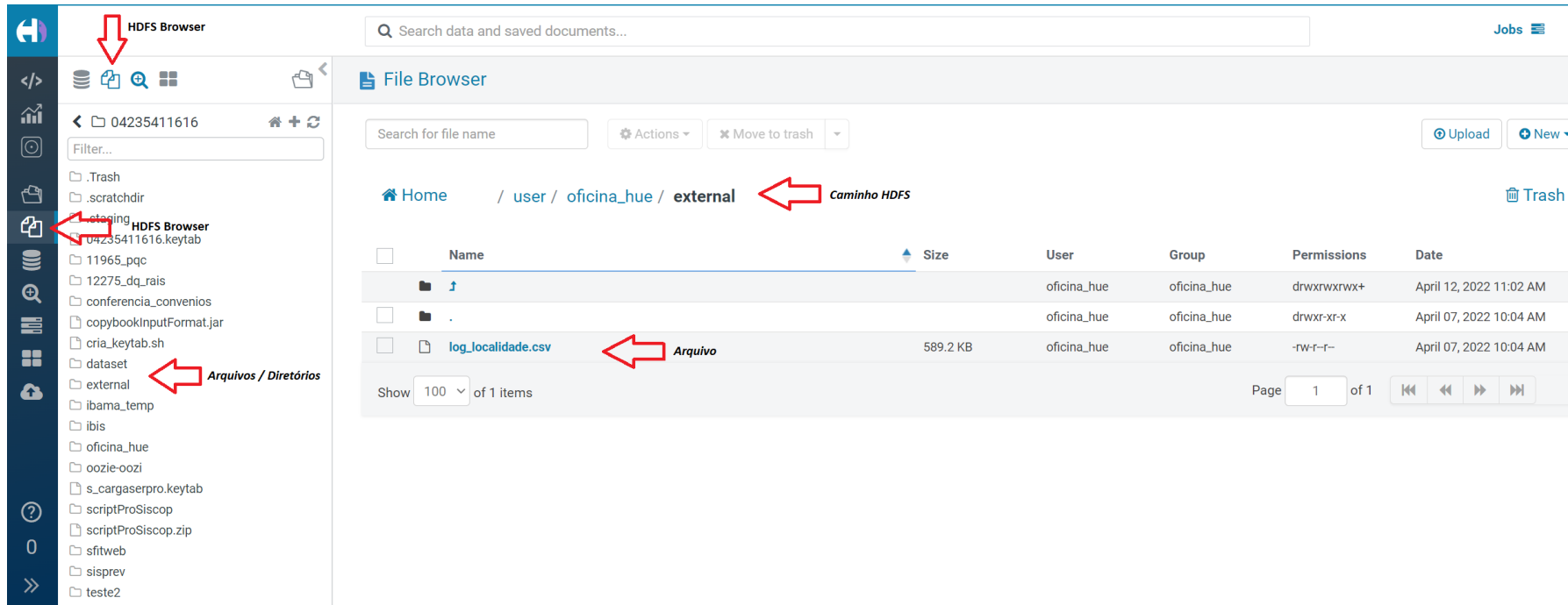


The image shows the HUE interface with a sidebar on the left containing various icons. Red arrows point from labels to specific icons in the sidebar:

- Editores** points to the code editor icon (a blue square with a white </> symbol).
- Oozie** points to the Oozie icon (a blue square with a white Oozie logo).
- Documentos** points to the document icon (a blue square with a white document icon).
- Arquivos** points to the file icon (a blue square with a white file icon).
- Databases** points to the database icon (a blue square with a white database icon).
- Wizard Import** points to the wizard icon (a blue square with a white wizard icon).
- Profile / Sign out** points to the profile icon (a blue square with a white profile icon).

The main interface shows the HUE logo, a search bar, and a list of tables under the 'oficina_hue' database. The 'Query History' tab is active, showing a query executed 11 minutes ago. The 'Saved Queries' tab is also visible. The right sidebar shows a list of functions under the 'impala' dropdown, including Aggregate, Analytic, Bit, Conditional, Date, Mathematical, Misc, String, and Type Conversion.

HDFS Browser, que permite visualizar, apagar, incluir, criar, compartilhar e movimentar arquivos e diretórios no HDFS através da interface web.



The screenshot displays the HDFS Browser web interface. On the left sidebar, the 'HDFS Browser' tab is selected, indicated by a red arrow. Below it, the 'File Browser' section shows a list of files and directories. A red arrow points to the 'external' directory, which is highlighted. Another red arrow points to the 'log_localidade.csv' file, which is also highlighted. The main area shows the contents of the 'external' directory, including a table of files and directories. The table has columns for Name, Size, User, Group, Permissions, and Date. The file 'log_localidade.csv' is listed with a size of 589.2 KB and permissions of -rw-r--r--. A red arrow points to the file name in the table, with the label 'Arquivo' next to it. The breadcrumb path at the top of the main area is 'Home / user / oficina_hue / external', with a red arrow pointing to 'external' and the label 'Caminho HDFS' next to it. The bottom of the interface shows a pagination bar with 'Page 1 of 1' and navigation buttons.

HDFS Browser

Search data and saved documents...

File Browser

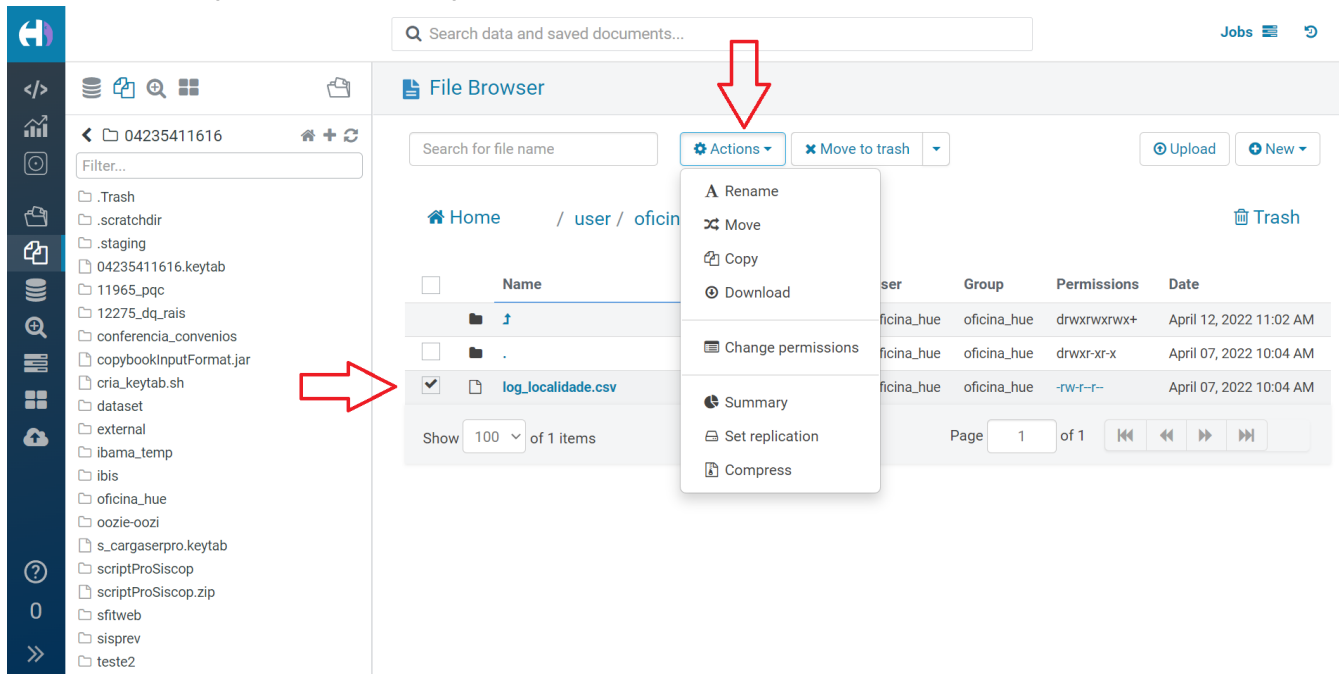
Search for file name Actions Move to trash Upload New

Home / user / oficina_hue / external Caminho HDFS

Name	Size	User	Group	Permissions	Date
log_localidade.csv	589.2 KB	oficina_hue	oficina_hue	-rw-r--r--	April 07, 2022 10:04 AM

Show 100 of 1 items Page 1 of 1

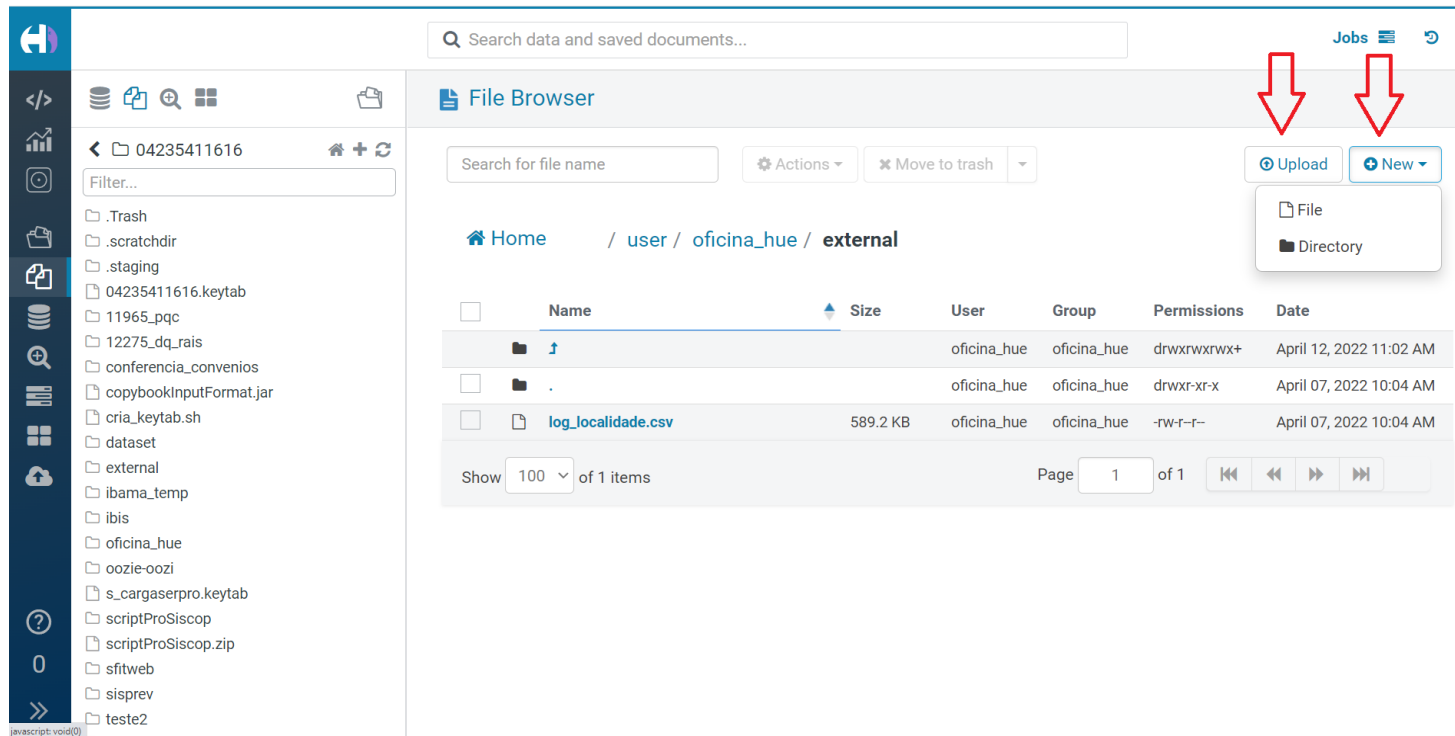
Por esta janela o usuário poderá, dependendo das permissões nos arquivos, realizar as ações disponíveis na dropdown “Ações” (renomear, mover, copiar, baixar, alterar permissões, alterar número de réplicas e compactar) para os arquivos selecionados através das *checkbox*, localizadas à esquerda de cada arquivo / diretório.



The screenshot displays the HDFS Browser interface. On the left, a sidebar shows a tree view of the file system. The main area is titled 'File Browser' and shows a search bar, a filter input, and a list of files. The file 'log_localidade.csv' is selected, and its context menu is open, showing options like Rename, Move, Copy, Download, Change permissions, Summary, Set replication, and Compress. A table of files is visible in the background.

Name	ser	Group	Permissions	Date
oficina_hue	oficina_hue	drwxrwxrwx+	April 12, 2022 11:02 AM	
oficina_hue	oficina_hue	drwxr-xr-x	April 07, 2022 10:04 AM	
oficina_hue	oficina_hue	-rw-r--r--	April 07, 2022 10:04 AM	

Poderá também criar diretórios e arquivos através da dropbox “Novo” e fazer upload de arquivos no diretório corrente pelo botão carregar.



The screenshot displays the HDFS Browser interface. On the left is a sidebar with various icons for navigation. The main area shows the 'File Browser' view for the path `/user/oficina_hue/external`. A search bar is at the top, and a table lists files and directories. Two red arrows point to the 'Upload' and 'New' buttons in the top right corner. The 'New' button has a dropdown menu with options for 'File' and 'Directory'.

Search data and saved documents...

Jobs

File Browser

Search for file name

Actions

Move to trash

Upload

New

File

Directory

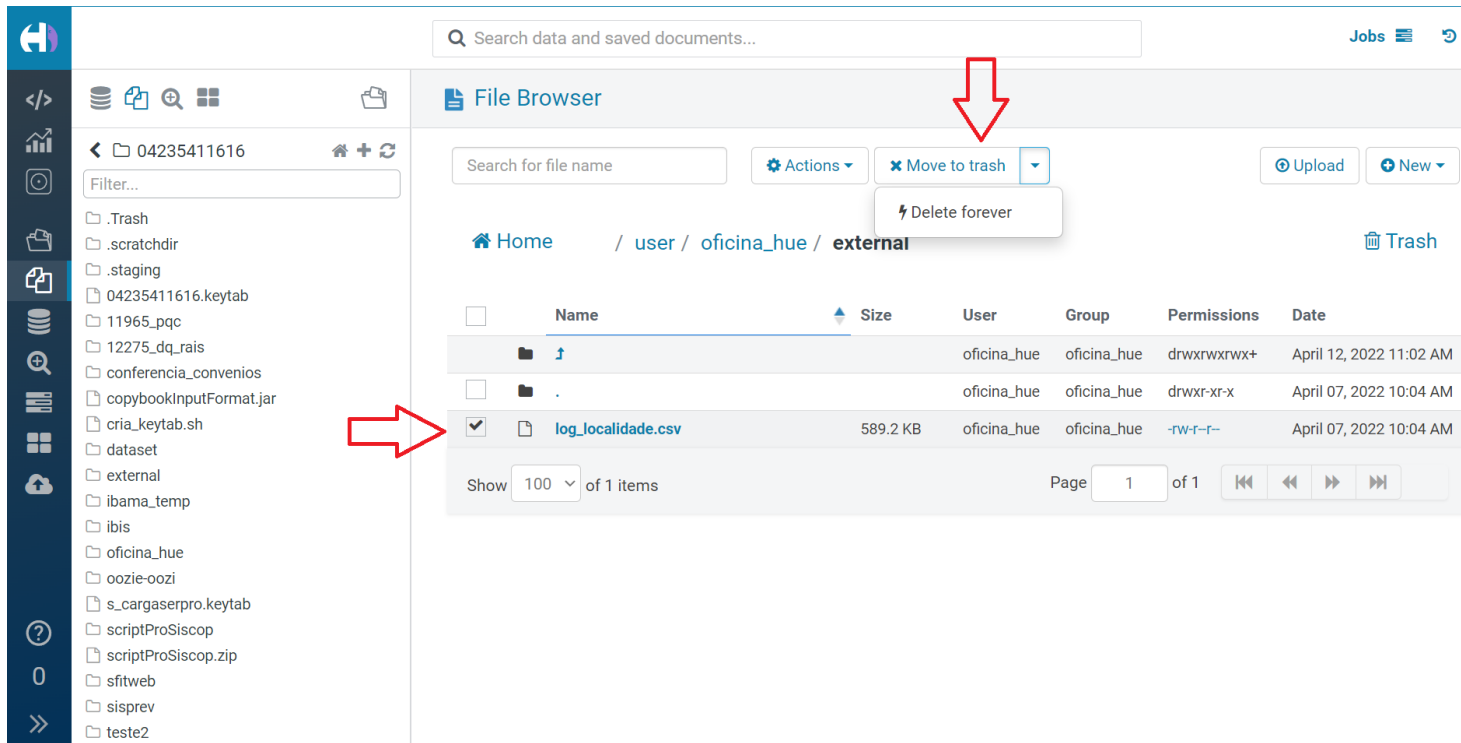
Home / user / oficina_hue / external

	Name	Size	User	Group	Permissions	Date
	↑		oficina_hue	oficina_hue	drwxrwxrwx+	April 12, 2022 11:02 AM
	.		oficina_hue	oficina_hue	drwxr-xr-x	April 07, 2022 10:04 AM
	log_localidade.csv	589.2 KB	oficina_hue	oficina_hue	-rw-r--	April 07, 2022 10:04 AM

Show 100 of 1 items

Page 1 of 1

Para apagar arquivos ou diretórios basta selecionar o referente pela checkbox e utilizar a dropbox “Move to trash”.

