

Visual Flow

User Guide

September 2022

Version 1.0.0

Document Revisions

Date	Version	Document Changes
12/08/2020	0.9.2	Initial Draft
22/04/2021	0.9.3	Pipeline Operators
26/04/2021	0.9.4	Job Operators
07/05/2021	0.9.5	Project Name, Project Operations
25/05/2021	0.9.6	Project Name in document
09/07/2021	0.9.7	Pipeline Operators, Job Operations, Storages
24/10/2021	0.9.8	Jobs and Pipelines statuses, Custom container, Storages
13/11/2021	0.9.9	New Data Storages
22/11/2021	0.9.10	New Data Storage
15/12/2021	0.9.11	Logs levels
29/12/2021	0.9.12	Cache stage, Job Operations
21/01/2022	0.9.13	New Data Storage, Stage descriptions
04/02/2022	0.9.14	Description of STDOUT storage
18/02/2022	0.9.15	Truncate mode
04/03/2022	0.9.16	Full query mode in Transformer stage
18/03/2022	0.9.17	Sort stage
01/04/2022	0.9.18	Password component
08/08/2022	0.9.19	Verification to ensure up-to-date condition
19/08/2022	0.9.19	Added Slice Stage and updates to 4.2.1 and 5.2
27/09/2022	1.0.0	Updates for release 19.09.2022 and new design

Revision Details

Version	Updates	Performed by
0.9.19	P.3.1 split into 3.1. Getting Started and 3.2. Create Project	Alesya Lemeshevskaya
0.9.19	P."4.2.. Create a job" updated to "4.2 .Create a job. Available Stages" and broken down to multiple sections per each stage	Alesya Lemeshevskaya
0.9.19	Added section 3 for Connections to p."3.3. Manage Project Settings"	Alesya Lemeshevskaya
0.9.19	Updated screenshots for Jobs list and Pipelines list	Alesya Lemeshevskaya
0.9.19	Style and verification updates to all chapters	Alesya Lemeshevskaya
0.9.19	Added Slice Stage and updates to p.4.2.1. Read Stage and p.5.2. Create a pipeline	Alesya Lemeshevskaya

1.0.0	Added User profile menu to p.3.1. Getting started	Alesya Lemeshevskaya
1.0.0	Updated all screens as per new design	Alesya Lemeshevskaya
1.0.0	Added Job History to jobs overview actions	Alesya Lemeshevskaya
1.0.0	Added chapter 5.5. Scheduling a pipeline	Alesya Lemeshevskaya
1.0.0	Added about Job History to p.4.4. Job Execution	Alesya Lemeshevskaya
1.0.0	Added about copying stages on canvas to p.4.3. Job Designer Functions Overview	Alesya Lemeshevskaya
1.0.0	Updates to p.4.2.1. Read Stage. Added statement about selecting existing connection	Alesya Lemeshevskaya

Table of Contents

1	Introduction.....	5
	1.1 ...Terminology.....	5
	1.2....Scope and Purpose.....	6
	1.3 ...Process Overview.....	6
2	Roles and Authorizations.....	7
3	Project Operations.....	8
	3.1 ...Getting started.....	8
	3.2 ...Create Project.....	9
	3.3 ...Project Overview.....	11
	3.4 ...Manage Project Settings.....	12
4	Job Operations	14
	4.1 ...Jobs Overview.....	14
	4.2 ...Create a Job. Available Stages.....	15
	4.2.1..Read Stage.....	16
	4.2.2. Write Stage.....	20
	4.2.3. Group By Stage.....	23
	4.2.4. Remove Duplicates Stage.....	23
	4.2.5. Filter Stage.....	24
	4.2.6. Transformer Stage.....	24
	4.2.7. Sort Stage.....	25
	4.2.7. Slice Stage.....	25
	4.2.9. Join Stage.....	25
	4.2.10. Change Data Capture Stage.....	26
	4.2.11. Union Stage.....	26
	4.2.12. Cache Stage.....	27
	4.3 ...Job Designer Functions Overview.....	28
	4.4 ...Job Execution.....	28
5	Pipeline Operations.....	30
	5.1 ...Pipelines Overview.....	30
	5.2 ...Create a Pipeline.....	31
	5.3 ...Pipeline Designer Functions Overview.....	36
	5.4 ...Pipeline Execution.....	37
	5.5 ...Scheduling a pipeline.....	38

1. Introduction

1.1. Terminology

ETL is an abbreviation for *extract, transform, load*, three database functions combined into one tool to pull data out of one database, transform it and place it into another database.

- **Extract** is the process of *reading data* from a database. In this stage the data is collected, often from multiple and different types of sources.
- **Transform** is the process of *converting the extracted data* from its previous form into the form needed to place it into another database.
- **Load** is the process of *writing the data* into the target database.

Job is a chain of individual stages linked together. It describes the flow of data from a data source to a data target. Usually, a stage has a minimum of one data input and/or one data output. However, some stages can accept more than one data input and output to more than one stage.

In Visual Flow the available stages are:

- Read
- Write
- Group By
- Remove duplicated
- Filter
- Transformer
- Sort
- Slice
- Join
- Union
- Change Data Capture
- Cache

Pipeline is a compound of multiple jobs and can be run. In Visual Flow, user can use such stages as:

- Job
- Pipeline
- Container
- Notification

1.2. Scope and Purpose

Visual Flow web application is the ETL tool designed for effective data manipulation via convenient and user-friendly interface. The tool has the following capabilities:

- Can integrate data from heterogeneous sources:
 - ✓ AWS S3
 - ✓ DB2
 - ✓ Cassandra
 - ✓ Elasticsearch
 - ✓ IBM COS
 - ✓ Mongo
 - ✓ MSSQL
 - ✓ MySQL
 - ✓ Oracle
 - ✓ PostgreSQL
 - ✓ Redis
 - ✓ Redshift
- Leverage direct connectivity to enterprise applications as sources and targets
- Perform data processing and transformation
- Leverage metadata for analysis and maintenance

1.3. Process Overview

Visual Flow jobs and pipelines exist within a certain namespace (project) so the first step in the application would be to create a project or enter an existing project. Then user needs to enter Job Designer to create a job.

Job Designer is a graphical design interface used to create, maintain, execute and analyze jobs. Each job determines the data sources, the required transformations and destination of the data.

Pipeline designer is a graphical design interface aimed for managing pipelines. Designing a pipeline is similar to designing a job.

Important: When editing stages in the Configuration Panel, to save data user must click *Confirm* button, otherwise the data is lost.

Visual Flow key functions include but not limited to:

- ✓ Create project which serves as a namespace for jobs and/or pipelines
- ✓ Manage project settings
- ✓ User access management
- ✓ Create/maintain a job in Job Designer
- ✓ Job execution and logs analysis
- ✓ Create/maintain a pipeline in Pipeline Designer
- ✓ Pipeline execution
- ✓ Import/Export jobs and pipelines

2. Roles and authorizations

The following roles are available in the application:

- ✓ Viewer
- ✓ Operator
- ✓ Editor
- ✓ Administrator

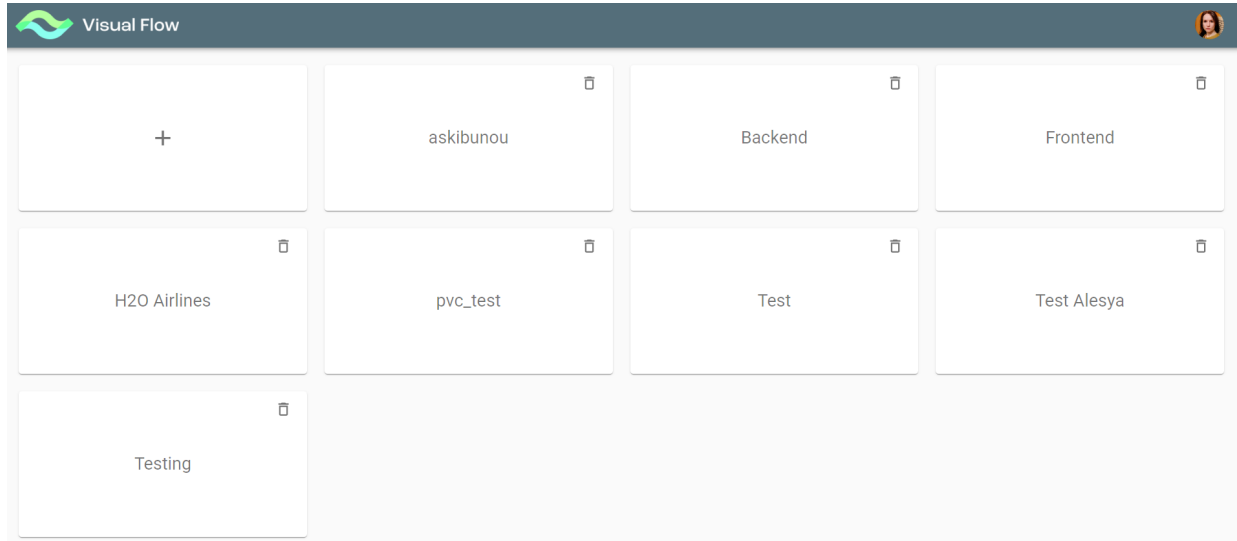
They can perform the below operations within the namespaces they are authorized to. Only a Super-admin user can create a workspace (project) and grant access to this project.

Role	Actions		
	Project Settings	Jobs	Pipelines
Viewer	View all	View all	View all
Operator	View all	View all / execute jobs	View all / execute pipelines
Editor	Edit all but Users and Roles	Edit / execute jobs	Edit / execute pipelines
Admin	Edit all	Edit / execute jobs	Edit / execute pipelines

3. Project Operations.

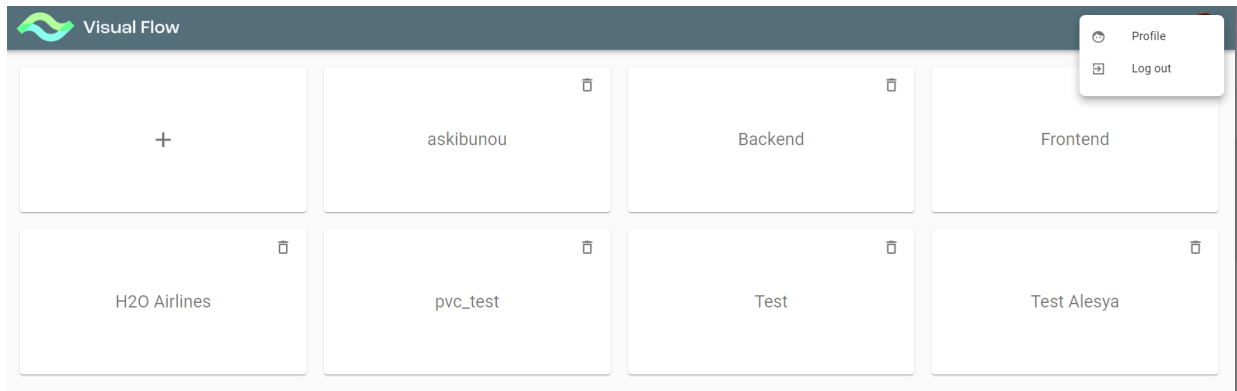
3.1. Getting started

Once you first logged on to the application you see the initial screen with all existing projects:



If you are not authorized to a certain project it is locked for you. So you see lock icon on its tile. Please contact project owners to get access to their projects.