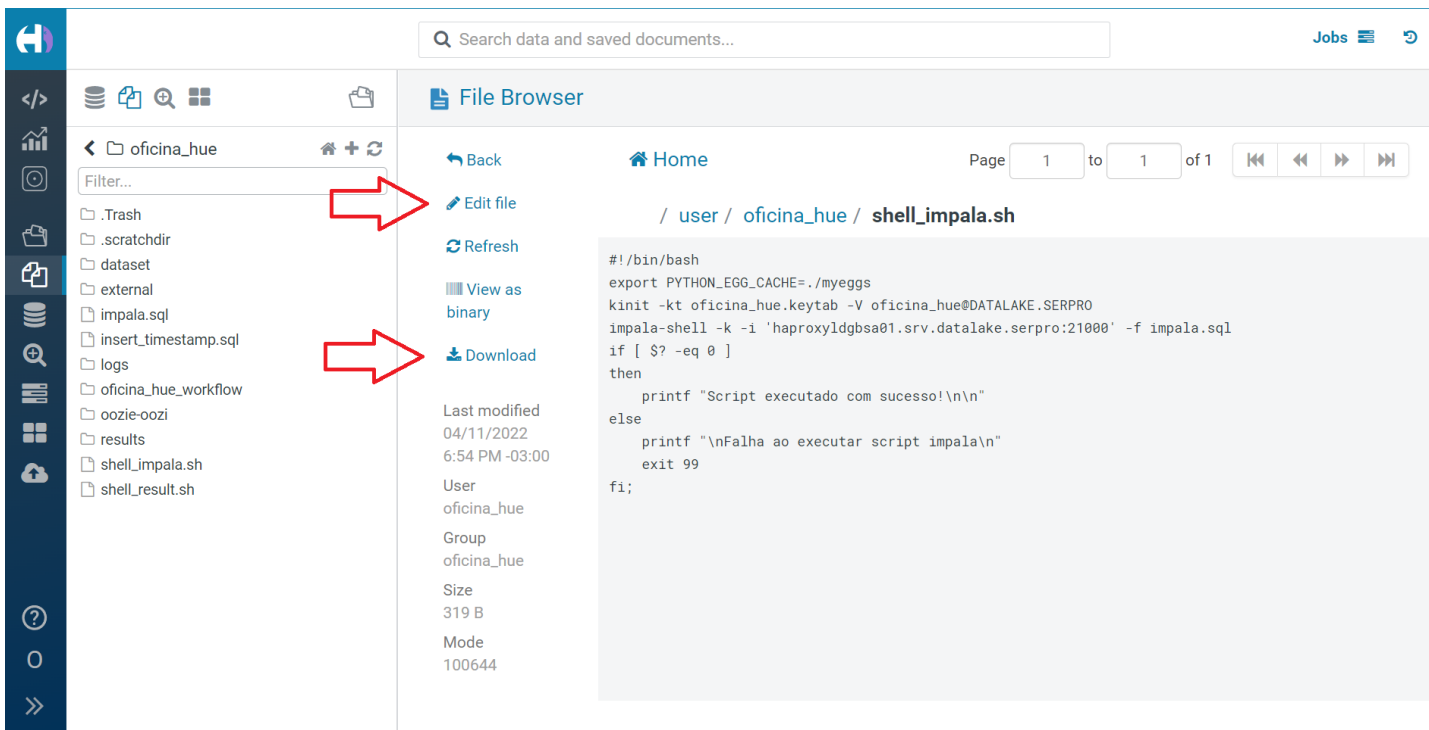


The screenshot displays the HDFS Browser interface. On the left is a sidebar with a navigation menu. The main area is titled 'File Browser' and shows a search bar, a breadcrumb path 'Home / user / oficina_hue / external', and a table of files. The file 'log_localidade.csv' is selected, indicated by a checked checkbox and a red arrow. A red arrow also points to the 'Move to trash' button in the 'Actions' dropdown menu. The table lists the following files:

	Name	Size	User	Group	Permissions	Date
<input type="checkbox"/>	↑		oficina_hue	oficina_hue	drwxrwxrwx+	April 12, 2022 11:02 AM
<input type="checkbox"/>	.		oficina_hue	oficina_hue	drwxr-xr-x	April 07, 2022 10:04 AM
<input checked="" type="checkbox"/>	log_localidade.csv	589.2 KB	oficina_hue	oficina_hue	-rw-r--r--	April 07, 2022 10:04 AM

At the bottom of the table, it shows 'Show 100 of 1 items' and 'Page 1 of 1'.

O HDFS Browser ainda permite que realizemos alterações nos arquivos localizados no HDFS diretamente no Browser e download do mesmo, caso prefira editar em sua máquina.



The screenshot displays the HDFS Browser interface. On the left, a sidebar contains navigation icons and a list of files and directories under the path `< oficina_hue`. The files listed are `.Trash`, `.scratchdir`, `dataset`, `external`, `impala.sql`, `insert_timestamp.sql`, `logs`, `oficina_hue_workflow`, `oozie-oozi`, `results`, `shell_impala.sh`, and `shell_result.sh`. Two red arrows point from the `Filter...` input field and the `shell_impala.sh` file to the right-hand pane. The right-hand pane, titled "File Browser", shows the selected file `shell_impala.sh` with its content displayed in a code editor. The file's metadata is shown below the code: Last modified 04/11/2022 6:54 PM -03:00, User oficina_hue, Group oficina_hue, Size 319 B, and Mode 100644. The code editor contains a shell script that sets environment variables, runs a kinit command, and executes an impala-shell command.

Search data and saved documents...

Jobs

File Browser

Home

Page 1 to 1 of 1

Back

Edit file

Refresh

View as binary

Download

Last modified 04/11/2022 6:54 PM -03:00

User oficina_hue

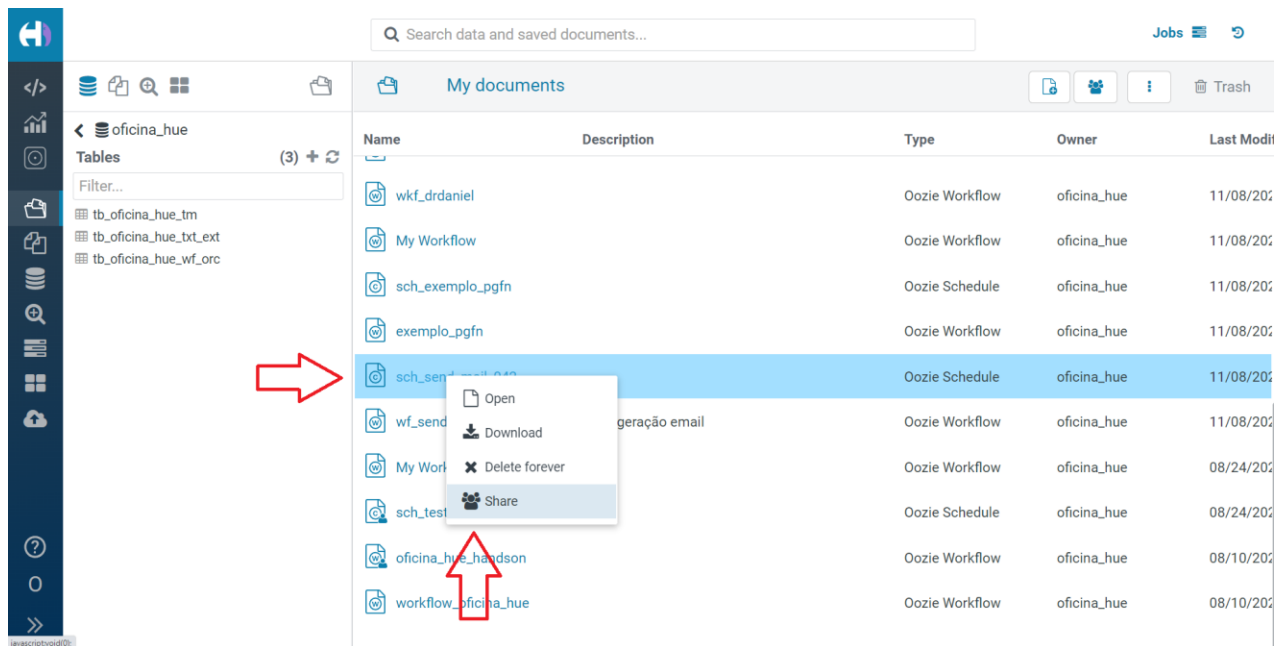
Group oficina_hue

Size 319 B

Mode 100644

```
#!/bin/bash
export PYTHON_EGG_CACHE=./myeggs
kinit -kt oficina_hue.keytab -V oficina_hue@DATA Lake.SERPRO
impala-shell -k -i 'haproxyldb01.srv.data lake.serpro:21000' -f impala.sql
if [ $? -eq 0 ]
then
    printf "Script executado com sucesso!\n\n"
else
    printf "\nFalha ao executar script impala\n"
    exit 99
fi;
```

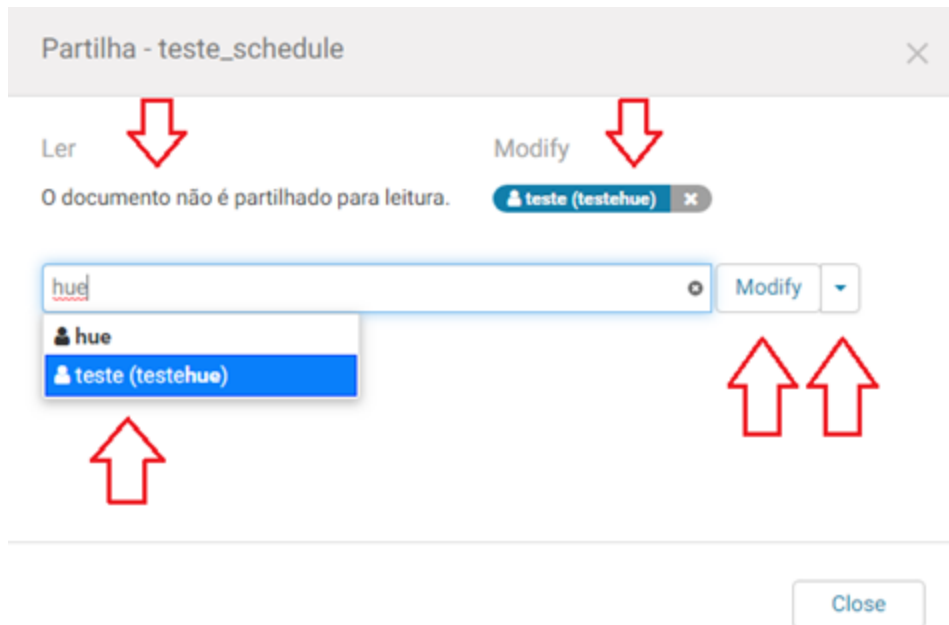
Podemos também realizar o compartilhamento de queries, workflows, schedules e outros documentos, com outros usuários do cluster. Acessando o menu principal, opção *Documents*, conseguimos visualizar os arquivos cuja edição foram feitas através das ferramentas do HUE (queries salvas, workflows oozie, schedules, diretórios, etc).



The screenshot displays the HDFS Browser interface. On the left, a sidebar shows the navigation menu with icons for Home, Tables, Documents, and other features. The main area is titled 'My documents' and contains a table of saved documents. A red arrow points to the 'Share' option in the context menu for the document 'sch_send_email'.

Name	Description	Type	Owner	Last Modified
wkf_drdaniel		Oozie Workflow	oficina_hue	11/08/2021
My Workflow		Oozie Workflow	oficina_hue	11/08/2021
sch_exemplo_pgfn		Oozie Schedule	oficina_hue	11/08/2021
exemplo_pgfn		Oozie Workflow	oficina_hue	11/08/2021
sch_send_email		Oozie Schedule	oficina_hue	11/08/2021
wf_send_email	geração email	Oozie Workflow	oficina_hue	11/08/2021
My Workflow		Oozie Workflow	oficina_hue	08/24/2021
sch_test		Oozie Schedule	oficina_hue	08/24/2021
oficina_hue_hadoop		Oozie Workflow	oficina_hue	08/10/2021
workflow_oficina_hue		Oozie Workflow	oficina_hue	08/10/2021

Na janela Partilha, o usuário deve digitar o nome do usuário (ou parte do nome) com quem deseja compartilhar o documento, e selecioná-lo quando mostrado. Para que seja efetivado o compartilhamento deve-se selecionar à direita qual tipo de permissão e clicar na descrição da permissão, para que a mesma apareça na parte superior da janela e seja assim efetivada.



Partilha - teste_schedule

Ler

O documento não é compartilhado para leitura.

Modify

teste (testehue)

hue

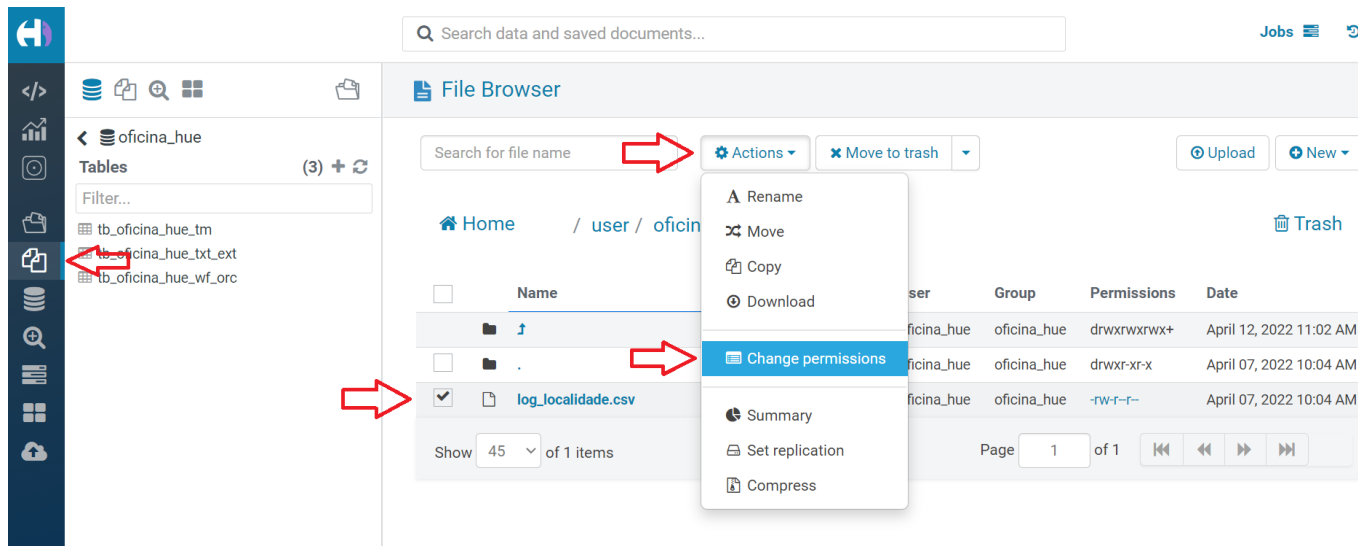
teste (testehue)

Modify

Close

O compartilhamento dos Documentos, conforme mostramos, pode ser realizado para usuários específicos, já o compartilhamento de arquivos dentro do HDFS só se realiza a nível de grupos ou abertura total do arquivo/diretório.

Para o compartilhamento de arquivos no HDFS devemos acessar o menu principal, opção *Files*, que nos levará para o diretório home do usuário. Selecionamos o arquivo a ser compartilhado através do checkbox à esquerda do arquivo e em seguida acessamos a dropdown *Ações* (*Actions*) selecionando a opção *Alterar Permissões* (*Change permissions*).



The screenshot displays the HDFS Browser interface. On the left, a sidebar shows a tree view with 'oficina_hue' selected, containing 'Tables' and 'Filter...'. The main area shows the 'File Browser' for the path '/ user / oficina_hue'. A file 'log_localidade.csv' is selected, and the 'Actions' dropdown menu is open, highlighting 'Change permissions'. Red arrows indicate the path: from the sidebar to the file, then to the 'Actions' menu, and finally to 'Change permissions'.

Name	ser	Group	Permissions	Date
...	oficina_hue	oficina_hue	drwxrwxrwx+	April 12, 2022 11:02 AM
...	oficina_hue	oficina_hue	drwxr-xr-x	April 07, 2022 10:04 AM
...	oficina_hue	oficina_hue	-rwxr-xr-x	April 07, 2022 10:04 AM

Por default, o usuário proprietário possui todas as permissões no arquivo, e os demais apenas permissão de leitura. Lembramos que mesmo que para leitura, os demais usuários ou o grupo do usuário corrente devem possuir permissão no diretório em que o arquivo se encontra, assim como nos demais diretórios caso ele não esteja no diretório raiz do usuário.

Para garantirmos que a permissão seja dada também aos diretórios de nível inferior ao diretório do arquivo, basta selecionar a opção Recursivo.

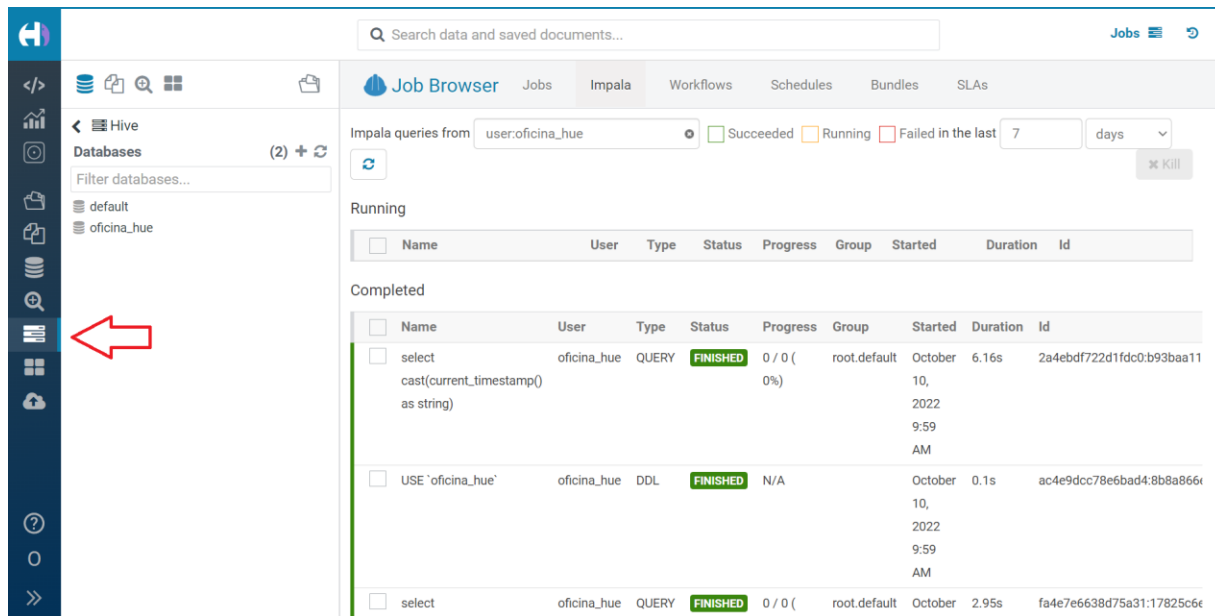
Change Permissions

	User	Group	Other
Read	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Write	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Execute	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Sticky			<input type="checkbox"/>
Recursive			<input type="checkbox"/>

Cancel

Submit

Utilizar o Hadoop implica na criação de muitos jobs e é bastante comum os desenvolvedores terem a necessidade de saber qual job está atualmente em execução no cluster do Hadoop ou qual o resultado de jobs finalizados (se houve sucesso ou falha na execução).



Search data and saved documents...

Jobs

Job Browser Jobs Impala Workflows Schedules Bundles SLAs

Impala queries from user:oficina_hue

Succeeded Running Failed in the last 7 days

Kill

Running

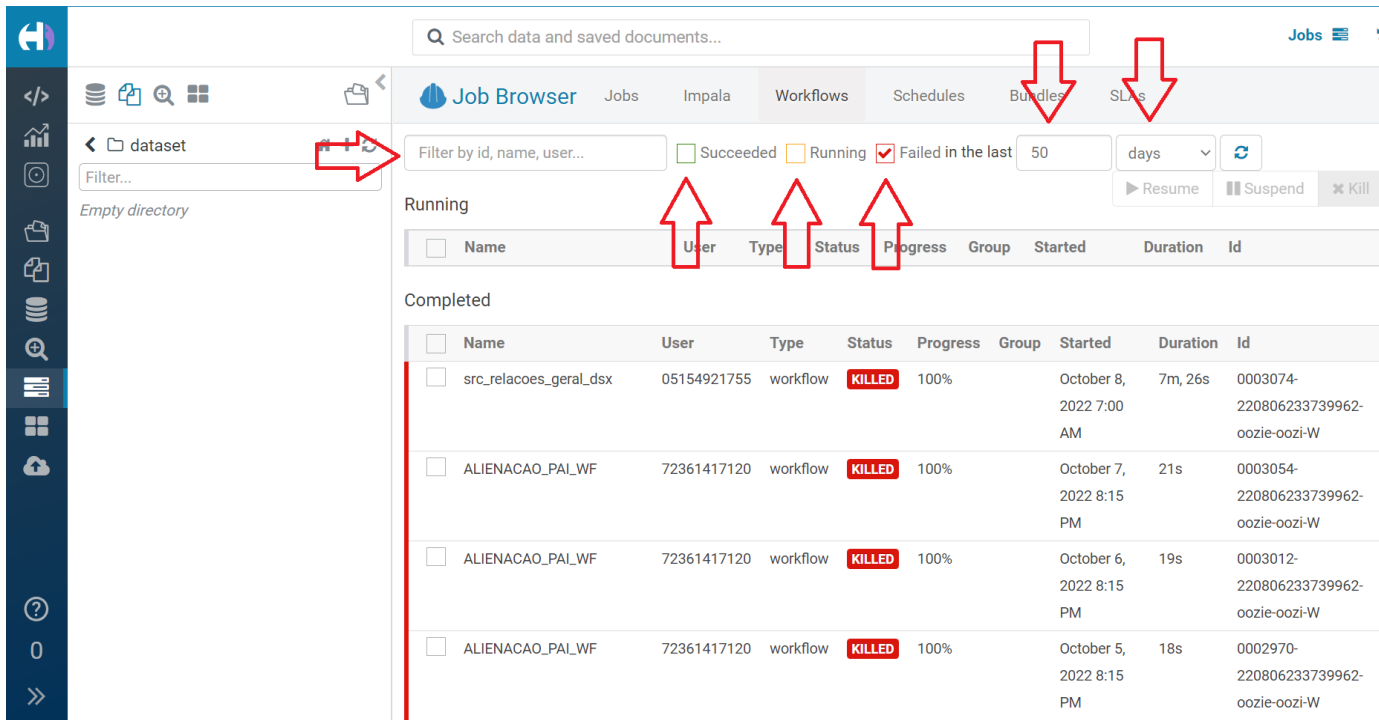
Name	User	Type	Status	Progress	Group	Started	Duration	Id
------	------	------	--------	----------	-------	---------	----------	----

Completed

Name	User	Type	Status	Progress	Group	Started	Duration	Id
select cast(current_timestamp() as string)	oficina_hue	QUERY	FINISHED	0 / 0 (0%)	root.default	October 10, 2022 9:59 AM	6.16s	2a4ebdf722d1fdc0:b93baa11
USE `oficina_hue`	oficina_hue	DDL	FINISHED	N/A		October 10, 2022 9:59 AM	0.1s	ac4e9dcc78e6bad4:8b8a866e
select	oficina_hue	QUERY	FINISHED	0 / 0 (0%)	root.default	October 10, 2022 9:59 AM	2.95s	fa4e7e6638d75a31:17825c6e

Por meio do navegador de jobs, podemos acessar todas as informações relacionadas ao job. Para isso, existe a opção no Hue que permite listar jobs e seus status.

Podemos ainda filtrar os jobs por status, busca textual e período em que foi executado.



The screenshot shows the Job Browser interface with the following elements:

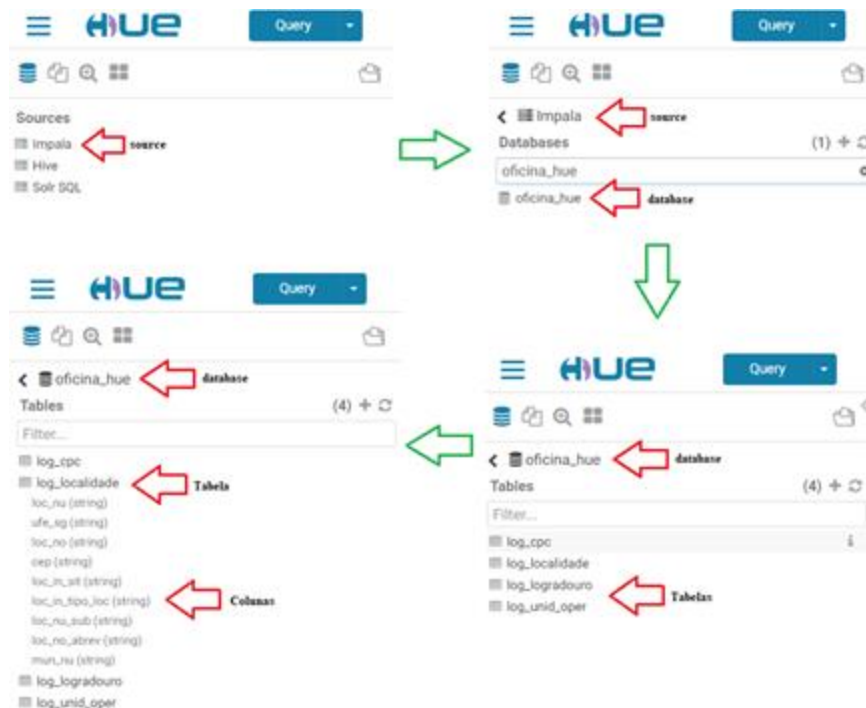
- Search Bar:** "Search data and saved documents..."
- Navigation Tabs:** Job Browser, Jobs, Impala, Workflows, Schedules, Bundles, SLAs.
- Left Sidebar:** Contains icons for various data management actions.
- Dataset View:** Shows a "dataset" directory with a "Filter..." input and "Empty directory" text.
- Filter Section:** Includes a "Filter by id, name, user..." input, status filters (Succeeded, Running, Failed in the last 50 days), and action buttons (Resume, Suspend, Kill).
- Job List:** Divided into "Running" and "Completed" sections. The "Completed" section shows a table of jobs with columns: Name, User, Type, Status, Progress, Group, Started, Duration, and Id.

Red arrows highlight the following features:

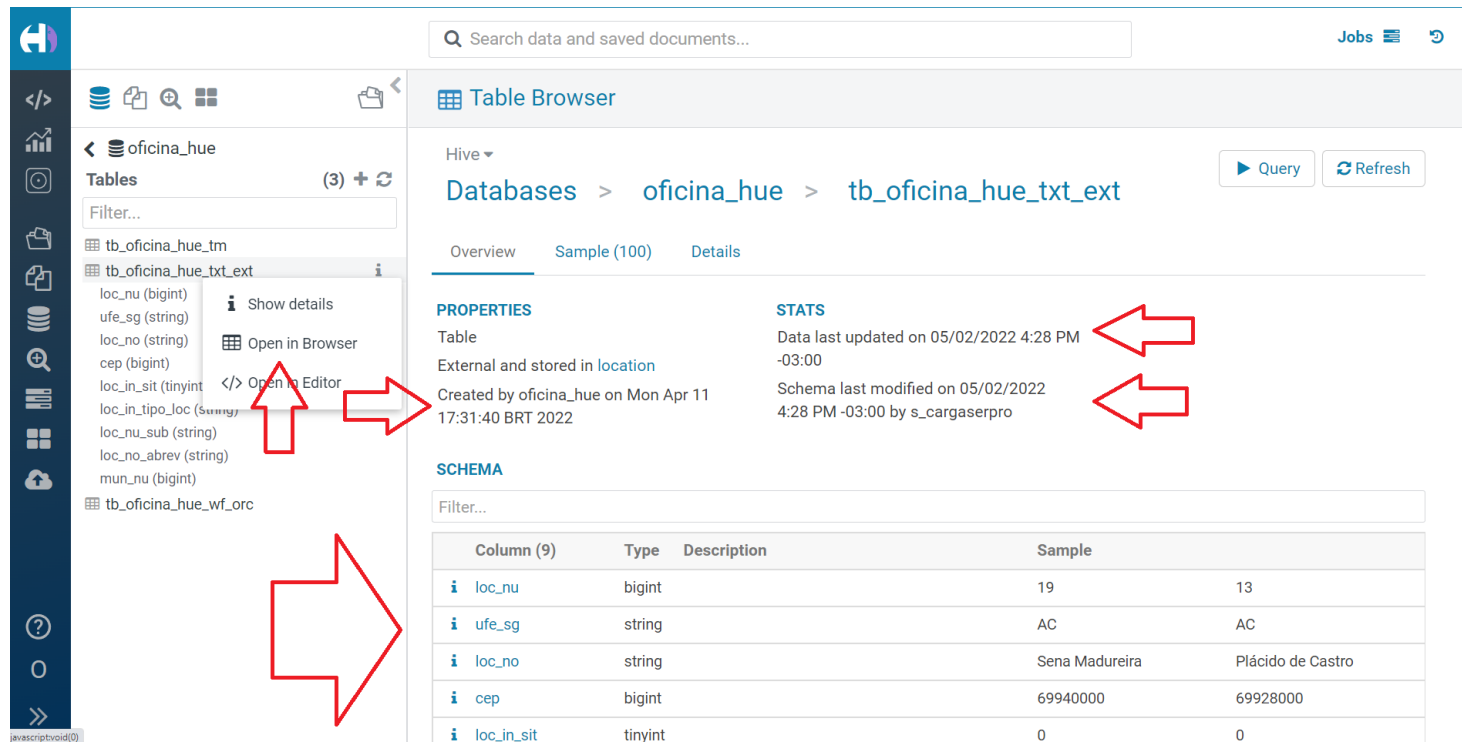
- The "Filter by id, name, user..." input field.
- The "Failed in the last 50 days" status filter.
- The "User", "Type", and "Status" columns in the job list header.

Name	User	Type	Status	Progress	Group	Started	Duration	Id
src_relacoes_geral_dsx	05154921755	workflow	KILLED	100%		October 8, 2022 7:00 AM	7m, 26s	0003074-220806233739962-oozie-oozi-W
ALIENACAO_PAI_WF	72361417120	workflow	KILLED	100%		October 7, 2022 8:15 PM	21s	0003054-220806233739962-oozie-oozi-W
ALIENACAO_PAI_WF	72361417120	workflow	KILLED	100%		October 6, 2022 8:15 PM	19s	0003012-220806233739962-oozie-oozi-W
ALIENACAO_PAI_WF	72361417120	workflow	KILLED	100%		October 5, 2022 8:15 PM	18s	0002970-220806233739962-oozie-oozi-W

Através do HUE conseguimos visualizar de forma gráfica os “databases” existentes no HDFS, e seus componentes. Tal visualização é separada por fonte de dados capaz de apresentar as informações contidas nos databases (ex: hive, impala, solr, etc) até o nível de coluna.



Podemos visualizar também várias informações importantes sobre cada tabela clicando com o botão direito do mouse sobre a tabela desejada e selecionando a opção *Table Browser*.



The screenshot shows the Database Browser interface. On the left, a sidebar lists tables under the 'oficina_hue' database. The table 'tb_oficina_hue_txt_ext' is selected, and a context menu is open with the 'Open in Browser' option highlighted. A red arrow points from this option to the 'Table Browser' view on the right. The 'Table Browser' view displays the table's properties, statistics, and schema. Red arrows point to the 'Stats' section, specifically to the 'Data last updated' and 'Schema last modified' information.

Table Browser

Hive ▾

Databases > oficina_hue > tb_oficina_hue_txt_ext

Query Refresh

Overview Sample (100) Details

PROPERTIES

Table

External and stored in [location](#)

Created by oficina_hue on Mon Apr 11 17:31:40 BRT 2022

STATS

Data last updated on 05/02/2022 4:28 PM -03:00

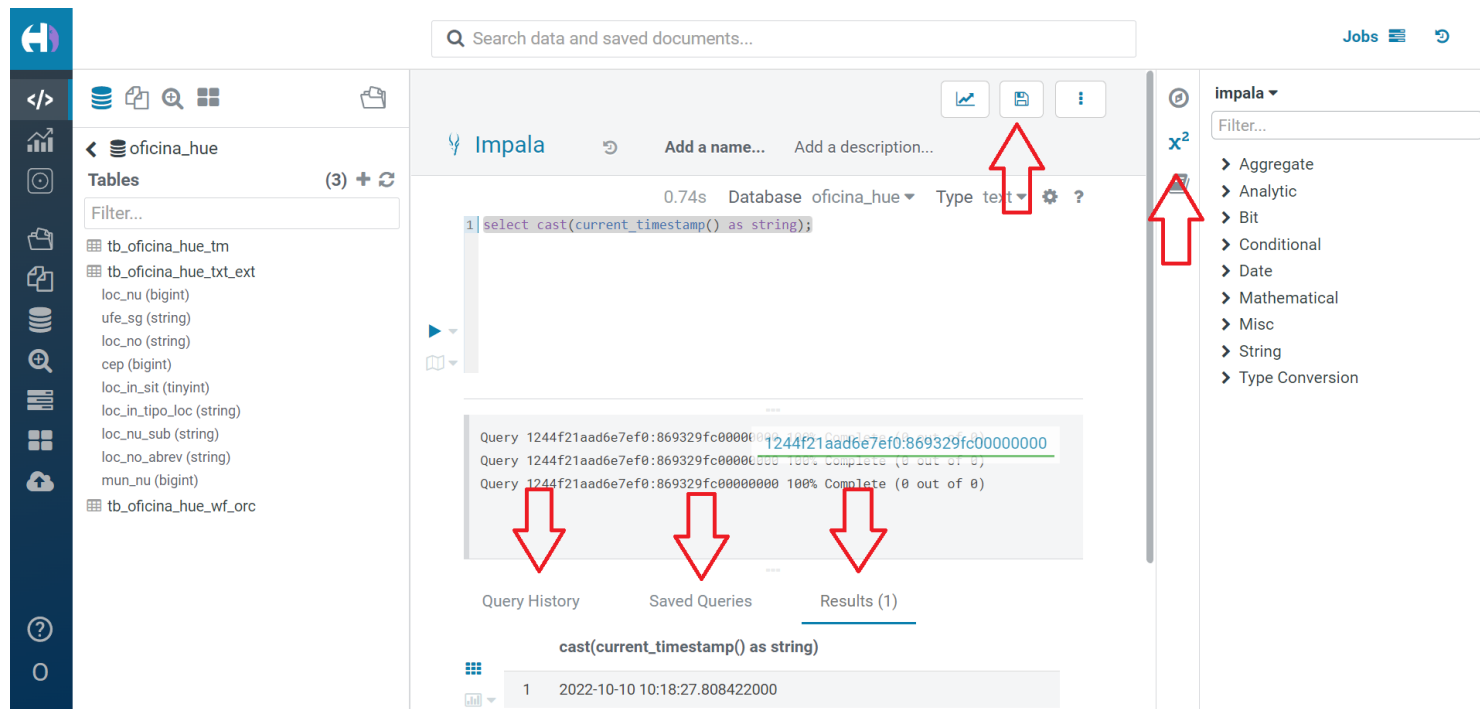
Schema last modified on 05/02/2022 4:28 PM -03:00 by s_cargaserpro

SCHEMA

Filter...