If you click on user icon in top right corner you get to user profile menu:

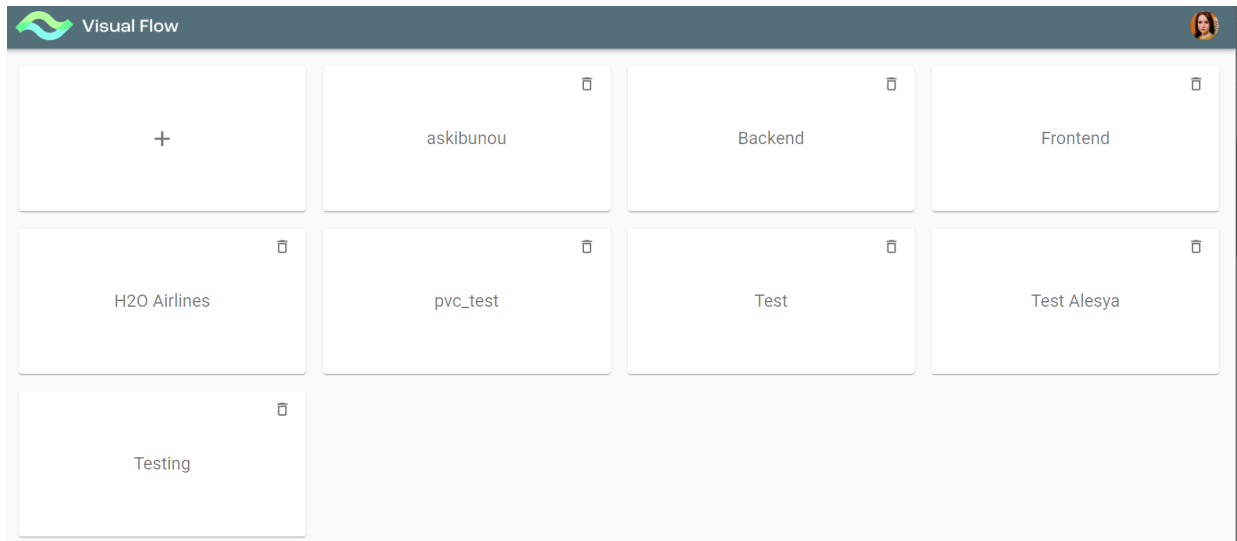


Here you can view your user profile or log out.

3.2. Create a Project

To create a project you need to push “+” button.

Note: this is the action of super-admin user only. The button is not visible for the application roles (Viewer, Operator, Editor, Admin).

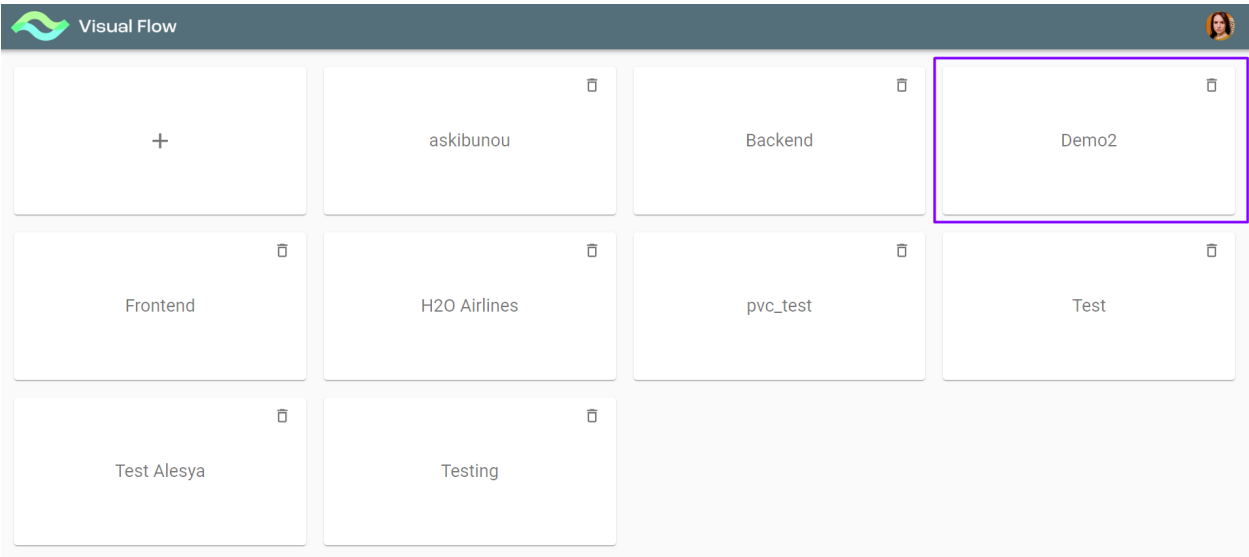


With “+” button pushed you get to *Create Project Form* to enter project basic settings:

- Project Name
- Project Description
- Requests (CPU/Memory)
- Limits (CPU/Memory)

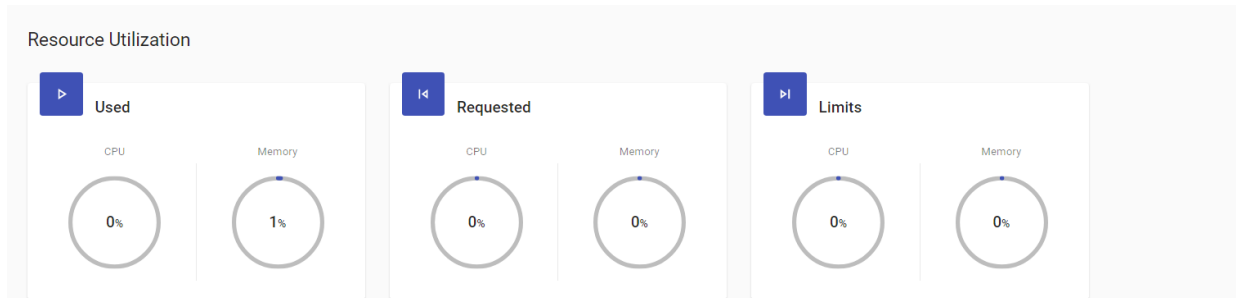
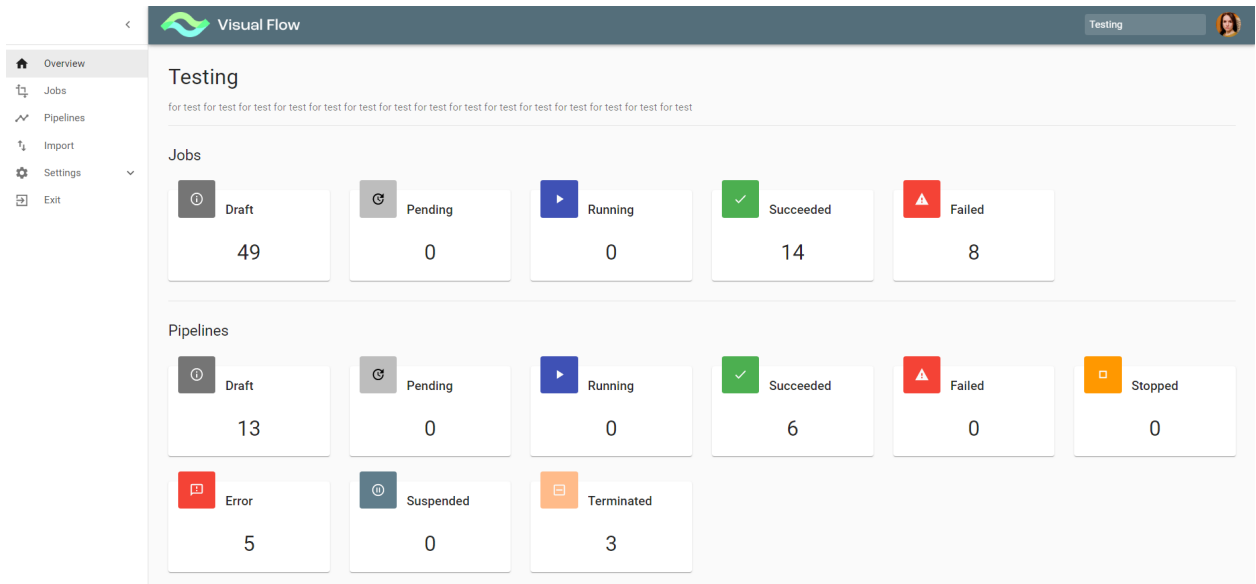
After saving *Create Project Form* the project is created under the given name and then can be found on

initial screen:



3.3. Project Overview

The screen contains project left menu and displays information about the project jobs, pipelines and their resource utilization (applicable for running jobs).



3.4. Manage Project Settings

Settings submenu contains:

- Basic
- Parameters
- Connections
- Users and Roles

1) *Basic* is already there after project creation. *Edit* button turns on the edit mode for updates.

The screenshot shows the 'Visual Flow' application interface. On the left is a sidebar menu with options: Overview, Jobs, Pipelines, Import, Settings (highlighted), Basic, Parameters, Connections, Users/Roles, and Exit. The main content area is titled 'View Project' and contains a form for project configuration. The form includes fields for 'Name' (set to 'Demo2'), 'Description', 'Requests' (with sub-fields for CPU and Memory, both set to 1), and 'Limits' (also with sub-fields for CPU and Memory, both set to 1). An 'Edit' button is visible in the top right corner of the form.

2) *Parameters* serve to store values required for the entire project, e.g. JDBC connection, DB2 credentials or table schema can be the same for multiple jobs within a project and therefore stored at the project level. *Edit* button turns on the edit mode for updates.

The screenshot shows the 'View Project Parameters' page. It features a search bar at the top. Below it is a table of parameters. The parameters are listed in two columns: the parameter name and its value. The values are masked with asterisks for sensitive information like passwords and access keys.

Parameter Name	Value
123	*****
Access_key_aws	3U0AQ6AT58UEHFCT0BQ5
Bucket_aws	visualflow-gomel
ENT_PWD	*****
ENT_URL	jdbc:postgresql://172.28.0.80:5432/enterprisedb_dev
ENT_USR	migration_dev

3) *Connections* option enables user to manage connections to a storage.

Here you see a list of all existing connections with their name/storage type and available actions (view, edit, delete, ping). Also you can create a new connection with *Create Connection* button.

Manage Connections

Q Search...

Storage

CREATE CONNECTION

Storage	Name	
DB2	db2	<div><div></div><div></div><div></div><div></div></div>
Storage	Name	
Cassandra	dd	<div><div></div><div></div><div></div><div></div></div>
Storage	Name	
Mongo	ddd	<div><div></div><div></div><div></div><div></div></div>
Storage	Name	
Cassandra	gg	<div><div></div><div></div><div></div><div></div></div>
Storage	Name	
	rrrrr	<div><div></div><div></div><div></div><div></div></div>
Storage	Name	
MySQL	test	<div><div></div><div></div><div></div><div></div></div>
Storage	Name	
DB2	test_1	<div><div></div><div></div><div></div><div></div></div>

4) *User and Roles* is meant for user access management or to view user access depending on your authorization.

Users cannot change their own roles, this operation can be done by *Admin* or *Super-admin* only.

If you try to change your role you will get error message: "You cannot change your role".

Edit button and therefore Edit mode is only available for admin within the project or super-admin.

Users and Roles

Q Search...