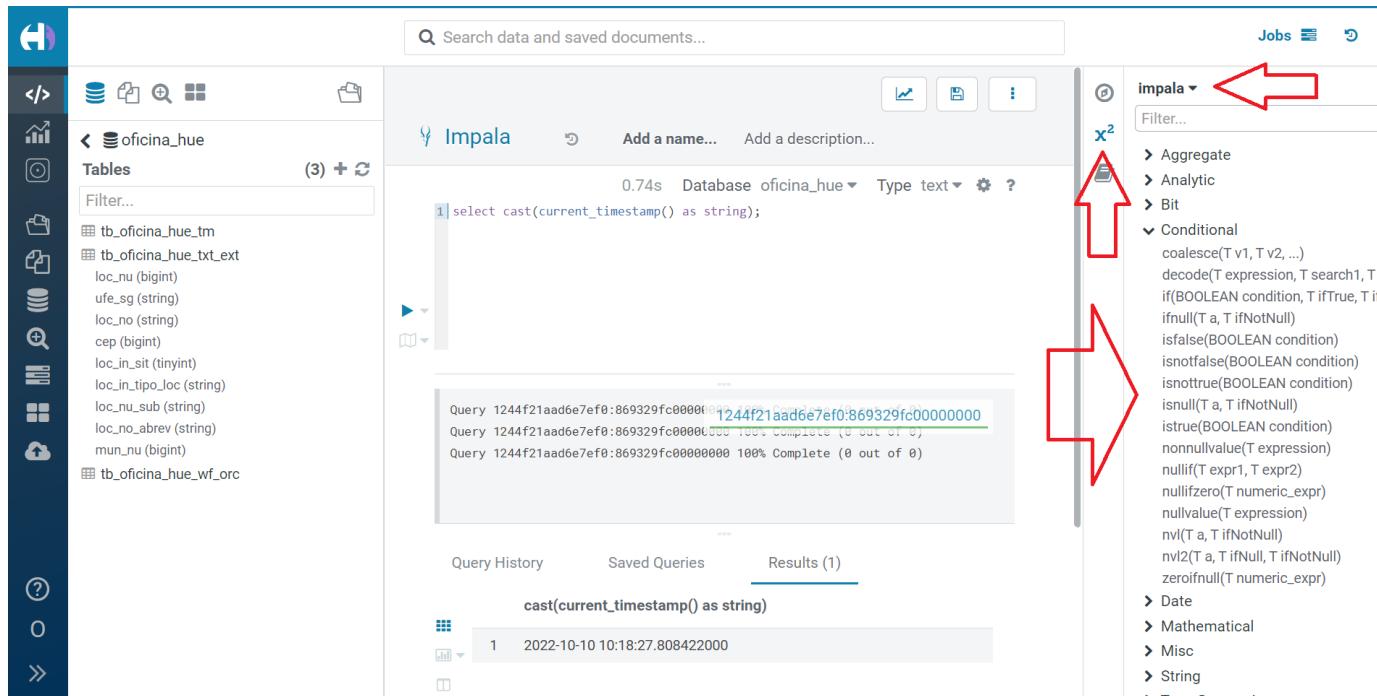
Column (9)	Type	Description	Sample
i loc_nu	bigint		19 13
i ufe_sg	string		AC AC
i loc_no	string		Sena Madureira Plácido de Castro
i cep	bigint		69940000 69928000
i loc_in_sit	tinyint		0 0

O Query Editor além de permitir a criação de consultas às fontes de dados conectadas através do HUE, permite também que salvemos estas consultas no HDFS e mantém um histórico de todas as consultas executadas pelo usuário.



The screenshot displays the HUE Query Editor interface. On the left, a sidebar shows a file tree for the 'oficina_hue' database, listing tables like 'tb_oficina_hue_tm' and 'tb_oficina_hue_wf_orc'. The main area contains a query editor with the SQL statement: `select cast(current_timestamp() as string);`. Below the editor, a 'Query History' panel shows three entries, each with a unique query ID and a '100% Complete' status. Red arrows point to these entries. To the right, a 'Jobs' panel lists various functions under the 'impala' category, including 'Aggregate', 'Analytic', 'Bit', 'Conditional', 'Date', 'Mathematical', 'Misc', 'String', and 'Type Conversion'. Another red arrow points to this panel. At the top right, a search bar is labeled 'Search data and saved documents...'. The bottom of the interface shows the 'Results (1)' tab with a single row of data: `cast(current_timestamp() as string)` and the timestamp `2022-10-10 10:18:27.808422000`.

Possui uma aba lateral à direita que nos auxilia na utilização de funções de manipulação de dados, mostrando os nomes das mesmas e a sintaxe a ser utilizada.



The screenshot displays the Query Editor interface. On the left, a sidebar shows a list of tables under the 'oficina_hue' database. The main area shows an Impala query editor with a search bar and a list of saved queries. On the right, a function library is visible, listing various functions such as Aggregate, Analytic, Bit, Conditional, Date, Mathematical, Misc, and String. Red arrows highlight the function library and the 'impala' dropdown menu.

Search data and saved documents...

Jobs

impala

Filter...

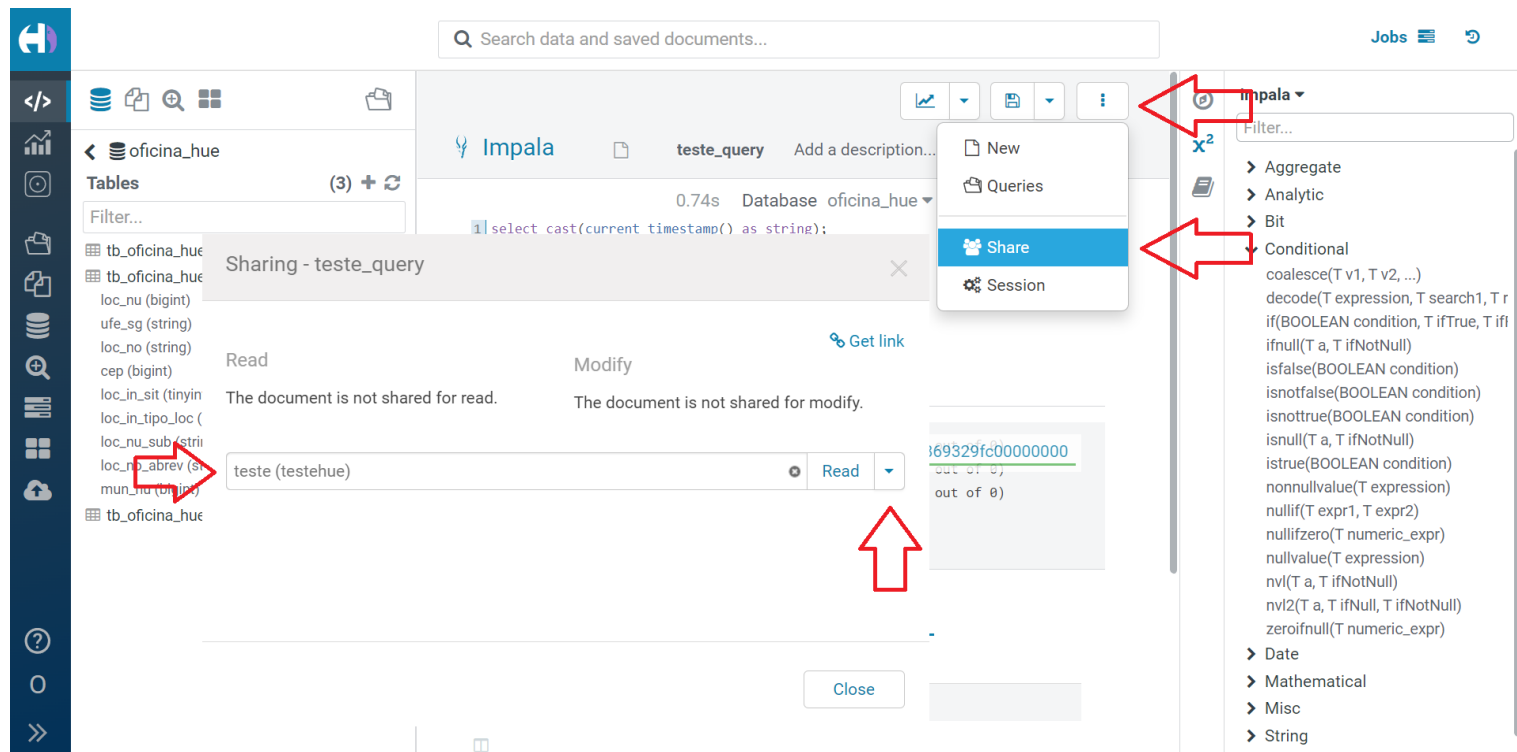
- > Aggregate
- > Analytic
- > Bit
- > Conditional
 - coalesce(T v1, T v2, ...)
 - decode(T expression, T search1, T r
 - if(BOOLEAN condition, T ifTrue, T if
 - ifnull(T a, T ifNotNull)
 - isfalse(BOOLEAN condition)
 - isnotfalse(BOOLEAN condition)
 - isnottrue(BOOLEAN condition)
 - isnull(T a, T ifNotNull)
 - istrue(BOOLEAN condition)
 - nonnullvalue(T expression)
 - nullif(T expr1, T expr2)
 - nullifzero(T numeric_expr)
 - nullvalue(T expression)
 - nvl(T a, T ifNotNull)
 - nvl2(T a, T ifNull, T ifNotNull)
 - zeroifnull(T numeric_expr)
- > Date
- > Mathematical
- > Misc
- > String

Query History Saved Queries Results (1)

cast(current_timestamp()) as string

Query ID	Query Text
1	2022-10-10 10:18:27.808422000

Existe ainda a possibilidade de o usuário compartilhar suas consultas salvas com outros usuários do cluster através do menu “Partilhar”. Basta digitar o nome do usuário na caixa de texto e selecionar qual tipo de permissão o usuário terá em sua consulta.



The screenshot displays the Query Editor interface. On the left, a sidebar shows a list of tables under the 'oficina_hue' database. A red arrow points to the 'loc_no' column in the 'loc_no' table. In the center, a 'Sharing - teste_query' dialog box is open. It contains a text input field with 'teste (testehue)' and a dropdown menu set to 'Read'. A red arrow points to the 'Read' dropdown. To the right, a 'Share' menu is open, showing options: 'New', 'Queries', 'Share' (highlighted with a red arrow), and 'Session'. The background shows a query editor with a query: `select cast(current timestamp() as string);` and a results pane showing a single row with the value '2023-09-29 10:00:00.000000000'.



X



IMPALA

- Escrito em C++
- Data Marts
- Boa escolha para análise interativa e ad-hoc, pela rápida resposta
- Boa escolha para ferramentas de Business Intelligence que permitem aos usuários alterar consultas
- Usa Parquet como o formato de arquivo aconselhado
- Processamento em memória podendo falhar por falta de recurso.
- Utiliza Massively Parallel Processing (MPP)

HIVE

- Escrito em Java
- Data Warehouse
- Boa escolha para consultas de longa duração que requerem transformações pesadas ou múltiplas junções
- Boa escolha para painéis que são predefinidos e não personalizáveis pelo visualizador
- Usa ORC como o formato de arquivo aconselhado
- Funciona melhor com JSON do que Impala
- Se ajusta ao recurso disponível, sendo mais tolerante a falhas.
- Utiliza modelo Map-reduce / Apache Tez

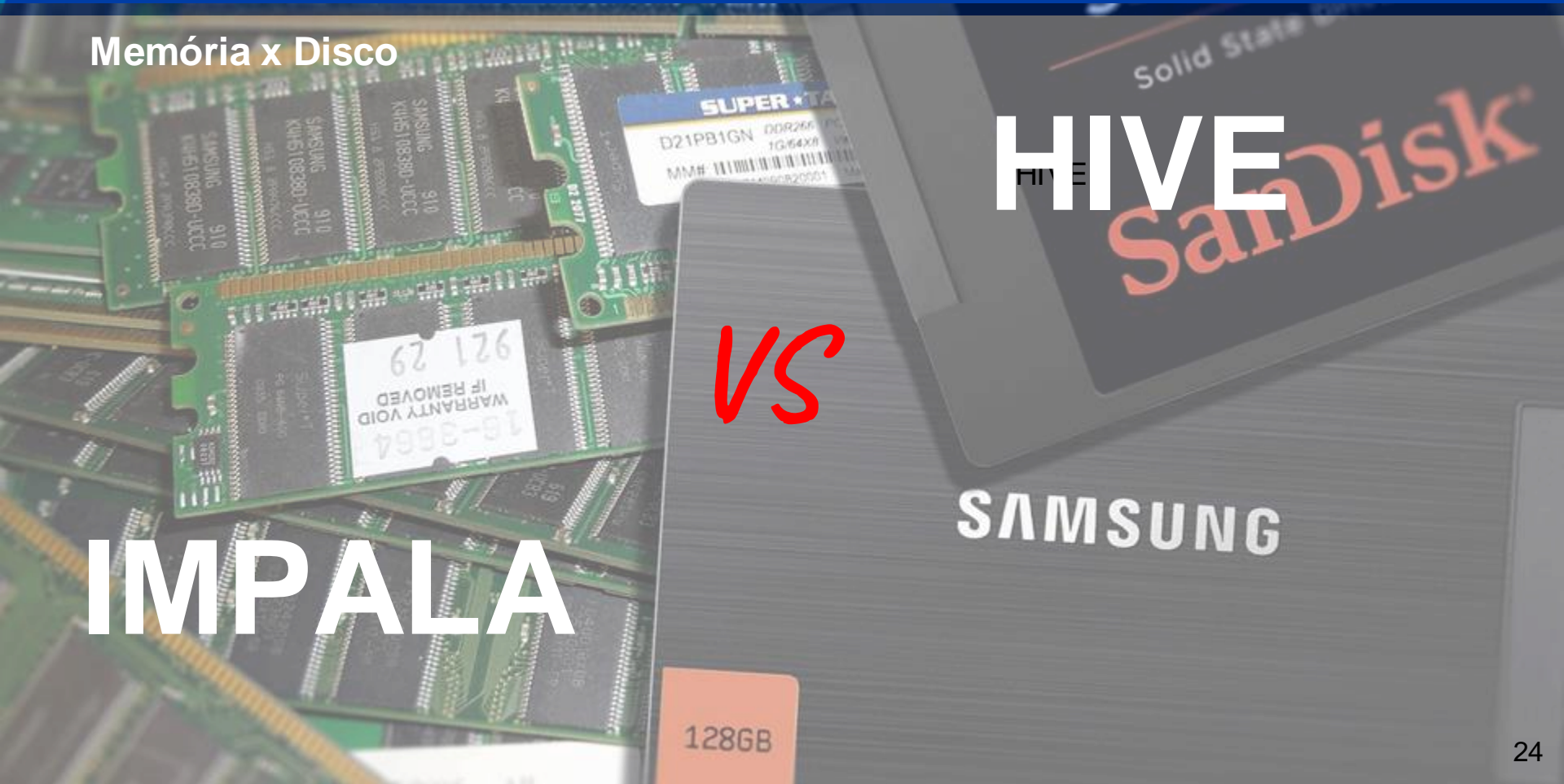
Memória x Disco

IMPALA

HIVE

VS

SAMSUNG



PARQUET

O Impala permite que você crie, gerencie e consulte tabelas Parquet. Parquet é um formato de arquivo binário, orientado a colunas, altamente eficiente para os tipos de consultas em grande escala. O formato Parquet é extremamente aconselhado para consultas que verificam colunas específicas em uma tabela, por exemplo, para consultar tabelas "largas" com muitas colunas ou para realizar operações de agregação como SUM () e AVG () que precisam processar a maioria ou todos os valores de uma coluna. Cada arquivo de dados Parquet escrito pelo Impala contém os valores de um conjunto de linhas (referido como "grupo de linhas").

Em um arquivo de dados, os valores de cada coluna são organizados de forma que sejam todos adjacentes, permitindo uma boa compactação para os valores dessa coluna. As consultas em uma tabela Parquet podem recuperar e analisar esses valores de qualquer coluna rapidamente e com o mínimo de I/O.



Sales			
Product	Customer	Date	Sale
Beer	Thomas	2011-11-25	2 GBP
Beer	Thomas	2011-11-25	2 GBP
Vodka	Thomas	2011-11-25	10 GBP
Whiskey	Christian	2011-11-25	5 GBP
Whiskey	Christian	2011-11-25	5 GBP
Vodka	Alexei	2011-11-25	10 GBP
Vodka	Alexei	2011-11-25	10 GBP

Product	
ID	Value
1	Beer
2	Beer
3	Vodka
4	Whiskey
5	Whiskey
6	Vodka
7	Vodka

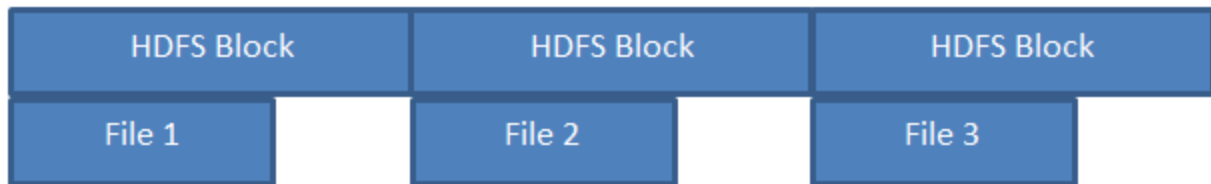
Customer	
ID	Customer
1	Thomas
2	Thomas
3	Thomas
4	Christian
5	Christian
6	Alexei
7	Alexei

Formato mais apropriado de arquivo

Normalmente, para grandes volumes de dados (vários gigabytes por tabela ou partição), o tipo de arquivo Parquet tem um desempenho melhor devido a sua combinação de layout de armazenamento em forma de coluna, grandes requisições de I/O, compressão e codificação.

Evite arquivos pequenos

- Arquivos menores que o blocos do HDFS ocuparão o bloco inteiro.
- Recuperar muitos arquivos tem um I/O maior que poucos, já que é lido bloco a bloco.
- Iniciar tarefas (tasks) para cada arquivo, em caso de pequenos arquivos, gerará um número muito maior de tasks, e com isso um custo muito maior de processamento.



Evite arquivos pequenos

- Sempre use `INSERT ... SELECT` para copiar grandes volumes de dados de uma tabela para outra dentro do Impala. Evite `INSERT ... VALUES` para qualquer volume substancial de dados ou tabelas de performance crítica, pois cada cláusula cria um pequeno arquivo de dados separado.
- Caso seu processo de ingestão utilize o conceito de delta, e cada delta ingerido seja com um baixo volume de dados, talvez seja necessário de tempos em tempos gerar novamente sua tabela definitiva de dados, pois a ingestão provavelmente tem gerado pequenos arquivos parquet em cada execução e poderá fazer com que sua performance nas consultas seja prejudicada. Para isso pode utilizar o conceito de `INSERT OVERWRITE TABLE`.

Granularidade das partições

- Quando você envia consultas que solicitam um valor específico ou um intervalo de valores para a coluna chave da partição, o Impala pode evitar a leitura de dados irrelevantes, resultando em uma grande economia de I/O em disco.
- Escolha uma estratégia de particionamento que insira pelo menos 256 MB de dados em cada partição para tirar vantagem do bulk de I/O do HDFS e das consultas distribuídas do Impala.
- Se você precisar reduzir o número de partições e aumentar o número de dados em cada partição, primeiro verifique as colunas chaves da partição que raramente são referenciadas ou que não são usadas por consultas críticas. Por exemplo, os dados de log do seu web site são particionados por ano, mês, dia e hora, mas a maioria das consultas são por dia, talvez então você precise apenas particionar por ano, mês e dia.

Números inteiros menores para as colunas chaves das partições e para os joins

- Embora seja tentador usar strings para particionar coluna chave, já que esses valores se tornam nomes de diretórios HDFS, você pode diminuir a utilização de memória usando valores numéricos como campo chave de partição comum tais como ANO, MÊS e DIA. Use o menor número inteiro que possua o intervalo apropriado de dados, geralmente TINYINT para MÊS e DIA, e SMALLINT para ANO.
- Use a função `EXTRACT()` para remover campos de data e hora individualmente de um valor `TIMESTAMP`, e `CAST()` para retornar valor inteiro. Ao realizar consultas, de preferência para a utilização de colunas numéricas para a realização dos JOINS entre as tabelas, permitindo um cruzamento das informações sem necessidade de grandes varreduras ou utilização de memória.

Minimize resultados se possível

Utilize as seguintes técnicas:

- Agregação: se você quer saber quantas linhas satisfazem uma determinada condição, o valor total de valores de alguma coluna, o menor ou maior valor, chame funções agregadas como COUNT(), SUM(), MIN() e MAX() na consulta para não precisar enviar o resultado para uma aplicação e fazer os cálculos na aplicação.
- Filtro: use os filtros aplicáveis na cláusula WHERE da consulta para eliminar linhas que não são relevantes.
- Cláusula LIMIT: se você precisa visualizar apenas alguns valores de um conjunto de dados, ou os valores superiores ou inferiores de uma consulta usando ORDER BY, inclua a cláusula LIMIT para reduzir o tamanho do retorno da consulta.

Gerenciamento de Memória em Consultas

Em certas situações (queries complexas e/ou que tendem a demandar grandes quantidades de dados em memória) o otimizador do Impala pode tentar reservar um volume de memória maior que o necessário para o processamento da consulta gerando dois tipos de problemas:

- Query não ser executada por não haver recurso de memória suficiente no pool. Neste caso a query sequer será admitida para execução.
- Query é admitida para execução, porém reservará mais recurso de memória que o necessário, o que pode ocasionar a não execução de outras consultas no mesmo pool.

Gerenciamento de Memória em Consultas

As soluções para resolvermos estes tipos de problema seriam:

- Gerenciar manualmente o recurso de memória a ser utilizado pelo Impala.
Ex: `SET mem_limit="1G";`
- Adaptar sua consulta para execução no Hive, uma vez que o mesmo é mais robusto por usar `map_reduce`.
- Particionar sua consulta para utilização de uma quantidade de memória menor.
Utilizar tabelas intermediárias, realizando assim uma consulta maior em várias etapas, reduzindo a utilização de memória em relação a uma única instrução sql mais complexa.

Melhorando Performance e Consumo de Recursos Utilizando Estatísticas

Tanto o Impala quanto o Hive se utilizam das estatísticas geradas em suas tabelas para definição de plano de acesso e execução das consultas no Hadoop. Uma prática essencial para que consiga extrair a melhor performance de suas consultas é sempre manter as estatísticas das tabelas atualizadas. Para tanto devemos nos atentar em executar a atualização das estatísticas sempre que realizarmos alguma modificação nos dados das tabelas a serem consumidas.

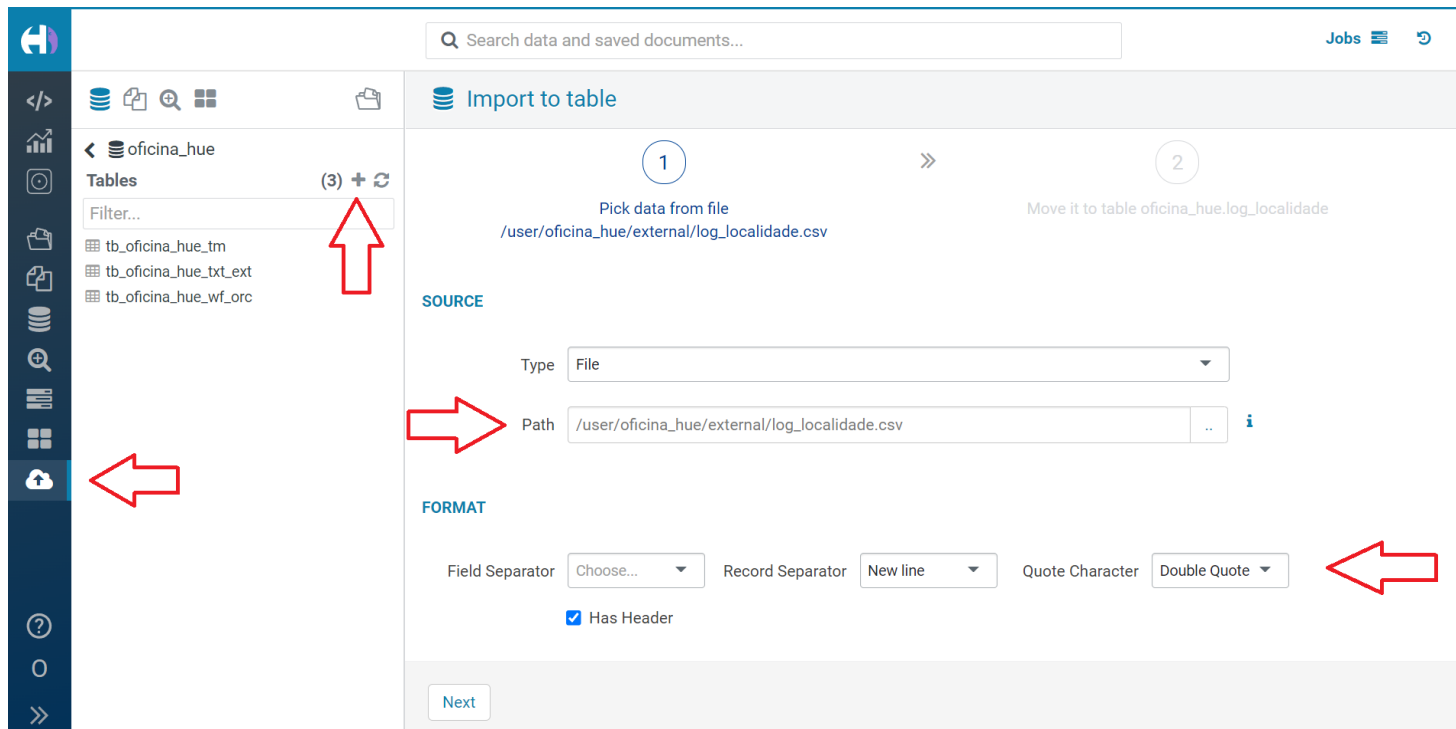
Comandos:

Hive: `Analyze table nome_da_tabela compute statistics;`

Impala: `compute stats nome_da_tabela;`

Hands on

Criando estruturas e ingerindo dados



Search data and saved documents...

Jobs

Import to table

1 Pick data from file
/user/oficina_hue/external/log_localidade.csv

2 Move it to table oficina_hue.log_localidade

SOURCE

Type File

Path /user/oficina_hue/external/log_localidade.csv

FORMAT

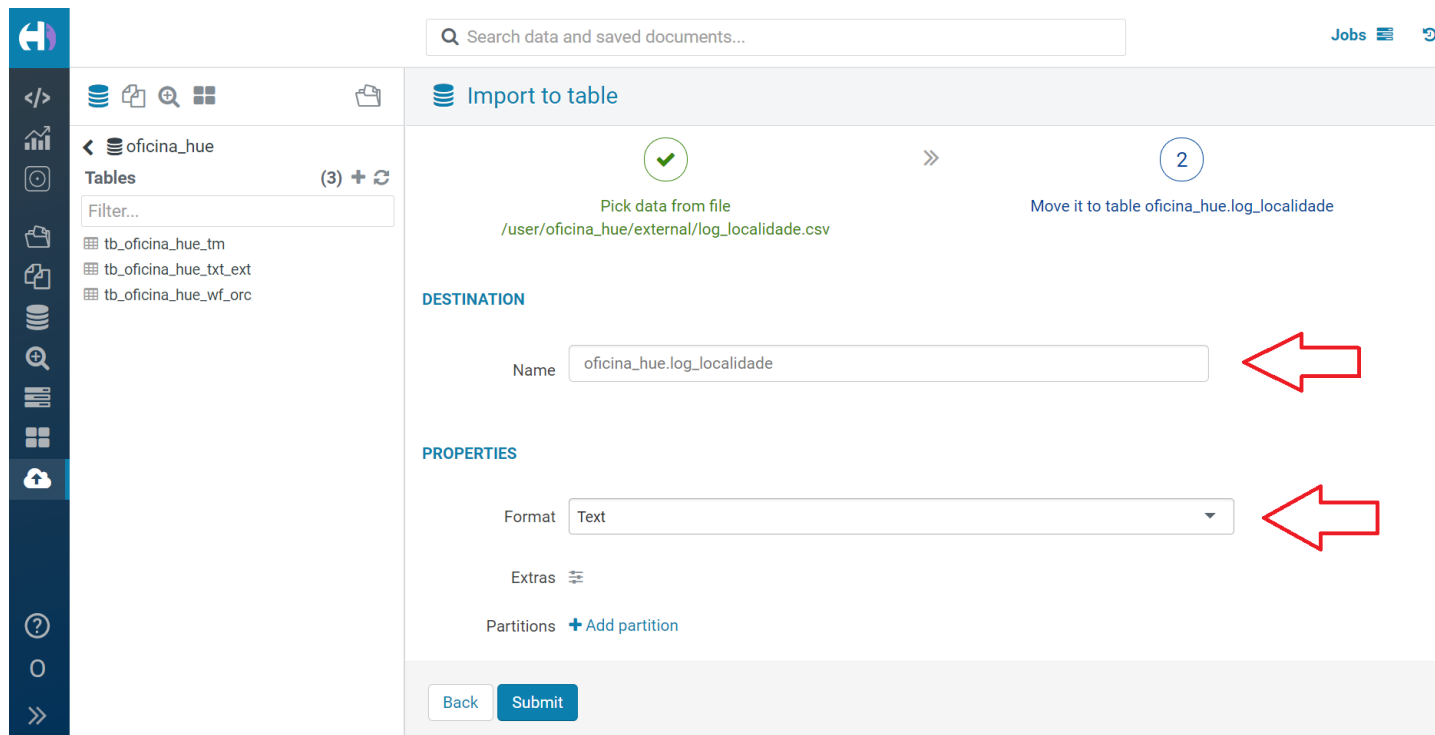
Field Separator Choose... Record Separator New line Quote Character Double Quote

☒ Has Header

Next

Hands on

Criando estruturas e ingerindo dados



The screenshot displays the Databricks Data Lake interface. On the left, a sidebar shows the workspace structure with a folder named 'oficina_hue' containing three tables: 'tb_oficina_hue_tm', 'tb_oficina_hue_txt_ext', and 'tb_oficina_hue_wf_orc'. The main area is titled 'Import to table' and shows a workflow with two steps: 1. 'Pick data from file' (source: '/user/oficina_hue/external/log_localidade.csv') and 2. 'Move it to table oficina_hue.log_localidade'. Below the workflow, the 'DESTINATION' section shows the table name 'oficina_hue.log_localidade' in a text input field, with a red arrow pointing to it. The 'PROPERTIES' section shows the 'Format' set to 'Text' in a dropdown menu, also with a red arrow pointing to it. At the bottom, there are 'Back' and 'Submit' buttons.

Hands on

Criando external table

```
CREATE EXTERNAL TABLE oficina_hue.tb_log_cpc_ext (cpc_nu bigint, ufe_sg string,  
loc_nu bigint, cpc_no string, cpc_endereco string, cep bigint)  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ';'   
LOCATION '/user/oficina_hue/external'   
TBLPROPERTIES("skip.header.line.count"="1")
```

External tables funcionam como uma estrutura (esqueleto) de uma tabela que aponta diretamente para o arquivo/local, dentro do HDFS, indicado no momento de sua criação.

Quando dropamos uma external table, os arquivos referenciados por ela continuam no diretório em que se encontram, sendo eliminada apenas a estrutura de acesso.

Hands on

Detalhe importante sobre internal (managed tables) e external tables:

Por definição, hoje no cluster de Brasília todas as tabelas criadas, independente de colocarmos de forma explícita que serão external, serão criadas como tal. Vale salientar que o conceito de external continuará valendo para as que explicitamente forem declaradas desta forma, logo se dropar uma tabela criada com o comando `CREATE EXTERNAL TABLE`, os dados continuarão no diretório cujo apontamento foi definido na criação da tabela, mas as tabelas external criadas sem explicitamente declaradas forem, quando dropadas, apagarão tanto a estrutura da mesma quanto os dados.

A exceção se dá para tabelas ACID (Atomicidade, Consistência, Isolamento e Durabilidade), estas serão criadas como tabelas internas gerenciadas pelo Hive (`TBLPROPERTIES('transactional'='true')`).

Hands on

Ingerindo/criando tabela linha de comando

```
CREATE TABLE oficina_hue.tb_log_cpc_txt( cpc_nu bigint, ufe_sg string, loc_nu bigint,  
cpc_no string, cpc_endereco string, cep bigint)  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ';'   
STORED AS TEXTFILE  
TBLPROPERTIES("skip.header.line.count"="1");
```

```
LOAD DATA INPATH '/user/oficina_hue/dataset' into table oficina_hue.tb_log_cpc_txt;
```

Quando utilizamos o comando *LOAD DATA* para inserirmos dados de um dado arquivo/diretório do HDFS para uma tabela específica, os arquivos são movidos para o diretório da tabela em questão e quando dropada a tabela, os arquivos serão literalmente apagados do HDFS.

Hands on

Criando tabela Parquet utilizando uma tabela texto.

```
CREATE TABLE oficina_hue.tb_log_cpc_pq STORED AS PARQUET as  
select * from oficina_hue.tb_log_cpc_txt;
```

Podemos realizar a criação de tabelas tipo parquet, tanto pelo Impala (preferencialmente) quanto pelo Hive, através da consulta via HQL em uma tabela já existente.

Hands on

Criando Views

```
CREATE VIEW vw_log_cpc as  
SELECT * FROM log_cpc  
WHERE ufe_sg = 'AL';
```

Views podem ser criadas para armazenamento de queries pré-definidas como de costume, mas podem representar também uma restrição para um grupo de usuários com menos privilégio.

Obs: para se conseguir criar uma view, o usuário necessita ter acesso à todos os dados cuja view faz referência.

Hands on

Views Hive x Impala (incompatibilidades)

Ao criar uma view, é importante observarmos em que editor (hive ou impala) a view está sendo criada, devido ao comportamento das definições desta view relativos tanto às estruturas de dados, quanto às funções nativas de cada editor. Isto vale tanto para views quanto para tabelas.

Por exemplo, quanto ao uso de tipos e/ou estruturas de dados que funcionam exclusivamente em um dos editores. O impala, não conseguirá ler uma view que possua colunas do tipo binary ou utilizem funções específicas do Hive.