ADD USER

<input type="checkbox"/> ID	Name	Role
<input type="checkbox"/> ashubenok	ashubenok-gomel-iba-by	vf-admin
<input type="checkbox"/> aapanasko	Alina Apanasko	vf-viewer
<input type="checkbox"/> agsadovsky	Alexander Sadovsky	vf-viewer
<input type="checkbox"/> alemeshevskaya-ibagroup-eu	alemeshevskaya-ibagroup-eu alemeshevskaya-ibagroup-eu	vf-viewer


Rows per page: 5 1-4 of 4

SAVE

CANCEL

4. Job operations

4.1. Jobs Overview

Clicking *Jobs* menu item leads to *Jobs Overview Screen* which allows you to see a list of jobs existing within a project. Some of the jobs can be used in pipelines, this is indicated by  icon. Jobs Overview Screen displays the following information:

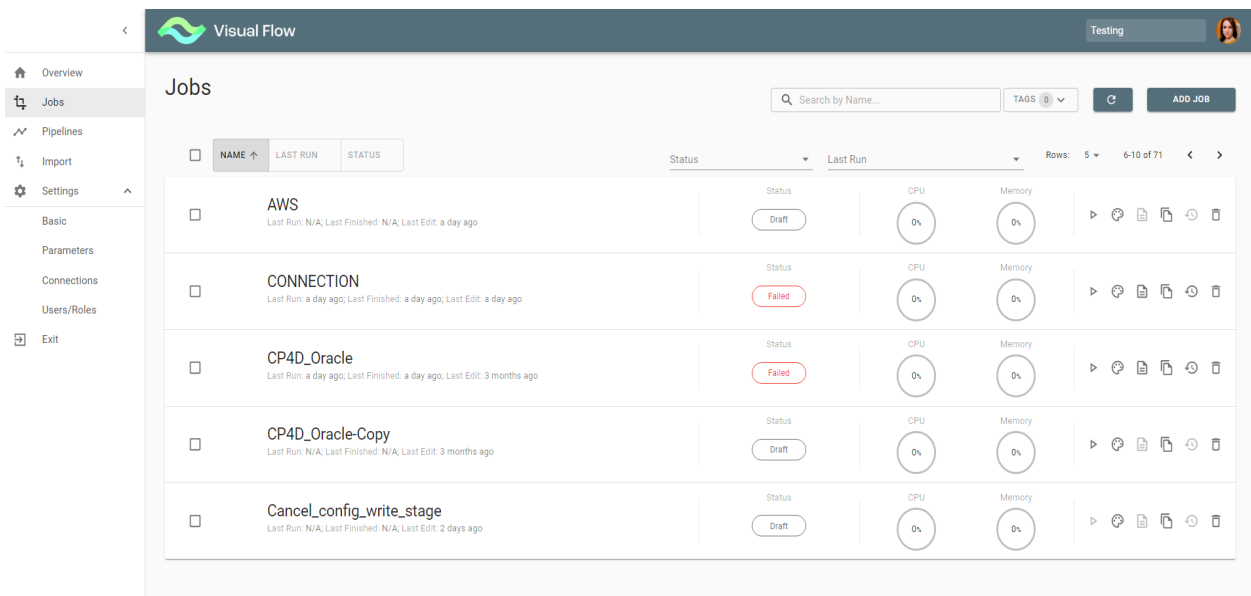
- Job Name
- Job Last run/Last finished/Last edit
- Job Status
- Resource Utilization (CPU/Memory)
- Available Actions (Run/Job Designer/Logs/Copy/Job History/ Delete)

Job has a certain status at various phases of execution:

- Draft
- Pending
- Running
- Succeeded
- Failed
- Unknown (This status appears very rarely in case of undefined error)

Notes:

- The actions availability and therefore visibility is depending on user authorizations
- The user cannot delete a job which compounds a pipeline



The screenshot displays the Visual Flow Jobs Overview screen. The interface includes a sidebar with navigation options: Overview, Jobs, Pipelines, Import, Settings, Basic, Parameters, Connections, Users/Roles, and Exit. The main content area shows a table of jobs with columns for Name, Last Run, Status, CPU, and Memory. The jobs listed are AWS (Draft), CONNECTION (Failed), CP4D_Oracle (Failed), CP4D_Oracle-Copy (Draft), and Cancel_config_write_stage (Draft). The top bar shows the Visual Flow logo, a search bar, and buttons for TAGS, C, and ADD JOB.

NAME	LAST RUN	STATUS	CPU	MEMORY	ACTIONS
AWS	Last Run: N/A; Last Finished: N/A; Last Edit: a day ago	Draft	0%	0%	Run, Copy, Job History, Delete
CONNECTION	Last Run: a day ago; Last Finished: a day ago; Last Edit: a day ago	Failed	0%	0%	Run, Copy, Job History, Delete
CP4D_Oracle	Last Run: a day ago; Last Finished: a day ago; Last Edit: 3 months ago	Failed	0%	0%	Run, Copy, Job History, Delete
CP4D_Oracle-Copy	Last Run: N/A; Last Finished: N/A; Last Edit: 3 months ago	Draft	0%	0%	Run, Copy, Job History, Delete
Cancel_config_write_stage	Last Run: N/A; Last Finished: N/A; Last Edit: 2 days ago	Draft	0%	0%	Run, Copy, Job History, Delete

4.2. Create a Job. Available Stages.

With *Add Job* button pushed user gets to *Job Designer* for creating a new job.

The screenshot shows the Job Designer interface. At the top, an orange header bar contains a back arrow and the text "Please enter name and save params". Below the header, the interface is divided into two main sections: a left configuration panel and a right main canvas. The left panel has a "PARAMS" tab selected. It contains several input fields: "Name" (with "Job_1" entered), "Tags", "Resources" (with a sub-section for "Driver Request Cores" set to "0.1", "Driver Cores" set to "1", and "Driver Memory" set to "1 GB"), "Executor Request Cores" (set to "0.1"), "Executor Cores" (set to "1"), "Executor Memory" (set to "1 GB"), "Executor Instances" (set to "2"), and "Shuffle Partitions" (set to "10"). At the bottom of the left panel are "CONFIRM" and "DISCARD" buttons. The right main canvas is currently empty. A green circular arrow icon is visible on the left side of the canvas, pointing towards the configuration panel.

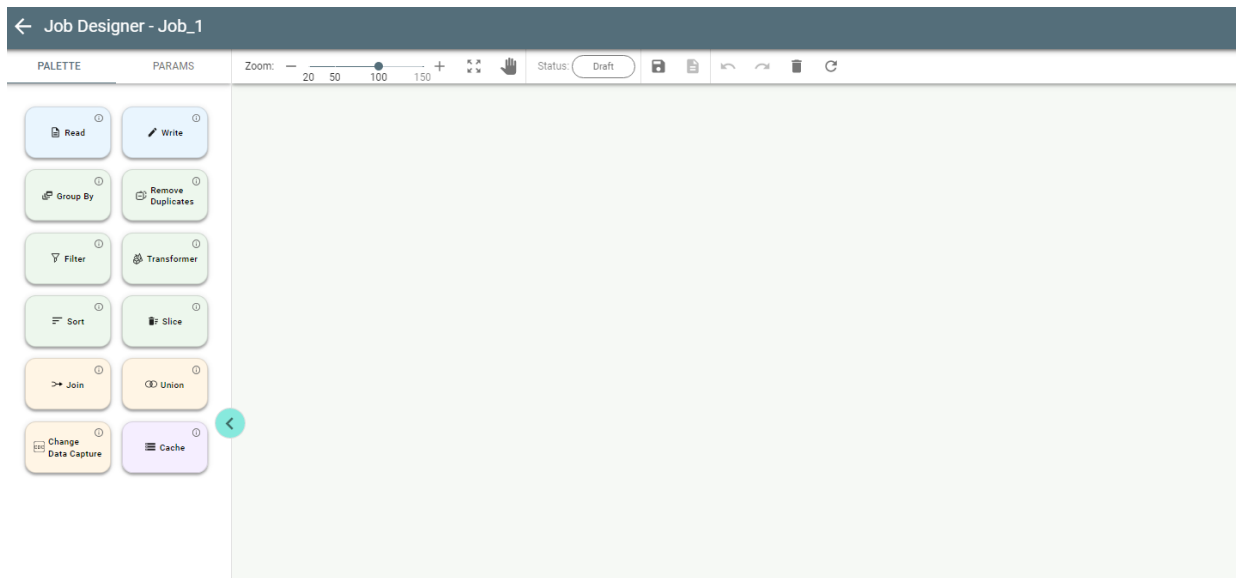
On the left configuration panel a name for the job must be provided. Tags can be used to classify your job.

Update parameters or keep their default values and then push *Confirm* on the panel:

This screenshot is similar to the previous one, but with an additional warning message at the top. The orange header bar now contains the text: "Your job has not been saved. If you leave Job Designer unsaved changes will be lost. Press 'SAVE' to save your job". To the right of this text are "SAVE" and "EXIT" buttons. The configuration panel on the left remains the same, with "Job_1" in the Name field and various resource parameters. The main canvas on the right is still empty. A green circular arrow icon is visible on the left side of the canvas, pointing towards the configuration panel.

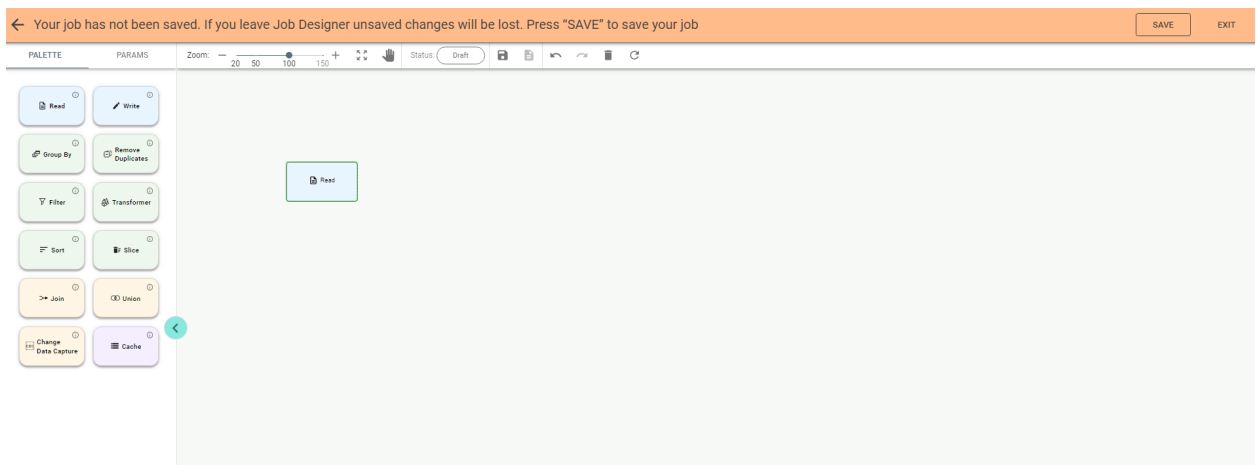
Save the job by pushing *Save* button on the *Job Designer* header.

Now you see *Palette* tab then with all available stages:

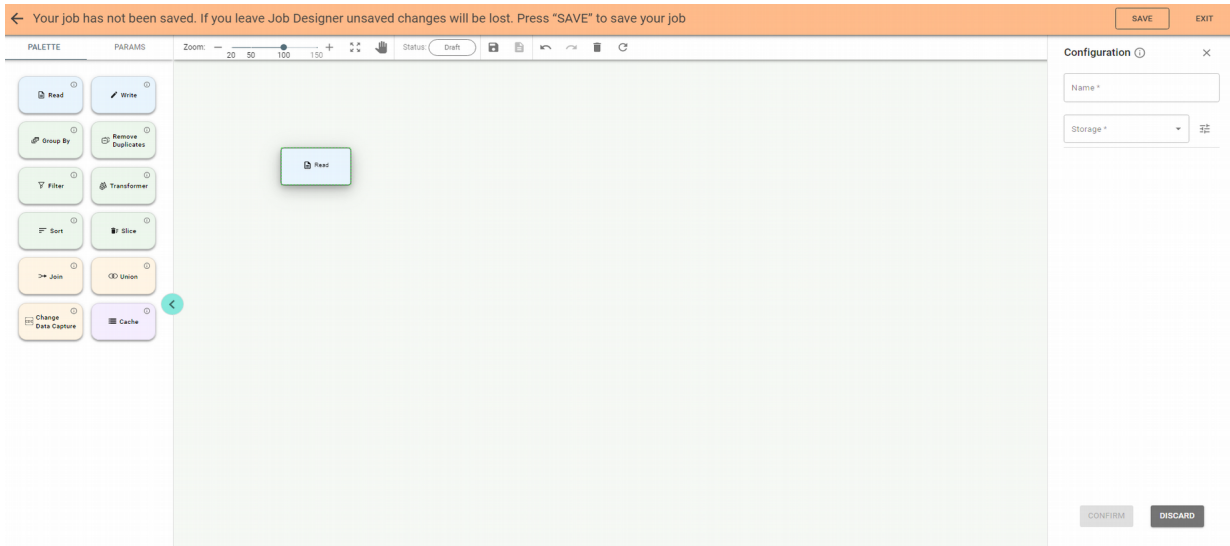


4.2.1. Read Stage.

You can start creating a job by dragging a stage to the canvas, e.g. you can drag *Read* stage:



Note: you can also add a stage by double clicking its tile on the palette.
Double-click on a stage on canvas opens the configuration panel on the right:



Enter a name for the stage and select *Storage* DB2 if you want to read data from DB2 table.
Available *Storage* values for Read stage are:

- ✓ AWS S3
- ✓ DB2
- ✓ Cassandra
- ✓ Elasticsearch
- ✓ IBM COS
- ✓ Mongo
- ✓ MySQL
- ✓ MSSQL
- ✓ Oracle
- ✓ PostgreSQL
- ✓ Redis
- ✓ Redshift
- ✓ Redshift-jdbc
- ✓ Local File

Important: you can select an existing connection with Params button near Storage field:

If you do so its parameters e.g. JDBC URL, user etc are filled automatically.