Hands on

Views / Tabelas Hive x Impala (incompatibilidades)

No Hive:

```
CREATE TABLE oficina_hue.teste_date(id_teste int, dt_teste date);
```

```
INSERT INTO oficina_hue.teste_date VALUES (1,'2021-01-01');
```

```
CREATE VIEW oficina_hue.vw_teste_date as
```

```
SELECT id_teste, date_format(dt_teste,'y') as year FROM oficina_hue.teste_date;
```

No Impala:

```
SELECT * FROM oficina_hue.vw_teste_date;
```

Hands on

Materialized views

Materialized Views não eram suportadas na versão do Hive que utilizávamos (Hive 2.1.1 - CDH 6.3.4), mas com a migração do cluster, já temos a versão Hive 3.1.3 - CDP 7.1.7 em produção, estando disponível à mesma. Para que a mesma funcione, deve-se criar a tabela cuja Materialized View irá realizar a consulta, utilizando o parâmetro 'transactional'='true' e um formato que permita que a tabela seja transacionada (Ex: ORC).

```
CREATE TABLE oficina_hue.teste_date(id_teste int, dt_teste date) stored as orc  
tblproperties('transactional'='true');
```

```
CREATE materialized VIEW oficina_hue.vw_teste_date as  
select id_teste, dt_teste from oficina_hue.teste_date;
```

Hands on

Materialized views

Para atualização dos dados das materialized views devemos usar o comando do exemplo abaixo:

```
ALTER MATERIALIZED VIEW oficina_hue.vw_teste_date REBUILD;
```

As tabelas ACID no hive suportam tanto insert, quanto update e delete, mas devemos salientar um detalhe. As materialized views criadas com base em tabelas que suportam as três operações de modificação de dados, não poderão ser lidas pelo Impala, o mesmo só aceita leitura de materialized views criadas em tabelas cujo ACID suportado seja apenas insert (tblproperties 'transactional_properties'='insert_only')

Considerando a versão Cloudera do condomínio atual, estes são os tipos de dados suportados, respectivamente, pelo hive e impala.



Hive	Impala
ARRAY*	ARRAY*
BIGINT	BIGINT
BOOLEAN	BOOLEAN
CHAR	CHAR
DATE	DATE
DECIMAL	DECIMAL
DOUBLE	DOUBLE
FLOAT	FLOAT
INT	INT
MAP*	MAP*
REAL	REAL
SMALLINT	SMALLINT
STRING	STRING
STRUCT*	STRUCT*
TIMESTAMP	TIMESTAMP
TINYINT	TINYINT
VARCHAR	VARCHAR
BINARY	-

Complex types

Também conhecidos como tipos aninhados (nested types), esses tipos de dados permitem a representação de vários valores de dados em uma única posição de linha/coluna. Eles diferem dos tipos de coluna conhecidas, como BIGINT e STRING, conhecidos como tipos escalares ou tipos primitivos, que representam um único valor de dados em uma determinada posição de linha/coluna.



Tipos dados Complex types

- **Array:** Um tipo complex type que pode representar um número arbitrário de elementos ordenados. Os elementos podem ser escalares ou outro tipo complexo (ARRAY, STRUCT ou MAP)
- **Map:** representa um conjunto arbitrário de pares de valores-chave. A parte chave deve ser, necessariamente, um tipo escalar, enquanto a parte do valor pode ser um escalar ou outro tipo complexo. Um exemplo simples de pensarmos em um uso para o tipo map, seria representarmos um conjunto de métricas para uma tabela de municípios, onde poderíamos ter informações como o código e nome do município como campos escalares e um campo do tipo map, contendo uma chave com o nome da métrica e um valor correspondente para o município (exemplos: 'Area:432244', 'Populacao':213685, 'PibPerCapta:R\$ 19.184,77')
- **Struct:** este tipo é útil para combinar duas tabelas relacionadas, de modo a minimizar a repetição. A maneira mais comum de representar esses dados é como um ARRAY de elementos STRUCT. O tipo STRUCT costuma ser mais útil como um item de um ARRAY ou a parte do valor do par de valores-chave em um MAP.

Manipulando complex types - Impala x Hive

O Impala usa a notação de ponto para se referir a nomes de elementos ou elementos dentro de tipos complexos e notação de join para fazer referência cruzada de colunas escalares com os elementos de tipos complexos na mesma linha, em vez da cláusula LATERAL VIEW e da função EXPLODE () do HiveQL.

Uma consulta como `SELECT * FROM banco.tabela` feita pelo hive sobre uma tabela que contenha colunas do tipo complex type é aceita normalmente. Já no impala, apesar de não termos erro de execução, o retorno será apenas dos dados referentes às colunas escalares.

A versão atual do Impala não suporta a criação de novos dados com colunas do tipo complex type. Assim, o caminho que demonstraremos aqui é criarmos uma tabela do tipo parquet e alimentá-la com o hive, para então realizarmos consultas via impala.

Hands on - Carregando dados Complex Type

```
[
  {
    "nome": "Joao",
    "sobrenome": "da Silva",
    "idade": "25",
    "endereco": {
      "logradouro": "Rua Abacaxi",
      "municipio": "Manaus",
      "uf": "AM",
      "cep": "11111-111"
    },
    "telefone": [
      {
        "tipo": "residencial",
        "numero": "92 1111-1234"
      },
      {
        "tipo": "celular",
        "numero": "92 3333-4567"
      }
    ]
  },
]
```

```
{
  "nome": "Maria",
  "sobrenome": "das Oliveiras",
  "idade": "29",
  "endereco": {
    "logradouro": "Rua Laranja",
    "municipio": "Feira de Santana",
    "uf": "BA",
    "cep": "22222-222"
  },
  "telefone": [
    {
      "tipo": "residencial",
      "numero": "75 2222-4321"
    },
    {
      "tipo": "celular",
      "numero": "75 4444-5678"
    }
  ]
},
```

```
{
  "nome": "Pedro",
  "sobrenome": "dos Santos",
  "idade": "40",
  "endereco": {
    "logradouro": "Rua Limoeiro",
    "municipio": "Pirenópolis",
    "uf": "GO",
    "cep": "33333-333"
  },
  "telefone": [
    {
      "tipo": "residencial",
      "numero": "62 3333-9876"
    },
    {
      "tipo": "celular",
      "numero": "62 5555-4321"
    }
  ]
},
]
```


Hands on - Carregando dados Complex Type

```
{ "nome": "Joao", "sobrenome": "da Silva", "idade": "25", "endereco": { "logradouro": "Rua Abacaxi", "municipio": "Manaus", "uf": "AM", "cep": "11111-111" }, "telefone": [ { "tipo": "residencial", "numero": "92 1111-1234" }, { "tipo": "celular", "numero": "92 3333-4567" } ] }
```

```
{ "nome": "Maria", "sobrenome": "das Oliveiras", "idade": "29", "endereco": { "logradouro": "Rua Laranja", "municipio": "Feira de Santana", "uf": "BA", "cep": "22222-222" }, "telefone": [ { "tipo": "residencial", "numero": "75 2222-4321" }, { "tipo": "celular", "numero": "75 4444-5678" } ] }
```





```
{ "nome": "Pedro", "sobrenome": "dos Santos", "idade": "40", "endereco": { "logradouro": "Rua Limoeiro", "municipio": "Pirenópolis", "uf": "GO", "cep": "33333-333" }, "telefone": [ { "tipo": "residencial", "numero": "62 3333-9876" }, { "tipo": "celular", "numero": "62 5555-4321" } ] }
```

Hands on - Carregando dados Complex Type

 Browser de ficheiros

 Voltar

 Início

Página 1 to 1 de 1    

 Editar ficheiro

 Refresh

 Visualizar como
binário

 Download

Última modificação
24/10/2022 14:59
-03:00

User
oficina_hue

Grupo
oficina_hue

Size
955 B

Modo
100644

/ user / oficina_hue / json2hadoop / pessoa-formatted.json

```
{ "nome": "Joao", "sobrenome": "da Silva", "idade": "25", "endereco": { "logradouro": "Rua Abacaxi", "municipio": "Manaus", "uf": "AM", "cep": "11111-111" }, "telefone": [ { "tipo": "residencial", "numero": "92 1111-1234" }, { "tipo": "celular", "numero": "92 3333-4567" } ] }
{ "nome": "Maria", "sobrenome": "das Oliveiras", "idade": "29", "endereco": { "logradouro": "Rua Laranja", "municipio": "F eira de Santana", "uf": "BA", "cep": "22222-222" }, "telefone": [ { "tipo": "residencial", "numero": "75 2222-4321" }, { "tipo": "celular", "numero": "75 4444-5678" } ] }
{ "nome": "Pedro", "sobrenome": "dos Santos", "idade": "40", "endereco": { "logradouro": "Rua Limoeiro", "municipio": "Pir enópolis", "uf": "GO", "cep": "33333-333" }, "telefone": [ { "tipo": "residencial", "numero": "62 3333-9876" }, { "tipo": "celular", "numero": "62 5555-4321" } ] }
```

```
DROP TABLE IF EXISTS oficina_hue.pessoa_stage;
```

```
DROP TABLE IF EXISTS oficina_hue.pessoa_stage;
```

```
CREATE EXTERNAL TABLE oficina_hue.pessoa_stage (
```

```
nome string
, sobrenome string
, idade string
, endereco struct<logradouro:string
, municipio:string
, uf:string
, cep:string>
, telefone array<struct<tipo:string
, numero:string>>
)
```

```
ROW FORMAT SERDE 'org.openx.data.jsonserde.JsonSerDe'
```

```
LOCATION 'hdfs://ldg-bsa-01-hdfs/user/oficina_hue/json2hadoop/';
```

```
DROP TABLE IF EXISTS oficina_hue.pessoa;
```

```
CREATE TABLE pessoa STORED AS PARQUET AS SELECT * FROM oficina_hue.pessoa_stage;
```

```
DROP TABLE IF EXISTS oficina_hue.pessoa_stage PURGE;
```

Obs.: Lembre-se que o impala pode levar um tempo para reconhecer os metadados criados pelo Hive. Assim, para vermos as mudanças imediatamente, devemos utilizar o comando `INVALIDATE METADATA`.

E para maior performance nas consultas, é sempre recomendado que geremos as estatísticas com o comando `COMPUTE STATS`.



Hands on - Carregando dados Complex Type

Experimente realizar a consulta abaixo no Hive e no Impala. Em seguida, compare os resultados.

```
select * from oficina_hue.pessoa;
```



Hive:

<u>pessoa.nome</u>	<u>pessoa.sobrenome</u>	<u>pessoa.idade</u>	<u>pessoa.endereco</u>	<u>pessoa.telefone</u>
Joao	da Silva	25	{"logradouro":"Rua Abacaxi","municipio":"Manaus","uf":"AM","cep":"11111-111"}	[{"tipo":"residencial","numero":"92 1111-1234"}, {"tipo":"celular","numero":"92 3333-4567"}]
Maria	das Oliveiras	29	{"logradouro":"Rua Laranja","municipio":"Feira de Santana","uf":"BA","cep":"22222-222"}	[{"tipo":"residencial","numero":"75 2222-4321"}, {"tipo":"celular","numero":"75 4444-5678"}]
Pedro	dos Santos	40	{"logradouro":"Rua Limoeiro","municipio":"Pirenópolis","uf":"GO","cep":"33333-333"}	[{"tipo":"residencial","numero":"62 3333-9876"}, {"tipo":"celular","numero":"62 5555-4321"}]

Impala:

<u>nome</u>	<u>sobrenome</u>	<u>idade</u>
Joao	da Silva	25
Maria	das Oliveiras	29
Pedro	dos Santos	40

Hive:

```
FROM
  pessoa
  LATERAL VIEW explode(telefone) telefone as
  telefones;
```

Impala:

```
INVALIDATE METADATA oficina_hue.pessoa;  
COMPUTE STATS oficina_hue.pessoa;
```

```
SELECT
    p.nome
    ,p.sobrenome
    ,p.idade
    ,p.endereco.logradouro
    ,p.endereco.municipio
    ,p.endereco.uf
    ,p.endereco.cep
    ,t.tipo
    ,t.numero
FROM
    pessoa p
    ,p.telefone t ;
```



Realizar atividades localizadas no documento **Exercícios**, subtítulo **Ingestão de Dados**, enviado por email no início da Oficina.

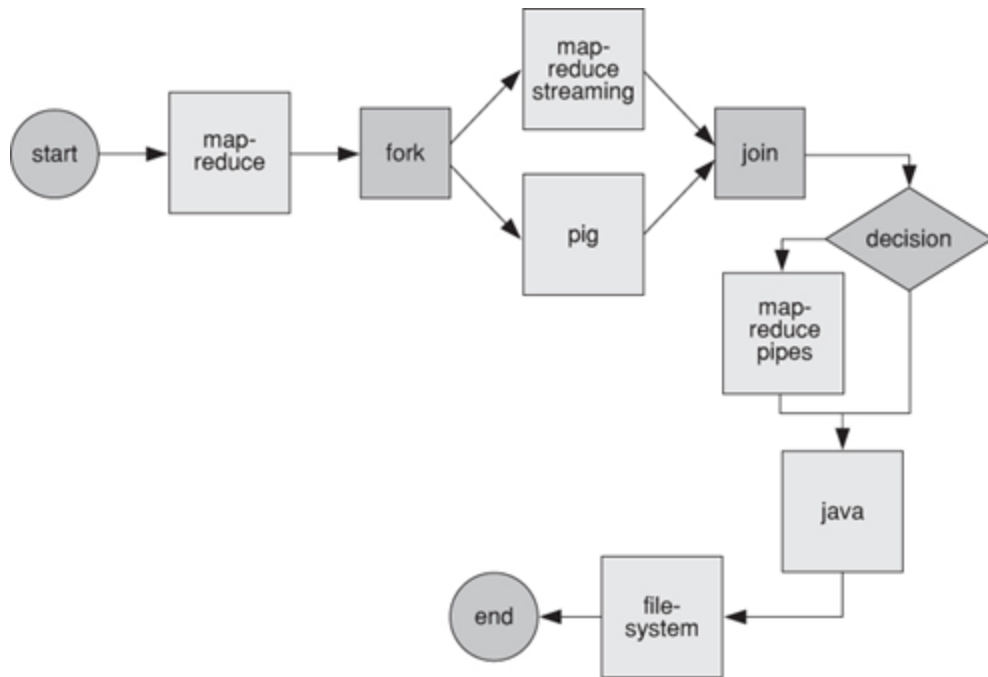
Realizar atividades localizadas no documento **Exercícios**, subtítulo **Ingestão de Dados**, enviado por email no início da Oficina.



Workflows Oozie

O Oozie é um mecanismo de workflow, baseado em servidor, especializado na execução de jobs com ações que executam tarefas em um cluster Hadoop.

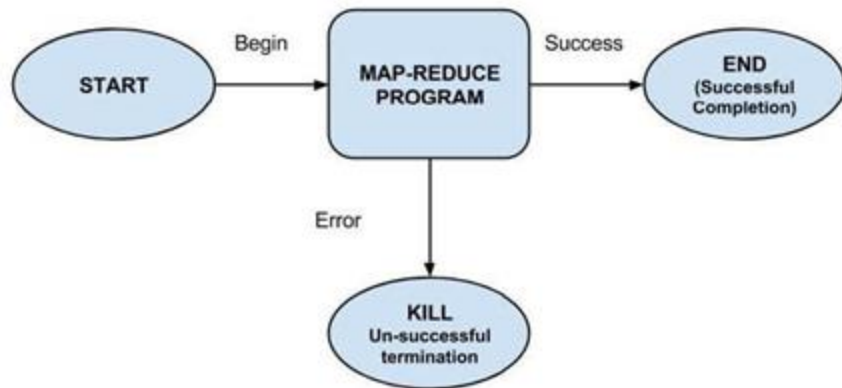
Para fins do Oozie, um Workflow é uma coleção de ações organizadas em um DAG (Direct Acyclic Graph) de controle de dependência.



Workflows Oozie

As definições de workflows do Oozie são escritas em hPDL (uma linguagem de definição de processo XML semelhante ao JBOSS JBPM jPDL).

As ações dos workflows do Oozie iniciam trabalhos em sistemas remotos do Hadoop. Após a conclusão da ação, os sistemas remotos chamam o Oozie de volta para notificar a conclusão da ação, neste ponto o Oozie prossegue para a próxima ação no workflow.



Workflows Oozie

Os workflows do Oozie contém nós de fluxo de controle e nós de ação.

Os nós de fluxo de controle definem o início e o fim de um workflow (nós de início, término e falha) e fornecem um mecanismo para controlar o caminho de execução do workflow (nós de decisão, bifurcação e junção).

Os nós de ação são o mecanismo pelo qual um fluxo de trabalho dispara a execução de uma tarefa de computação / processamento. O Oozie fornece suporte para diferentes tipos de ações: mapreduce do Hadoop, sistema de arquivos do Hadoop, SSH, HTTP, e-mail e subfluxo do Oozie. O Oozie pode ser estendido para oferecer suporte a tipos adicionais de ações.

Hands on

Workflow na prática

Por definição, tudo que criamos e salvamos com a interface do Hue é definido como Documento (queries, workflows, schedules, scripts, etc)

Para acessarmos a interface de construção de workflows do Oozie, bem como os fluxos já criados, devemos acionar o menu Documents, localizado no menu principal do Hue.

The screenshot shows the Hue web interface. On the left sidebar, the 'Documents' icon (a folder) is highlighted with a red arrow. In the top navigation bar, the 'My documents' link is also highlighted with a red arrow. The main content area displays a table of saved documents.

Name	Description	Type	Owner	Last Modified
wkf_drdaniel		Oozie Workflow	oficina_hue	11/08/2020
My Workflow		Oozie Workflow	oficina_hue	11/08/2020
sch_exemplo_pgfn		Oozie Schedule	oficina_hue	11/08/2020
exemplo_pgfn		Oozie Workflow	oficina_hue	11/08/2020
sch_send_mail_042		Oozie Schedule	oficina_hue	11/08/2020
wf_send_mail_042	teste geração email	Oozie Workflow	oficina_hue	11/08/2020
My Workflow		Oozie Workflow	oficina_hue	08/24/2020
sch_teste_email		Oozie Schedule	oficina_hue	08/24/2020
oficina_hue_handson		Oozie Workflow	oficina_hue	08/10/2020
workflow_oficina_hue		Oozie Workflow	oficina_hue	08/10/2020

Hands on

Workflow na prática

Para criação de um novo workflow, devemos acessar o botão *Novo Documento* localizado na parte superior direita da janela do Hue, e selecionar a opção *Workflow*.

The screenshot shows the Hue web interface. On the left is a sidebar with navigation icons. The main area is titled 'My documents' and contains a table with columns 'Name', 'Description', 'Type', and 'Last Modified'. The table lists several documents, including 'wkf_drdaniel', 'My Workflow', 'sch_exemplo_pgfn', 'exemplo_pgfn', and 'sch_send_mail_042'. A red arrow points to the 'New Document' button (represented by a document icon) in the top right corner. Another red arrow points to the 'Workflow' option in the dropdown menu that appears after clicking the button.

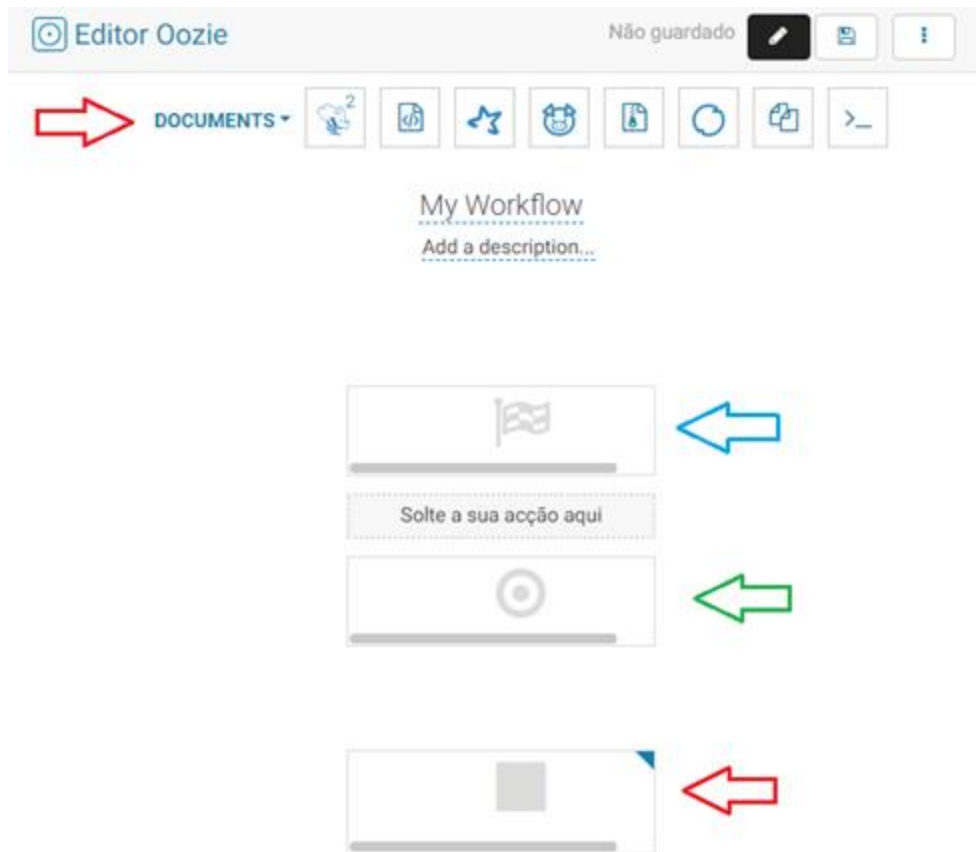
Name	Description	Type	Last Modified
wkf_drdaniel			
My Workflow		Oozie Workflow	
sch_exemplo_pgfn		Oozie Schedule	
exemplo_pgfn		Oozie Workflow	
sch_send_mail_042		Oozie Schedule	

- Hive Query
- Impala Query
- Workflow**
- Schedule
- Bundle
- Dashboard
- New folder

Hands on

Workflow na prática

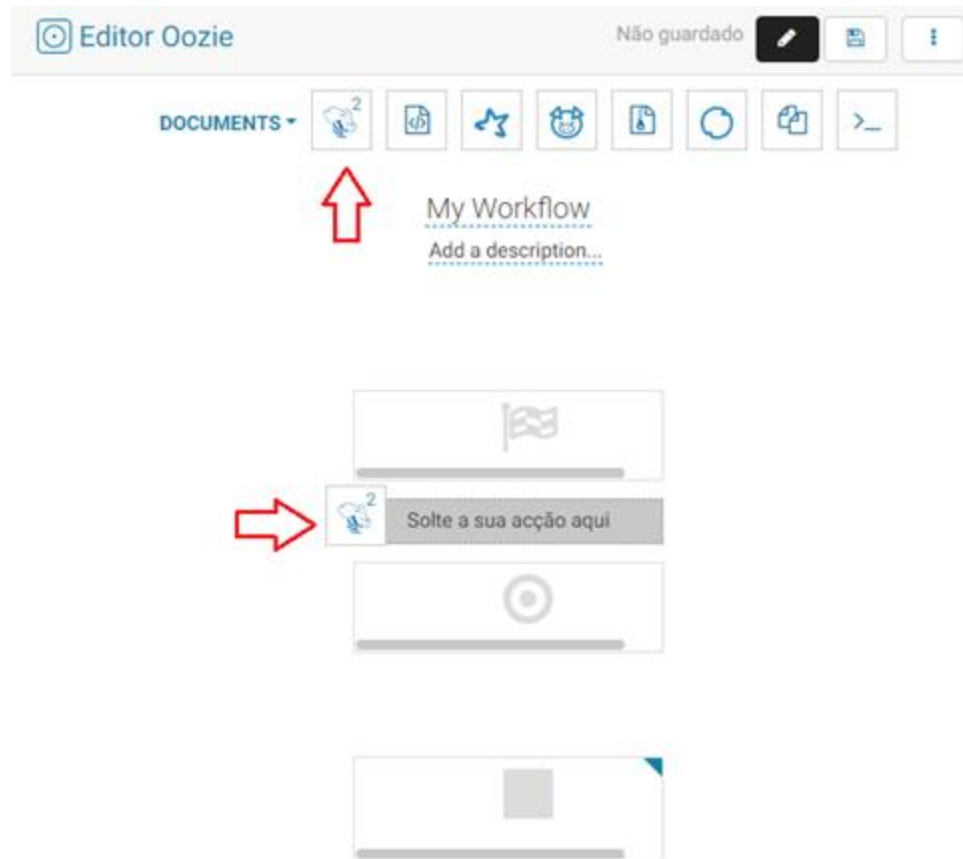
O Hue já trará uma estrutura básica de workflow com o ponto de início do mesmo, o ponto de finalização em caso de sucesso e o ponto de parada em caso de erro. Uma lista de documentos/ações que podemos implementar no fluxo também nos é apresentada.



Hands on

Workflow na prática

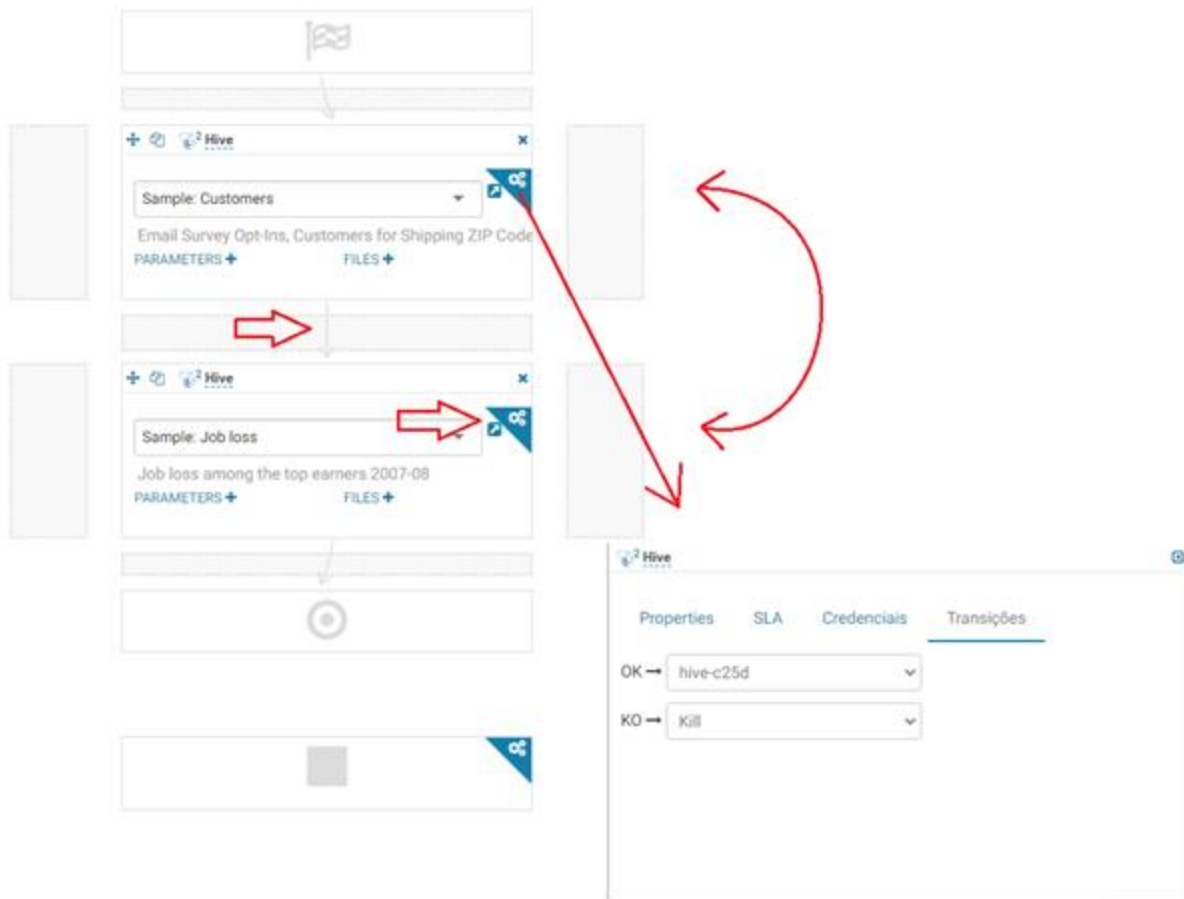
O framework trabalha com o esquema drag & drop. Você seleciona a ação que deseja e a arrasta para o local do workflow em que ela deve ser executada.



Hands on

Workflow na prática

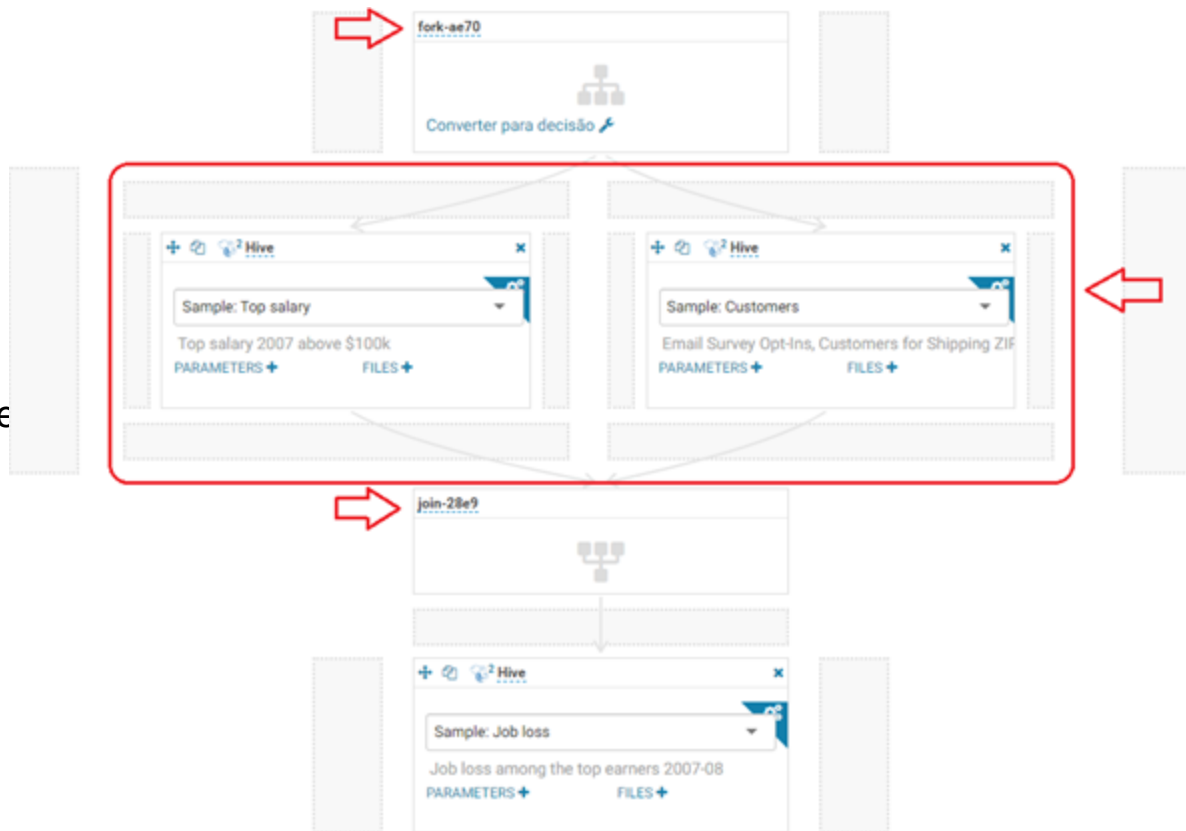
As dependências são criadas de acordo com o local em que se “solta” a ação que será executada, e pode ser vista clicando no ícone de engrenagens de cada etapa do fluxo, na aba *Transições*.



Hands on

Workflow na prática

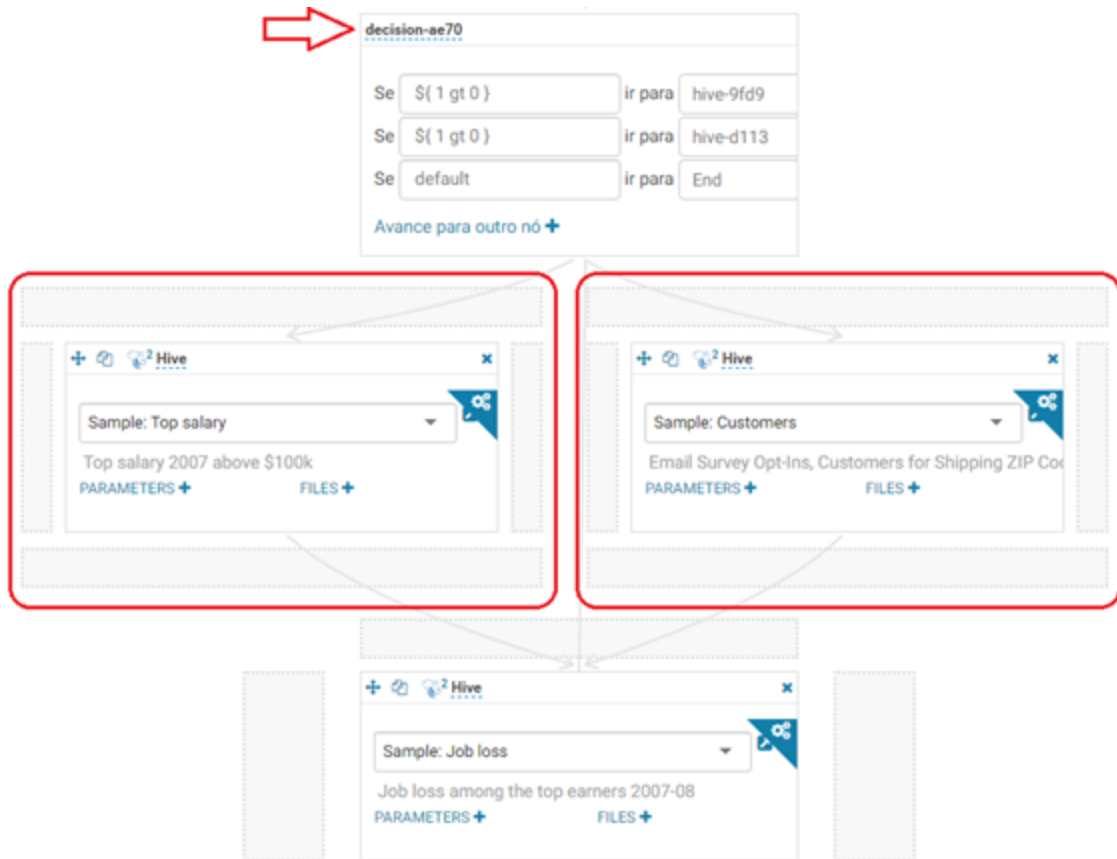
Quando houver necessidade de se executar duas ou mais ações em paralelo, ou que necessite de uma ação de decisão, deve-se arrastar a ação a ser executada para caixa lateral à ação já incluída no workflow. Isso fará com que sejam geradas duas novas etapas no fluxo, no caso de paralelismo serão geradas etapas de *fork* e *join*, unicamente para controle do fluxo.