But now suppose you don't have previously created connection so fill required parameters for DB2 *Storage* manually:

Configuration ⓘ

×

Name *

Read_tab1

Storage *

DB2

⌵

⌵

JDBC URL

jdbc:db2://10.224.0.52

⌵

×

User

db2inst1

⌵

×

Password

.....

👁

⌵

×

Custom SQL

False

⌵

×

Schema

db2inst1

⌵

×

Table

user_tab1

⌵

×

CertData (optional)

⌵

×

CONFIRM

DISCARD

Important: you can pick up a parameter value with *Parameters* ⌵ button on the right panel if it is previously created as a project parameter.

Configuration ⓘ

×

Name *

Read_tab1

Storage *

DB2

⌵

⌵

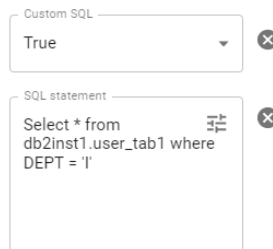
JDBC URL

jdbc:db2://10.224.0.52

⌵

×

Read Stage for DB2 storage has an option *Custom SQL* to read data with SQL statement (e.g. *select * from table where field = value*). If you select *Custom SQL - True* you need to enter the SQL statement and specify schema:



Custom SQL

True

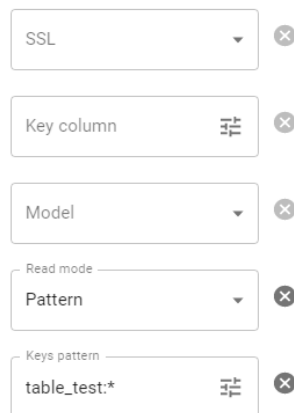
SQL statement

```
Select * from
db2inst1.user_tab1 where
DEPT = 'I'
```

For *Redis* source you need to define *Key column*, *Model*, *SSL*, *Read mode*, *Keys pattern* or *Table* fields in the Configuration of Read stage.

- ✓ *Key column*. Column name to store the hash key.
- ✓ *Model (binary, hash)* defines the *Redis* model used to persist DataFrame. By default it is *hash*.
- ✓ *Read mode* (key, pattern) defines the way the read operation is handled. If *Key* is selected, then the read is done based on *table* field. In case of *Pattern* the provided pattern (option *Keys Pattern*) dictates what Redis keys to be read.

Keys pattern. If pattern ends with * (e.g., “table: *”), all keys from the pattern are read. If one pattern is defined (e.g. “table: first value”) then only one key is read.



SSL

Key column

Model

Read mode

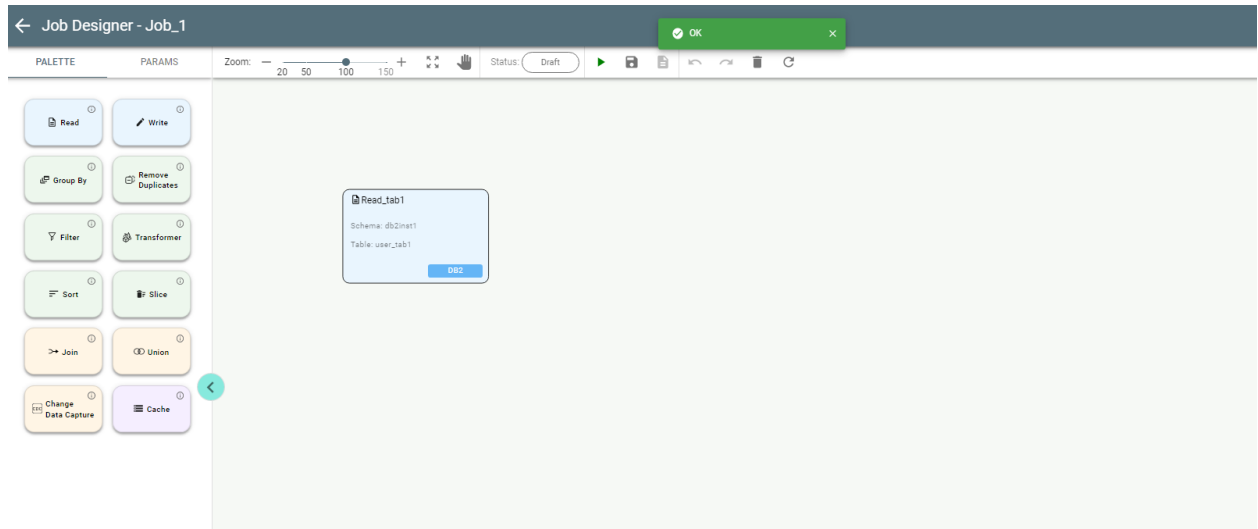
Pattern

Keys pattern

table_test:*

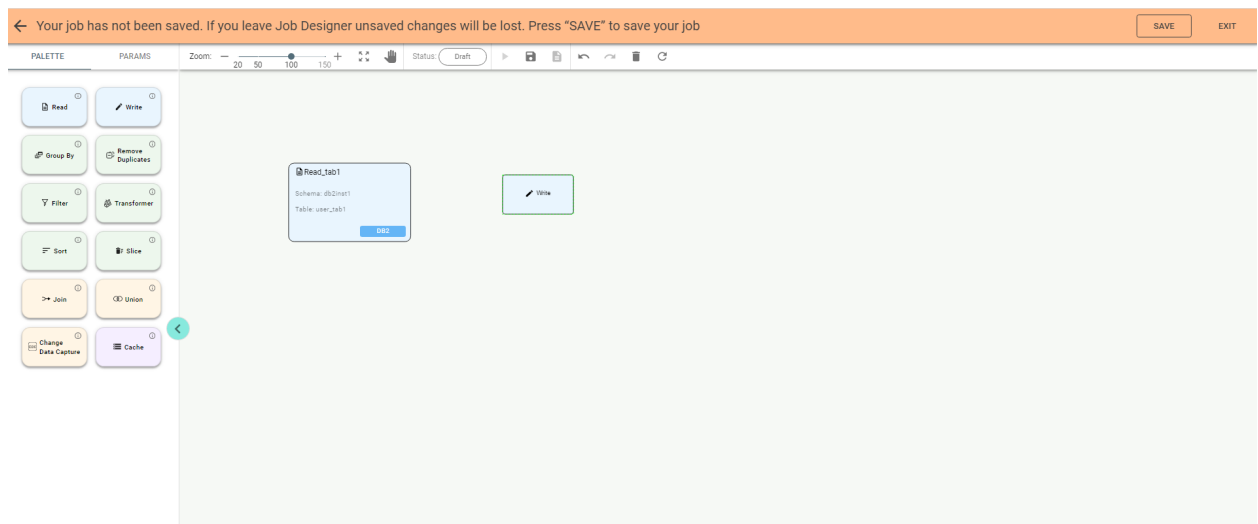
Save the stage by pushing *Confirm* button on the configuration panel. Push *Save* button on the header If you want to save the job at this step.

When the first stage of the job is configured the canvas looks like this:



4.2.2. Write Stage.

Now drag another stage, e.g. *Write* stage:



Enter a name for the stage and select *Storage IBM COS* if you want to post data from DB2 table to Cloud Object Storage file. Fill required parameters for IBM COS *Storage*.

Available *Storage* values for Write stage are:

- ✓ AWS S3
- ✓ DB2
- ✓ Cassandra
- ✓ Elasticsearch
- ✓ IBM COS
- ✓ Mongo
- ✓ MSSQL
- ✓ MySQL
- ✓ Oracle

- ✓ PostgreSQL
- ✓ Redis
- ✓ Redshift
- ✓ Redshift-jdbc
- ✓ Local File
- ✓ STDOUT

IBM COS Storage has two options of *Authentication type*: *HMAC* and *IAM*.

If *HMAC* is selected you should fill *accessKey* and *secretKey*.

If *IAM* is selected *iamApiKey* and *iamServiceId* should be entered.

For *IBM COS* and *AWS S3* storages function *Partition By* can be used in *Write* stage. It partitions the output on the file system by given columns. If specified, the output is laid out on the file system similar to Hive's partitioning scheme.

As an example, when we partition a dataset by year and then month, the directory layout looks like:

- year=2016/month=01/

- year=2016/month=02/

In a case of importing table data with *Write* stage to *Cassandra* source from another storage the table layout for output must be previously created in *Cassandra*. Columns, key fields, data types of the fields must be specified in the table.

Important: All the above points must match the imported table.

If the column names have uppercase characters in the imported table when data is output to *Cassandra*, the job will fail as in *Cassandra* column names are stored in lowercase only. The issue can be resolved with *Transformer* stage.

The output recorded to *STDOUT* storage can be seen in Logs. The number of records to be shown in Logs can be specified in *Quantity* field. Available range is from 1 to 2147483631 records.

For *Redis* source in *Write* stage the following fields must be specified: *Key column*, *Model*, *SSL*, *TTL*, *Table* and *Write mode*.

✓ *Key column*. For writing it specifies the unique column used as Redis key. By default the key is auto-generated.

✓ *TTL*. Data expiration time in seconds. Data doesn't expire if TTL is negative or 0. By default it is 0. Positive value of *TTL* means number of seconds in which the data will be removed.

Key column

NAME

✕

SSL

False

▼

Model

Hash

▼

Table

table_test

✕

ttl

0

Write mode

Overwrite

▼

Important:

Write mode field defines how data is posted to its destination. Available values are:

- ✓ Overwrite
- ✓ Append
- ✓ Error if Exists

With *Overwrite* *Write mode Truncate mode* can be used for DB2, Oracle, MySQL, PostgreSQL, MSSQL, Redshift:

- ✓ *None*. No truncation occurs, but the target table is deleted and recreated. Note that all the indexes, constraints, etc. defined for this table will be lost.
- ✓ *Simple*. The standard truncation that deletes data from the target table but keeps the indexes, constraints and other modifiers intact. However note that if the target table has a primary key referenced as a foreign key in other tables, the truncation will fail.

To resolve this either use Cascade mode instead or drop constraints manually (outside of VF) prior to accessing the table with VF.

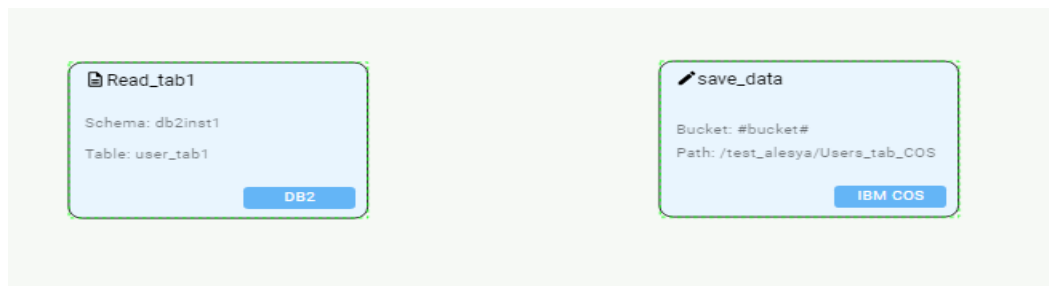
- ✓ *Cascade* (only for *Oracle* and *PostgreSQL*). The cascade truncation does not only delete the data from the target table, but also from other tables that use the target table's primary key as a foreign key constraint.

File format is to choose a format of destination file. Available formats are:

- ✓ CSV
- ✓ JSON
- ✓ Parquet
- ✓ ORC
- ✓ Text
- ✓ Avro

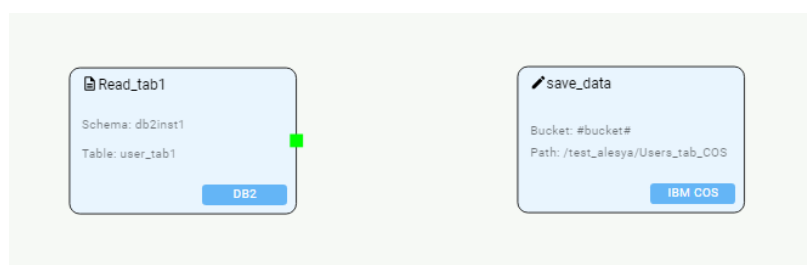
Confirm the stage by pushing *Confirm* on the panel.

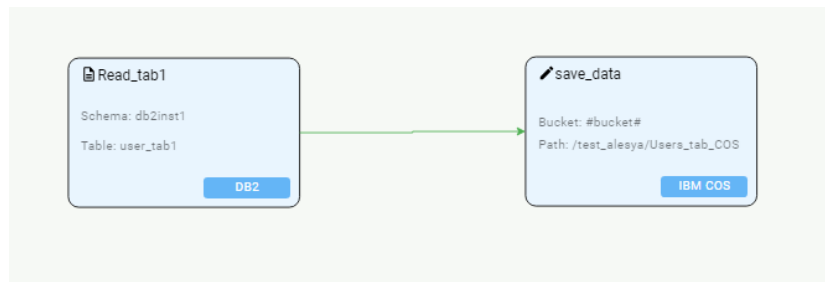
Now there are two stages to connect to each other.



Important:

To connect stages, hover his mouse on a stage edge until user sees a green rectangle. Click it and drag it to the border of another stage and its green rectangle. When you reach it a green arrow should appear.





Other stages available:

- ✓ Group By
- ✓ Remove duplicates
- ✓ Filter
- ✓ Transformer
- ✓ Sort
- ✓ Join
- ✓ Change data capture
- ✓ Union
- ✓ Cache

4.2.3. Group By Stage.

The stage allows grouping by column which must be specified in the Configuration panel. There is an option *Drop grouping columns* for removing grouping columns from the output. Also aggregate function can be added e.g., Count, Avg etc.

Configuration ⓘ

Name *
Dept I

☐ Drop grouping columns

Group By *
DEPT

Aggregate functions +
Count Column NAME

4.2.4. Remove Duplicates Stage.

Specify a key column for the operation. To specify more than one key, use comma or Enter. For Order By operation you need to specify column to sort by and sort order Asc or Desc. Default is Asc.

Configuration ⓘ

×

Name *

Remove_dups

Key *

Order By

+

CONFIRM

DISCARD

4.2.5. Filter Stage.

Enter any boolean expression. Two or more expressions can be combined using logical operators (AND, OR).

Examples: (column1 < 10) and (column2 between 10 and 25).

Configuration ⓘ

×

Name *

only-new-data

Filter *

operation = 1

4.2.6. Transformer Stage.

Transformer stage gives you the ability to modify columns to be written to some data storage later. Transformer mode defines the type of the SQL query (Spark SQL dialect) that is accepted and executed.

Simple mode only allows you to specify the part between SELECT and FROM. You can do things like these:

- col1, concat(col1, col2)
- count(*) as count
- a, b, row_number() OVER (PARTITION BY a ORDER BY b)
- col, exists(col, x -> x % 2 == 0)
- col, collect_list(col)

Syntax: <column_name_1> as <alias_1>, function(<column_name_2>) as <alias_2>

Full SQL mode allows you to write a full-blown Spark SQL query. In this case you have to specify table name manually or reference a table name from parameter.

Table name. The name of the table that you should use within the Spark SQL query. Applicable for *Full SQL* transformer mode only.

4.2.7. Sort Stage.

There are two types of sorting available: *Full sort* or *Sort within partitions*.

Full sort sorts DataFrame by the specified column(s).

Sort within partitions sorts each DataFrame partition by the specified column(s). In this case the order of the output data is not guaranteed because the data is ordered at partition level.

Select column(s) to sort by and sort order (default value is Asc).

Available sort options:

- asc
- asc nulls first
- asc nulls last
- desc
- desc nulls first
- desc nulls last

4.2.8. Slice Stage.

Slice stage allows you to remove unnecessary columns from your data stream.

There are 2 modes of Slice stage:

- *Keep*
- *Drop*

Default mode is Drop. In this case you need to specify columns you would like to slice from your data flow.

With Keep mode selected you need to specify columns you want to keep.

4.2.9. Join Stage.

These are available types of join:

- *Inner* join. Transfers records from input data sets whose key columns contain equal values to the output data set. Records whose key columns do not contain equal values are dropped.
- *Left outer* join. Transfers all values from the left data set but transfers values from the right data set only where key columns match. The stage drops the key column from the right data set.
- *Right outer* join. Transfers all values from the right data set and transfers values from the left

data set and intermediate data sets only where key columns match. The stage drops the key column from the left and intermediate data sets.

- *Full outer* join. Transfers records in which the contents of the key columns are equal from the left and right input data sets to the output data set. It also transfers records whose key columns contain unequal values from both input data sets to the output data set.

- *Cross* join. Returns a result data set where each row from the first table is combined with each row from the second table.

Left semi join. Returns values from the left side of the relation that has a match with the right.

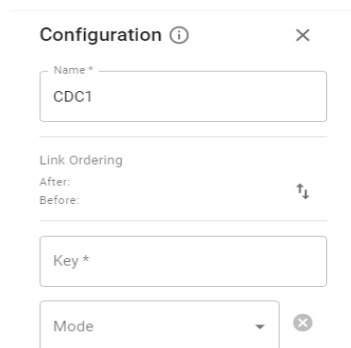
Left anti join. Returns values from the left relation that has no match with the right.

Link Ordering option allows you to specify which input link is regarded as the left link and which link is regarded as the right link. By default, the first link added is regarded as the left link, and the last one as the right link.

4.2.10. Change Data Capture Stage.

This stage is intended to find all differences between before (old) and after (new) datasets. Based on differences, CDC produces an additional column 'Operation', which indicates the state of the row from the old dataset considering it's presence/absence in the new one. CDC compares each row of the new and the old datasets based on key and columns to compare values and sets Operation value.

Note: old and new datasets must not contain duplicates (rows with the same key) based on key column(s). Old and new datasets columns to compare and key columns must be presented in both datasets with the same names. If there are duplicated rows at least in one of the dataset, the result of the CDC gets unpredictable.



The screenshot shows a 'Configuration' dialog box with a close button (X) in the top right corner. Inside the dialog, there is a 'Name *' text input field containing the text 'CDC1'. Below this is a 'Link Ordering' section with 'After:' and 'Before:' labels and a vertical double-headed arrow icon. Underneath is a 'Key *' text input field. At the bottom, there is a 'Mode' dropdown menu currently showing 'Mode' and a small 'X' icon to its right.

4.2.11. Union Stage.

You can union two datasets.

Note: Column's sequence, names, types are important for union operation.

Mode contains 2 options: *All values* and *Distinct values*.

Distinct values which is default eliminates duplicate records from the second dataset.

All values needs to be specified explicitly, and it tolerates duplicates from the second dataset.

Configuration ⓘ
×

Name *

Union_2

Mode *

All values
Distinct values

×

CONFIRM
DISCARD

4.2.12. Cache Stage.

Cache stage persists dataset in some storage. The storage type can be tweaked by specifying/combining parameters. The configuration gives you the ability to define:

- ✓ Whether to use memory.
- ✓ Whether to drop the RDD to disk if it falls out of memory.
- ✓ Whether to keep the data in memory in a serialized format.
- ✓ Whether to replicate the RDD partitions on multiple nodes.

Configuration ⓘ
×

Name *

Cache1

Use disk

True

Use memory

True

Use off heap

False

Deserialized

True

Replication

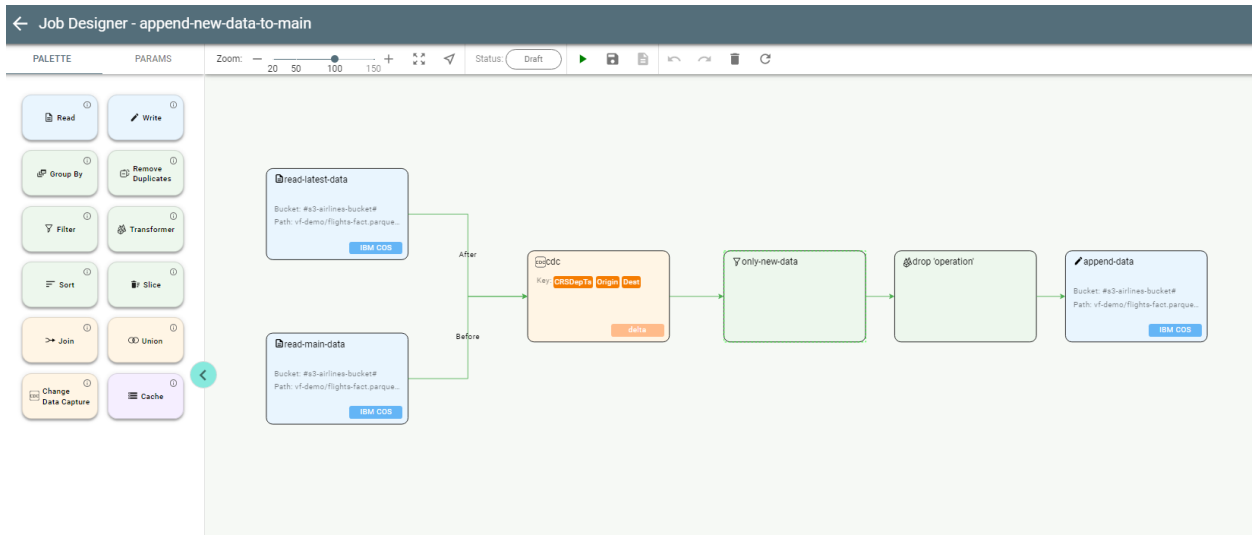
1

Save the job by pushing *Save* on the *Job Designer* header.