Hands on

Workflow na prática

No caso de condicional, é gerada uma etapa *decision*, onde o usuário pode definir em quais situações o fluxo deve seguir determinado caminho.



Hands on

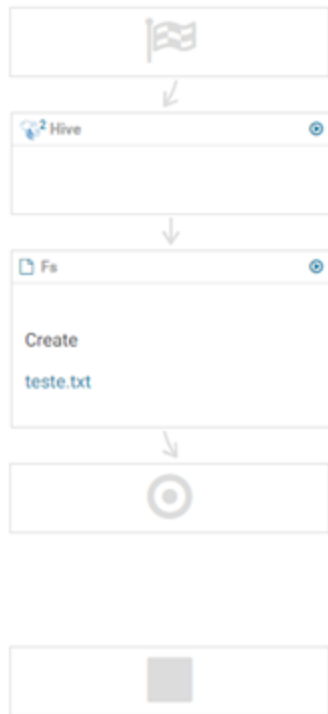
Workflow na prática

Para executar pontualmente o fluxo criado, deve-se primeiramente salvá-lo caso haja alterações recentes e clicar no botão play.

Editor Oozie



Teste



Hands on

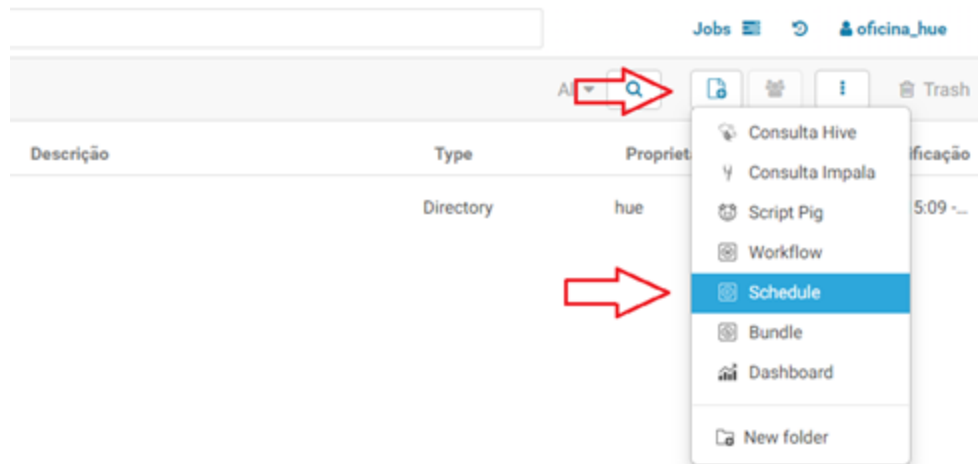
Scheduler

Após a criação de seus Workflows oozie, você pode ainda realizar agendamento de suas execuções utilizando o programa Schedule (Oozie Coordinator).

Através do Schedule poderá agendar a execução de seus fluxos com base em tempo, eventos ou dados.

Para criação de um novo fluxo, deve-se acessar o Menu principal, Browser Documents.

No canto superior direito irá encontrar um botão para criação de novos documentos, e nesta lista a opção schedule.

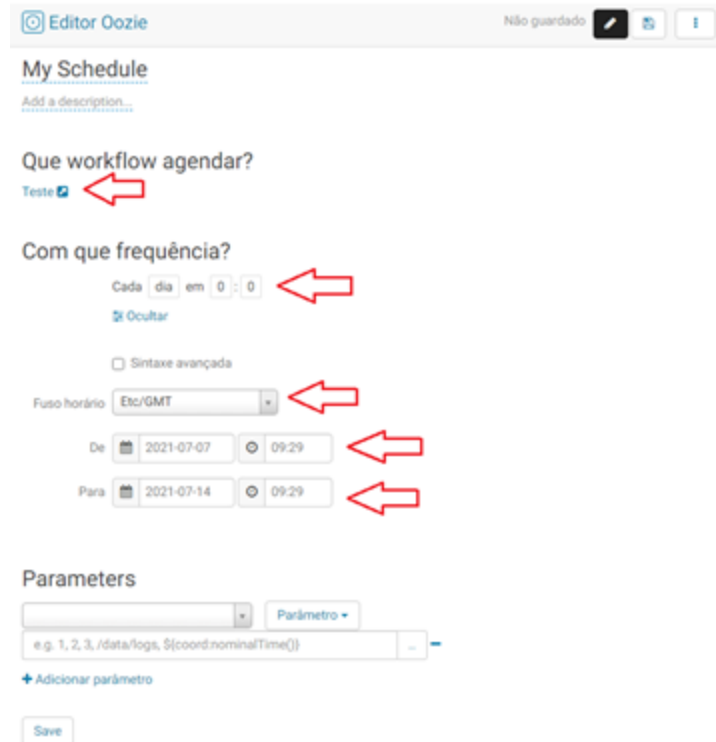


Hands on

Scheduler

Os dados básicos a serem fornecidos para criação do agendamento são:

- Workflow a ser agendado
- Frequência de execuções: Definição de qual período será utilizado para definição das execuções. ex: Fluxo deve ser executado diariamente às 12 e as 00 horas.
- Timezone– Timezone do aplicativo schedule. Fuso horário que deverá ser utilizado como parâmetro para as execuções do fluxo.
- Período de existência do schedule: Por qual período (data/hora início e fim) o schedule deverá existir. ex: Schedule de 01/01/2021 10:00 até 21/12/2021 10:00.



Editor Oozie Não guardado

My Schedule

Add a description...

Que workflow agendar?

Teste

Com que frequência?

Cada dia em 0:0

Ocultar

☐ Sintaxe avançada

Fuso horário: Etc/GMT

De: 2021-07-07 09:29

Para: 2021-07-14 09:29

Parameters

Parâmetro

e.g. 1, 2, 3, /data/logs, \$[coord.nominalTime()]

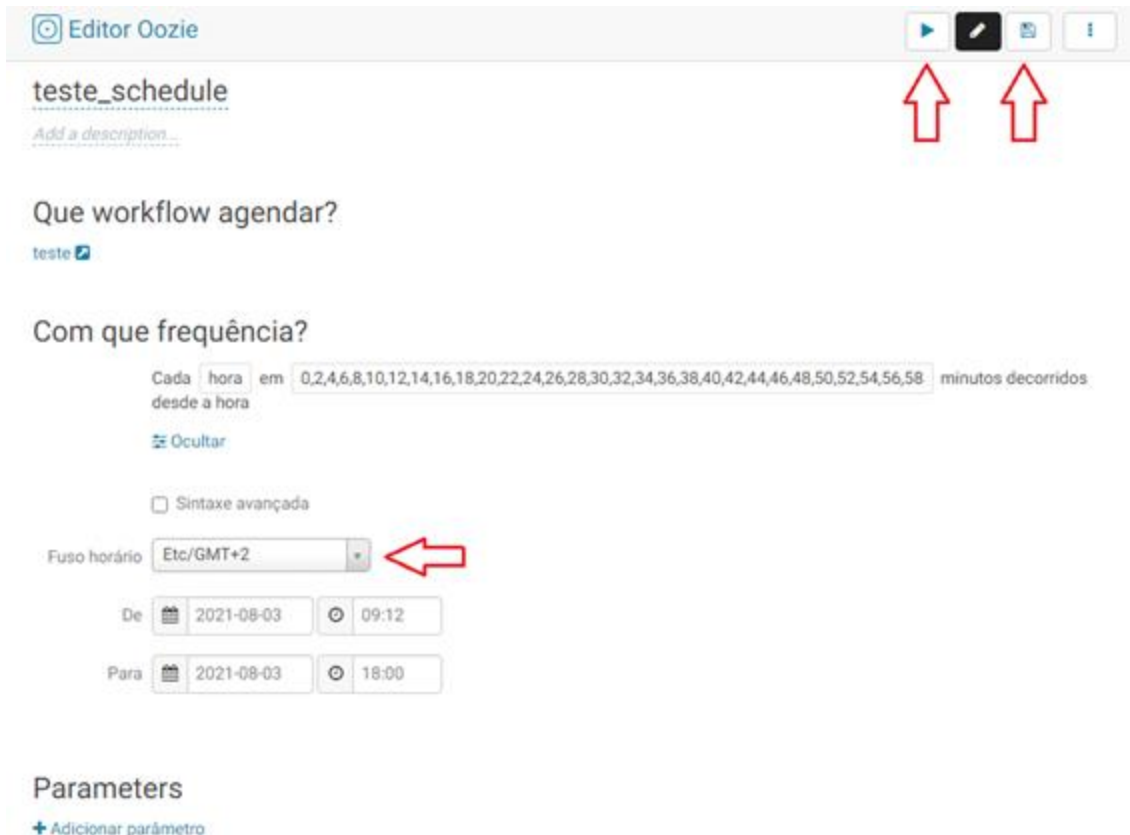
+ Adicionar parâmetro

Save

Hands on

Scheduler

Após salvar o schedule, para que o mesmo seja executado deve-se clicar no botão Play, localizado no canto superior direito da janela. Dessa forma o scheduler irá criar as instâncias de execução de acordo com as definições feitas pelo usuário no momento da criação do scheduler.



The screenshot shows the Oozie Editor interface. At the top, there is a toolbar with icons for Play, Edit, Save, and Help. Two red arrows point to the Play and Edit icons. Below the toolbar, the workflow name 'teste_schedule' is displayed, followed by a description field. The 'Que workflow agendar?' section shows 'teste' as the selected workflow. The 'Com que frequência?' section includes a dropdown for 'Cada hora em' with a list of intervals (0,2,4,6,8,10,12,14,16,18,20,22,24,26,28,30,32,34,36,38,40,42,44,46,48,50,52,54,56,58) and a 'minutos decorridos' field. There is an 'Ocultar' button. The 'Sintaxe avançada' checkbox is unchecked. The 'Fuso horário' dropdown is set to 'Etc/GMT+2', with a red arrow pointing to it. The 'De' and 'Para' date and time pickers are set to '2021-08-03' and '09:12' to '18:00' respectively. The 'Parameters' section at the bottom has an 'Adicionar parâmetro' button.

Hands on

Scheduler

Podemos ver na imagem seguinte o scheduler já iniciado, mostrando o percentual de execução do mesmo, quando foi enviado para execução e o status de execução de cada instância gerada. Podemos observar que foram geradas neste caso 12 instâncias de execução, mas reparem que ele já coloca como próxima execução o primeiro horário subsequente à última instância gerada.

teste_schedule

ID: 0000808-210710155507184-oozie-oozi-C

DOCUMENT: teste_schedule

TYPE: schedule

ESTADO: WAITING

USER: oficina_hue

PROGRESSO: 4%

ENVIADO: 03 Aug 2021 09:12:00

NEXT RUN: Tue, 03 Aug 2021 09:36:00

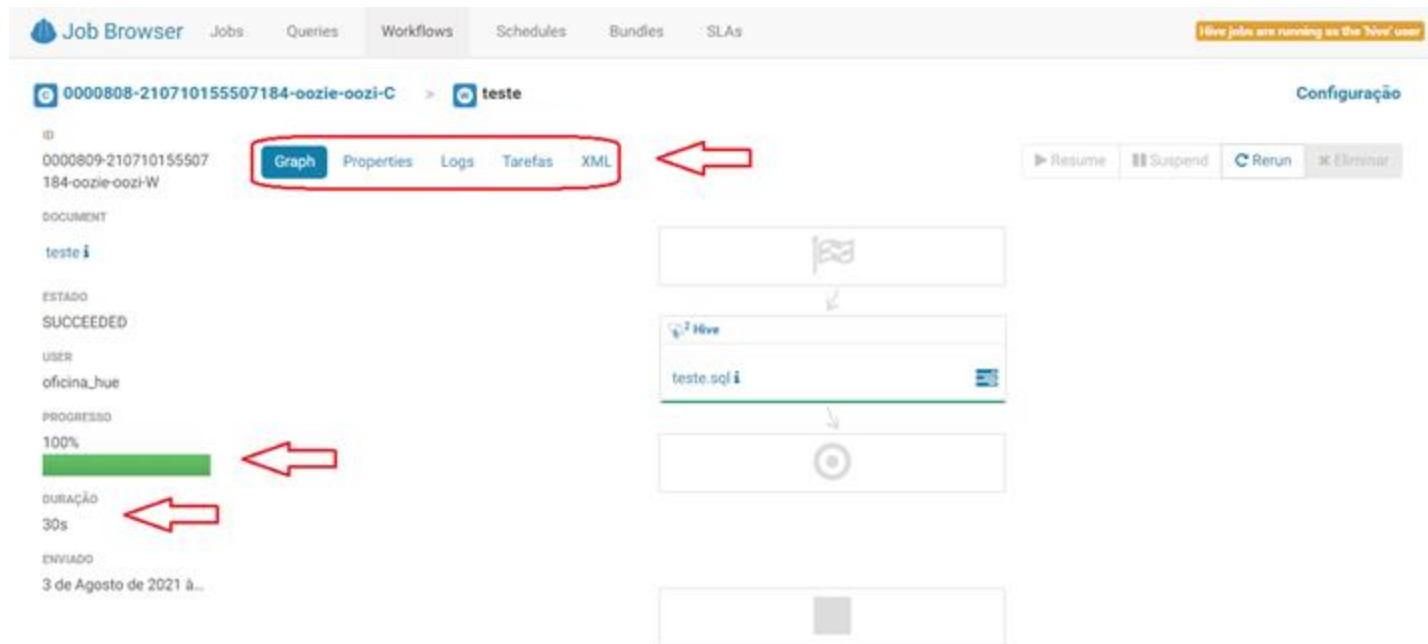
TOTAL ACTIONS: 12

Estado	Title	type	errorMessage	missingDependencies	number	errorCode	externalId	id	lastModifiedTime
WAITING	12-03 Aug 2021 09:34:00	schedule-task			12	-		0000808-210710155507184-oozie-oozi-Cj12	Tue, 03 Aug 2021 09:11:51
WAITING	11-03 Aug 2021 09:32:00	schedule-task			11	-		0000808-210710155507184-oozie-oozi-Cj11	Tue, 03 Aug 2021 09:11:51
WAITING	10-03 Aug 2021 09:30:00	schedule-task			10	-		0000808-210710155507184-oozie-oozi-Cj10	Tue, 03 Aug 2021 09:11:51
WAITING	9-03 Aug 2021 09:28:00	schedule-task			9	-		0000808-210710155507184-oozie-oozi-Cj9	Tue, 03 Aug 2021 09:11:51
WAITING	8-03 Aug 2021 09:26:00	schedule-task			8	-		0000808-210710155507184-oozie-oozi-Cj8	Tue, 03 Aug 2021 09:11:51
WAITING	7-03 Aug 2021 09:24:00	schedule-task			7	-		0000808-210710155507184-oozie-oozi-Cj7	Tue, 03 Aug 2021 09:11:51
WAITING	6-03 Aug 2021 09:22:00	schedule-task			6	-		0000808-210710155507184-oozie-oozi-Cj6	Tue, 03 Aug 2021 09:11:51
WAITING	5-03 Aug 2021 09:20:00	schedule-task			5	-		0000808-210710155507184-oozie-oozi-Cj5	Tue, 03 Aug 2021 09:11:51
WAITING	4-03 Aug 2021 09:18:00	schedule-task			4	-		0000808-210710155507184-oozie-oozi-Cj4	Tue, 03 Aug 2021 09:11:51
WAITING	3-03 Aug 2021 09:16:00	schedule-task			3	-		0000808-210710155507184-oozie-oozi-Cj3	Tue, 03 Aug 2021 09:11:51
WAITING	2-03 Aug 2021 09:14:00	schedule-task			2	-		0000808-210710155507184-oozie-oozi-Cj2	Tue, 03 Aug 2021 09:11:51
RUNNING	1-03 Aug 2021 09:12:00	schedule-task			1		0000809-210710155507184-oozie-oozi-W	0000808-210710155507184-oozie-oozi-Cj1	Tue, 03 Aug 2021 09:12:00

Hands on

Scheduler

Durante ou mesmo após a execução o usuário pode clicar em uma instância para verificar maiores detalhes da execução.



The screenshot displays the Oozie Job Browser interface. At the top, there is a navigation bar with tabs for Jobs, Queries, Workflows, Schedules, Bundles, and SLAs. A notification on the right states "1 Hive jobs are running as the 'hive' user". Below the navigation bar, the selected job is identified as "0000808-210710155507184-oozie-oozi-C" with a sub-tab "teste".

On the left side, the job details are listed:

- ID: 0000809-210710155507184-oozie-oozi-W
- DOCUMENT: teste
- ESTADO: SUCCEEDED
- USER: oficina_hue
- PROGRESSO: 100% (indicated by a green progress bar)
- DURAÇÃO: 30s
- ENVIADO: 3 de Agosto de 2021 à...

At the top of the details section, there are tabs for Graph, Properties, Logs, Tarefas, and XML. The "Graph" tab is highlighted with a red box and a red arrow pointing to it. To the right of these tabs are action buttons: Resume, Suspend, Rerun, and Eliminar.

On the right side, the workflow graph is visible, showing a sequence of steps: a start node, a "Hive" step labeled "teste.sql", and an end node. A red arrow points from the "Graph" tab to the workflow graph.

Realizar atividades localizadas no documento **Exercícios**, subtítulo **Criação de Workflow e Schedule**, enviado por email no início da Oficina.



Dúvidas ?



Instrutores

Bruno Mariz

bruno.oliveira@serpro.gov.br

Anderson Ceo


anderson.ceo@serpro.gov.br

SUPAI - Janeiro/2023

 /serprobrasil

 @serprobrasil

 @serpro

 /serpro

 serpro.gov.br