For newly created job as long as it is not yet run its status is *Draft*:

Status: Draft

Drag other stages according to the job flow from source to destination. See the job with more stages as example:



4.3. Job Designer functions overview

The following functions are available in *Job Designer*:

✓ Zoom operations:



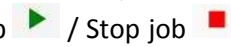
✓ Switch mode between Move and Select



✓ Show job status



✓ Run job / Stop job (for running)



✓ Save job



✓ See job logs



✓ Undo / Redo operation on canvas



✓ Remove element from canvas



✓ Refresh



Note: you can copy stages on canvas using mouse right-click menu:

✓ copy child node

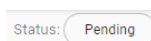
✓ paste child node

4.4. Job Execution

Push *Play* button to run the job:





Its status changes from *Draft* to *Pending*

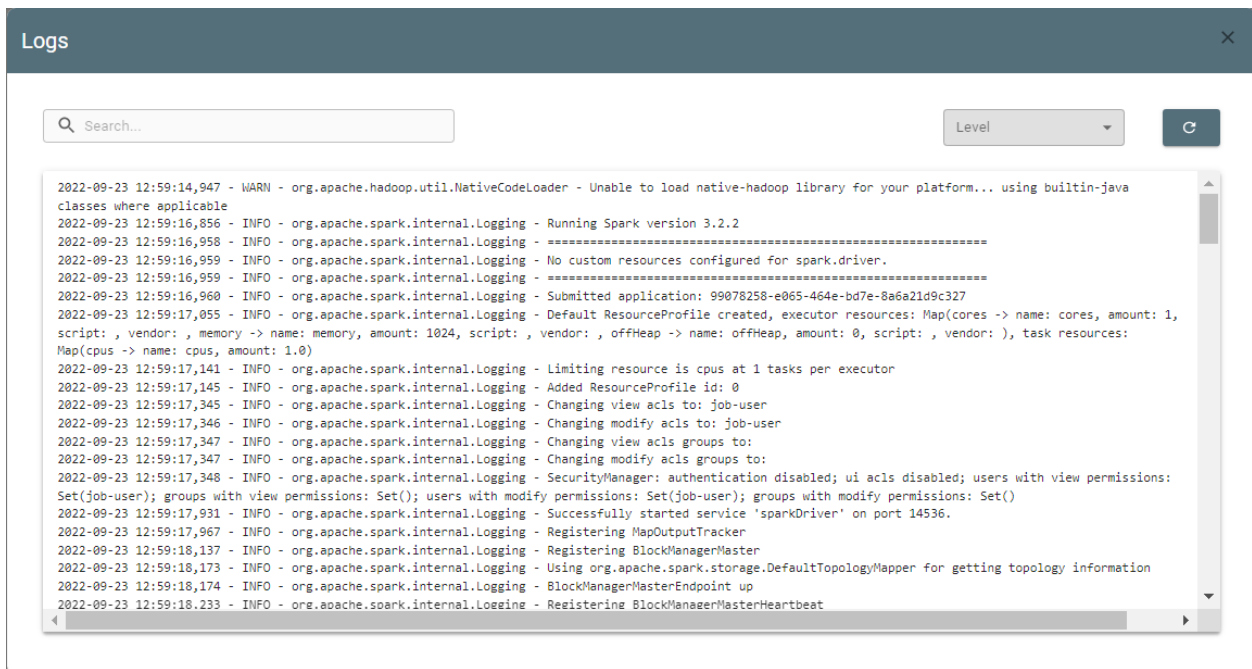


Push Refresh to update the status. It should turn to *Running*



While running, it can be interrupted with *Stop* button.  When a job is completed its status is *Succeeded or Failed*

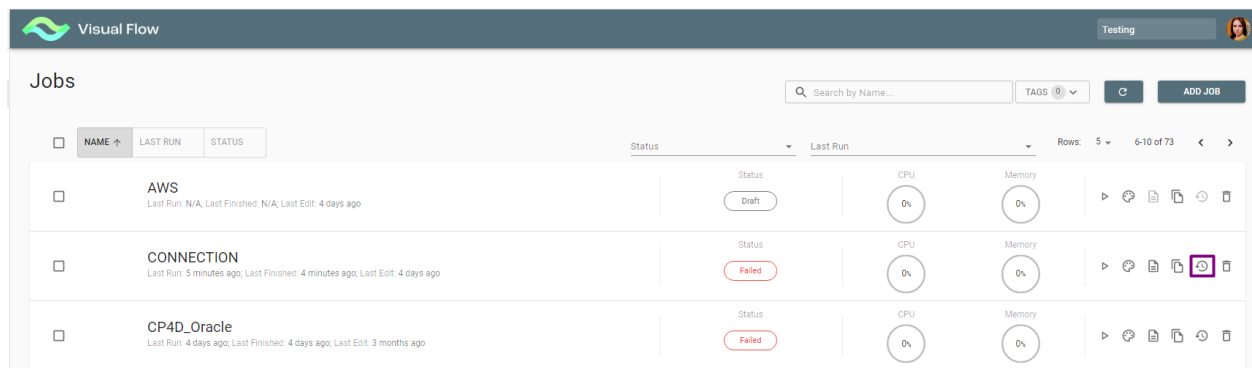
Use *Logs* button  to analyze job logs. It gets you to *Logs Screen*:



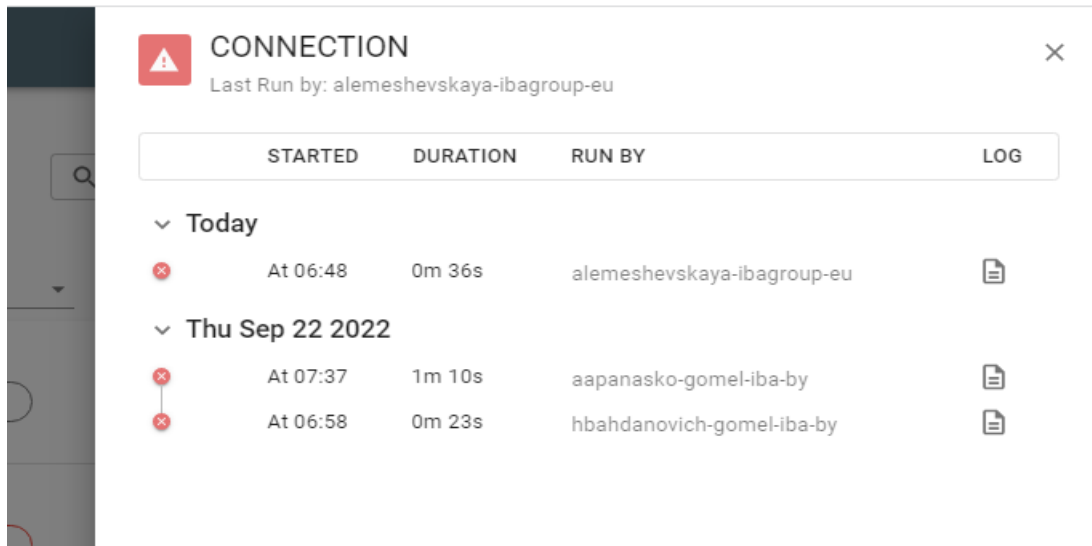
Logs Screen has several levels:




- ✓ WARNING
- ✓ INFO
- ✓ ERROR
- ✓ DEBUG
- ✓ RESULT

You can also view earlier job runs data with *Job History* button:



It takes you to *Job History* screen containing each job run data including logs:



	STARTED	DURATION	RUN BY	LOG
▼ Today	At 06:48	0m 36s	alesmeshevskaya-ibagroup-eu	
▼ Thu Sep 22 2022	At 07:37	1m 10s	aapanasko-gomel-iba-by	
	At 06:58	0m 23s	hbahdanovich-gomel-iba-by	

5. Pipeline Operations

5.1. Pipelines Overview

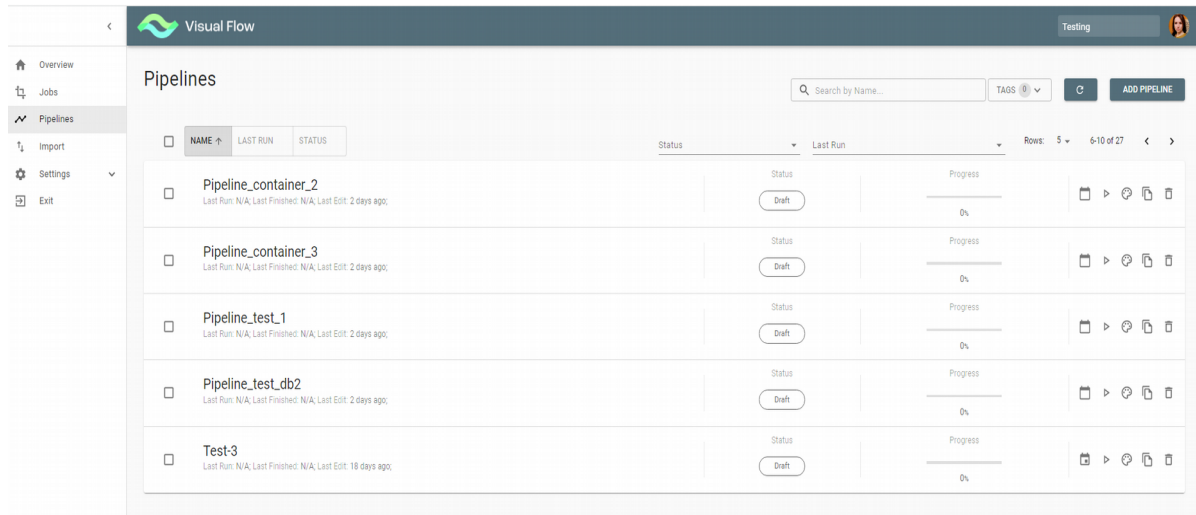
Clicking *Pipelines* menu item takes you to *Pipelines Overview Screen* which shows a list of pipelines existing within a project.

It displays the following information:

- Pipeline Name
- Checkbox for deleting/exporting multiple pipelines
- Pipeline Last run/Last finished/Last edit
- Pipeline Status
- Pipeline Progress
- Available Actions (Scheduling/Run/Pipeline Designer/Copy/Delete)

Pipeline has a certain status at various phases of execution:

- Draft
- Error (This status appears e.g. due to incorrectly entered data)
- Failed
- Pending
- Running
- Stopped
- Succeeded
- Suspended (This status can be reproduced via API)
- Terminated

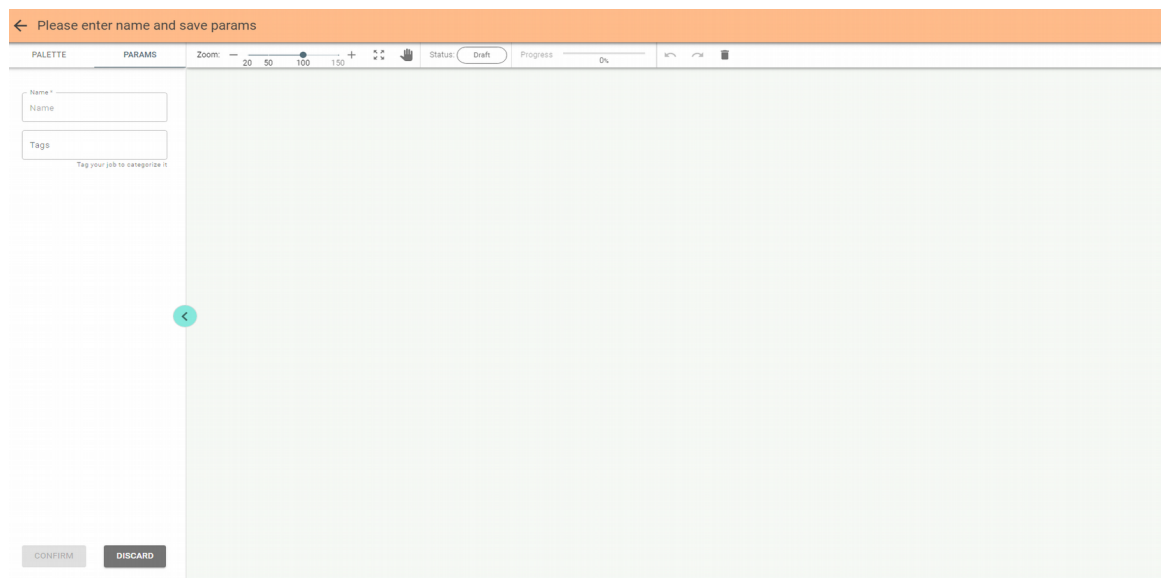


Note: the actions availability and therefore visibility is depending on user authorizations.

5.2. Create a Pipeline

With *Add Pipeline* button pushed you get to *Pipeline Designer* for creating a pipeline.

1) On the left configuration panel *Params* tab is open by default, where you can enter pipeline name and tags for the pipeline classification:



Once you added a name notifications become available for you:

← Please enter name and save params

PALETTE PARAMS Zoom: 20 50 100 150 Status: Draft Progress: 0%

Name *
Pipe1

Tags
Tag your job to categorize it

Notifications

Notify about failure ☐

Notify about success ☐

CONFIRM DISCARD

- 2) Set notifications to be notified on the pipeline events as you wish and push *Confirm* button.
- 3) Save the pipeline by pushing *Save* button on the *Pipeline Designer* header.
- 4) After saving the pipeline Palette tab with all available stages is opened by default:

← Pipeline Designer - Pipe1

PALETTE PARAMS Zoom: 20 50 100 150 Status: Draft Progress: 0%

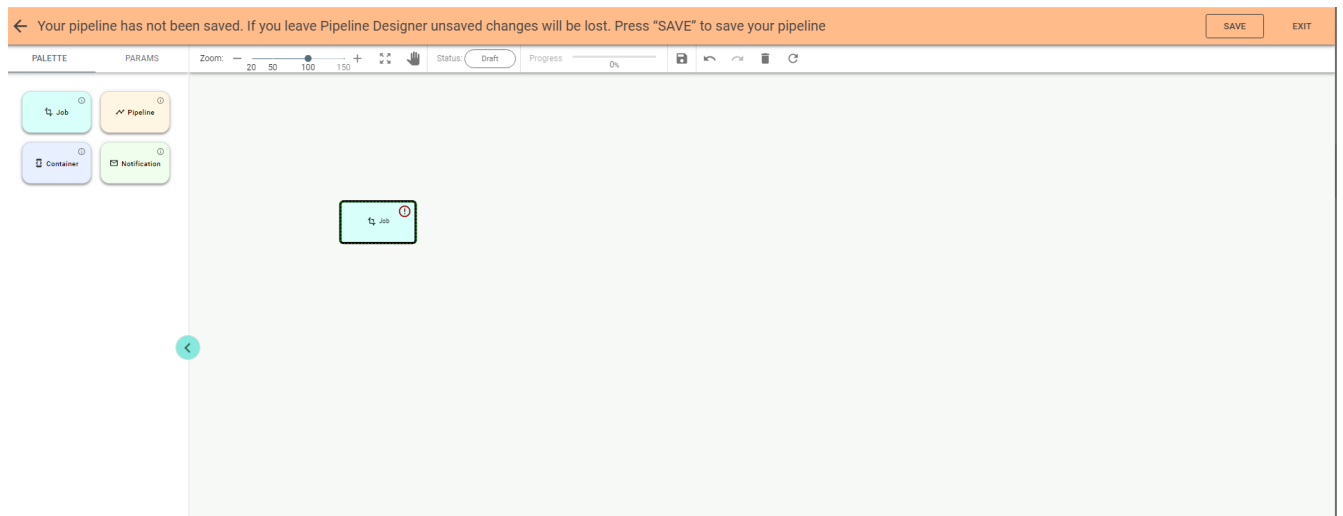
Job Pipeline

Container Notification

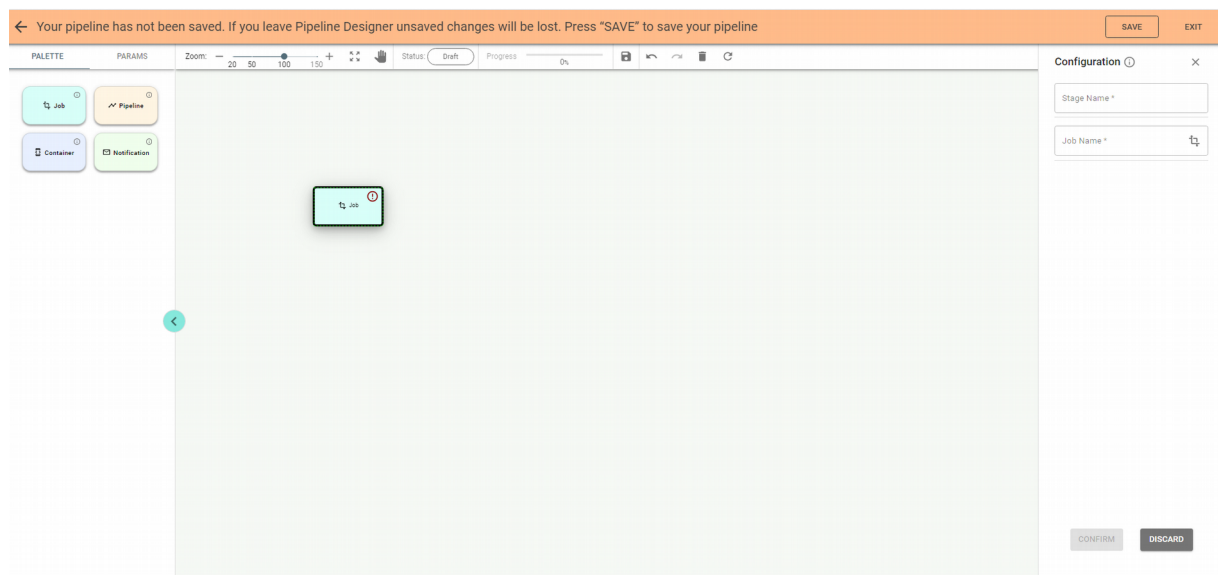
CONFIRM DISCARD


5) Pipeline is a combination of existing job stages and/or pipeline and/or notification stages and/or container stages. Notification stage is most often added to configuration to notify about job/pipeline stage failure/success.

Start creating a pipeline by dragging *Job* stage to the canvas:



6) Double-click on the stage opens the configuration panel on the right:



Enter a name for the stage and select a job from the list *with*  button.

Jobs

Search...

☐ 111111

☐ 11111test

☐ 1234

☐ 1234

☐ 1234-Copy

☐ AWS

☐ CONNECTION

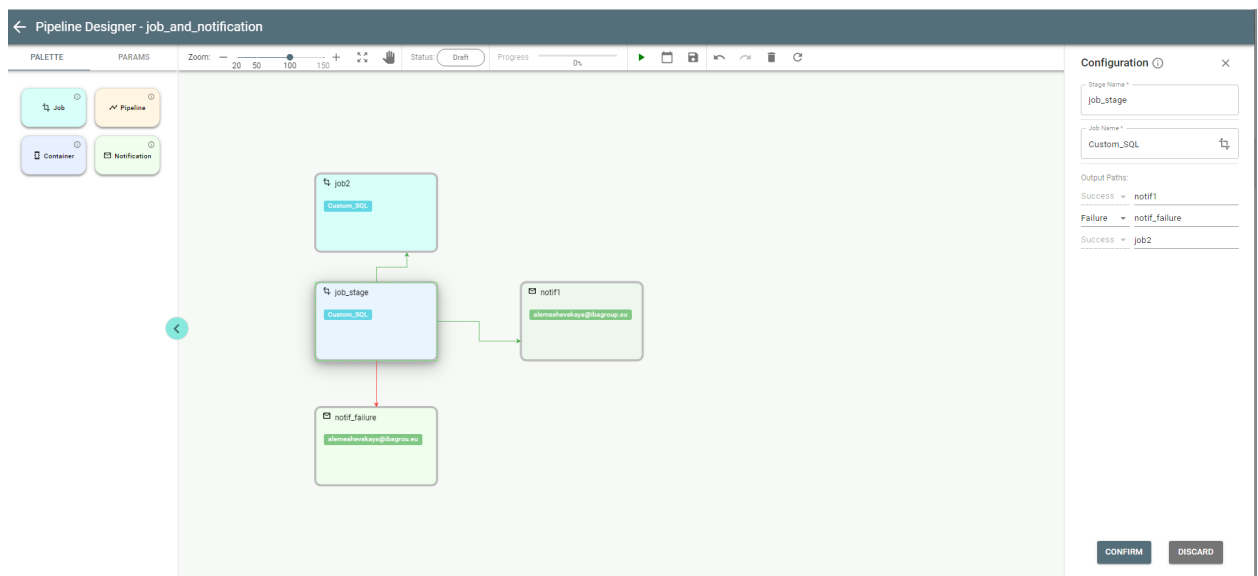
CONFIRM

DISCARD

7) Save the stage by pushing *Confirm* button on the panel. Push *Save* button on header If you want to save the pipeline at this step.

8) Drag and configure other stages. Connect them with the same manner as in Job Designer. Note: you can also add a stage by double clicking its tile on the palette.

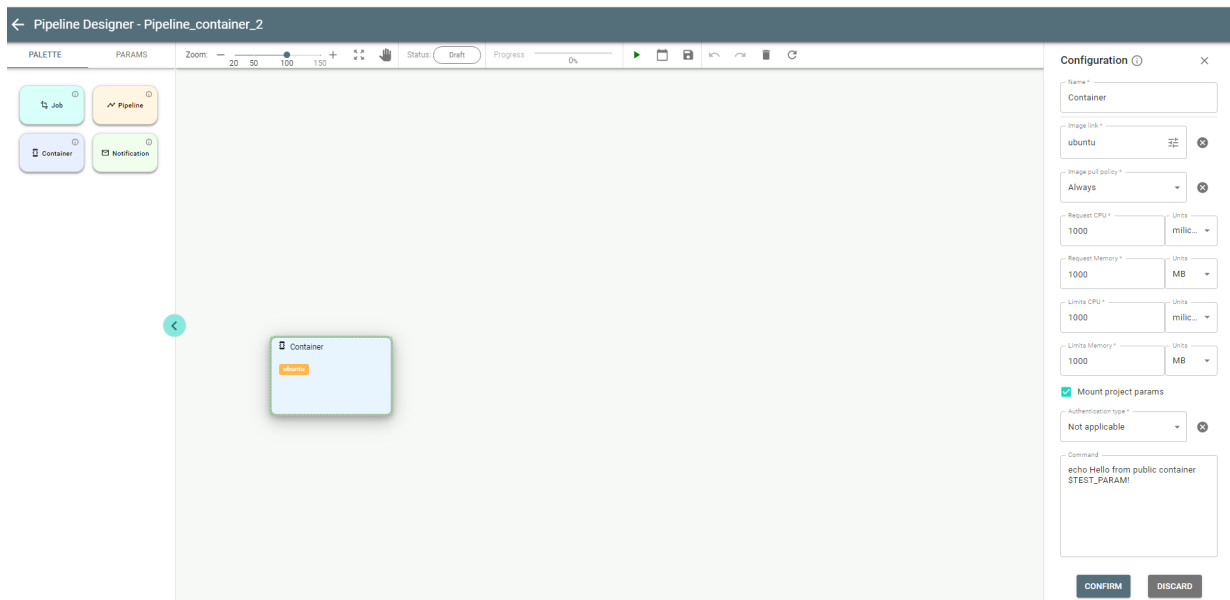
You can link stages based on success or failure of each stage. After connecting stages to each other you can choose Success or Failure link on the configuration panel. There can be only one connection for failure. See the example of configured pipeline:



A *Custom container* stage is required to run custom commands to execute any logic in a pipeline. You can use docker image instead of custom commands.

Start creating a pipeline by dragging *Container* stage to the canvas and enter parameters in the

Configuration panel:



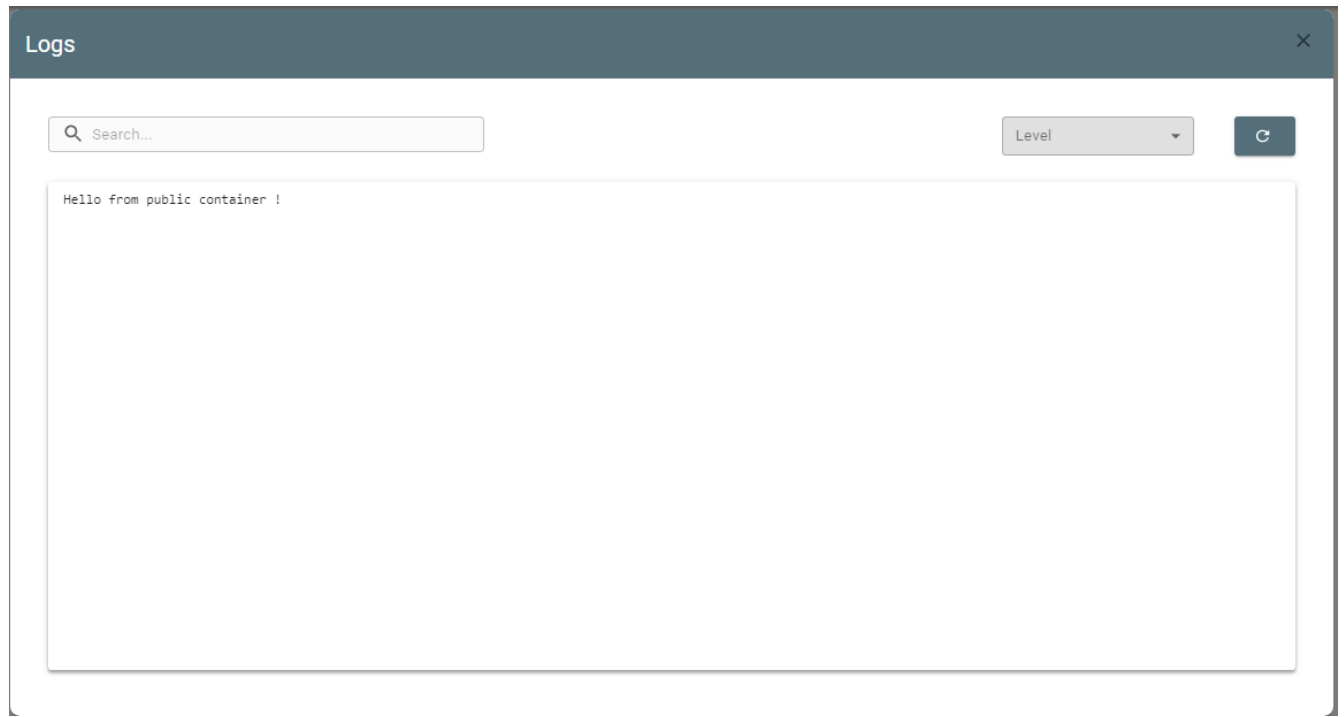
Container stage has the following fields in the Configuration:

- ✓ *Image link*. Docker image path.
Examples:
mysql,
mysql:latest,
bitnami/argo-cd:2.1.2,
localhost:5000/bitnami/argo-cd:2.1.2,
registry.redhat.io/rhel7:latest.
- ✓ *Image pull policy*. Defines when the image is pulled (downloaded).
Possible values:
 - *If not present* – is downloaded only if it does not exist locally;
 - *Always* – is downloaded before each start;
 - *Never* – is not downloaded, local copy is used.
- ✓ *Requests and Limits CPU*
- ✓ *Requests and Limits memory*
- ✓ *Mount project parameters*. Defines whether to mount all project parameters as environment variables inside the Pod.
- ✓ *Authentication type*
- ✓ *Authentication mode* can be one of these:
 - *Not applicable*: image pull secrets are not required as the image is pulled from the public registry;
 - *New*: create a new image pull secret on the fly by providing all necessary information;
 - *Provided*: use existing image pull secret by providing its name (Image pull secret name).
- ✓ *Image pull secret name*. Name of the secret to pull the image. Note that it must exist within the same k8s namespace as the current pipeline.
- ✓ *Username*
- ✓ *Password*
- ✓ *Registry*. Name of the registry for authentication.
- ✓ *Command*. The command to be executed once Pod is created.

Important:

Container stage has a *Logs* button  leading to Logs window.

If the pipeline completed successfully the logs display the message contained in *Command* field in the configuration of the *Container* stage.




Before the first run or after updating the status of the pipeline is *Draft*.
See each stage border painted in *Grey* color, which stands for *Draft*.


5.3. Pipeline Designer Functions Overview


The following functions are available in *Pipeline Designer*:


- ✓ Zoom functions:






- ✓ Move elements: 


- ✓ Move elements/screen: 

- ✓ Show pipeline status: 


- ✓ Show pipeline progress: 

- ✓ Run pipeline  / Stop pipeline  (for running)

- ✓ Save pipeline 

- ✓ Undo / Redo operation on canvas 

- ✓ Remove element from canvas 

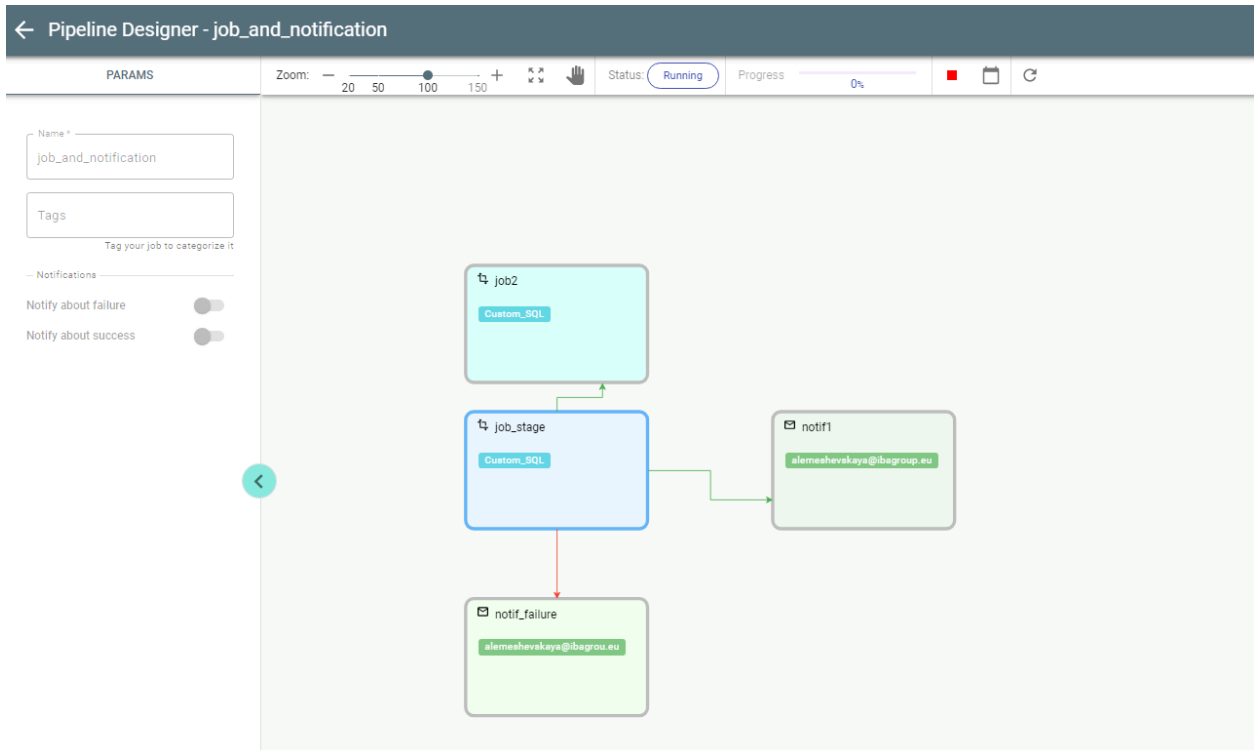
- ✓ Create cron schedule for pipeline 

- ✓ Refresh 

5.4. Pipeline Execution

If you run a pipeline, e.g. from the above example, its status changes from *Draft* to *Pending* and then to *Running*. Push Refresh to update the status.

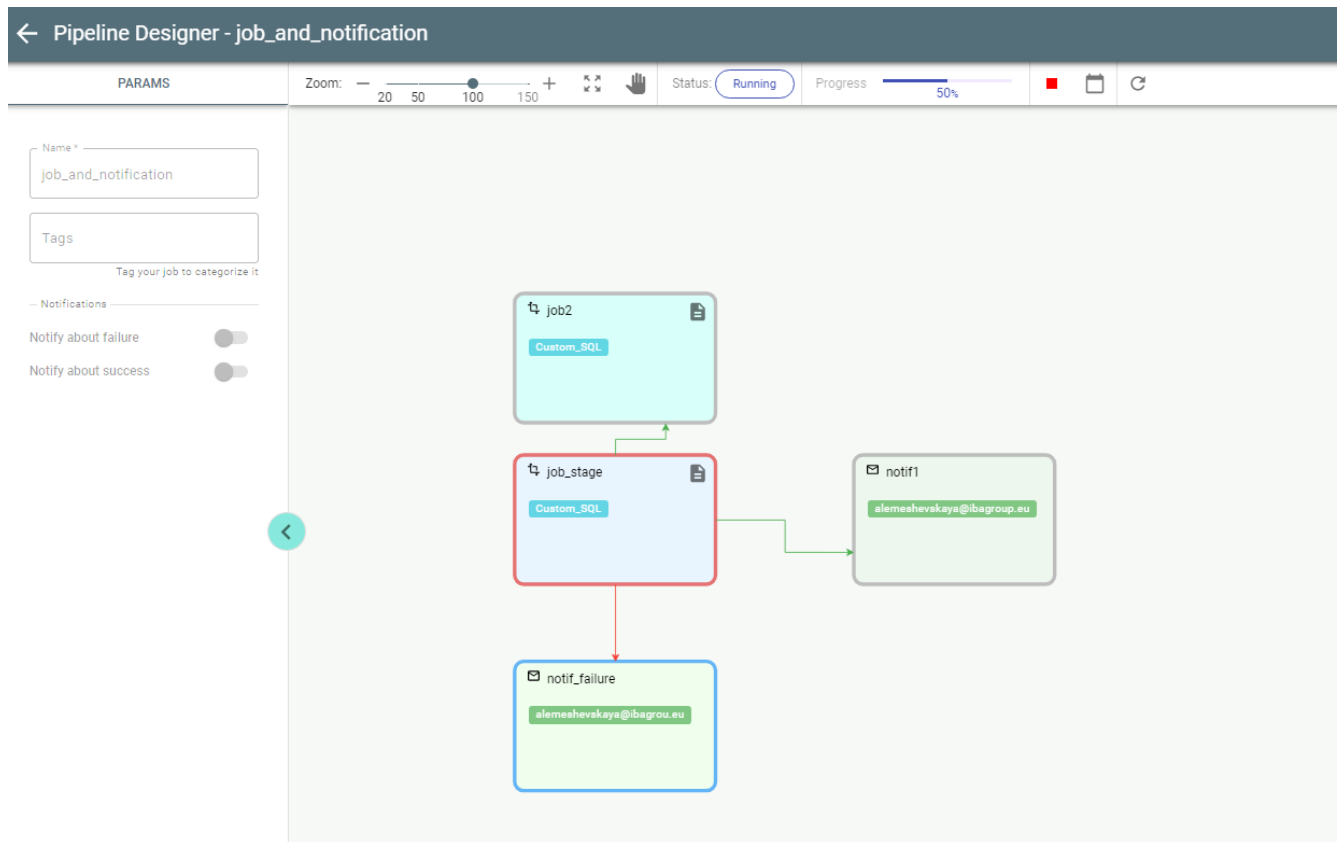
The border of the stage currently running is painted in *Blue*:




If a pipeline succeeded, all completed stages are painted in *Green* indicating success.

The stages configured for failure scenario (red arrow) remain *Grey* as *Draft* as they have not been executed.

If a pipeline fails, then *Red* border indicates the failed stage:



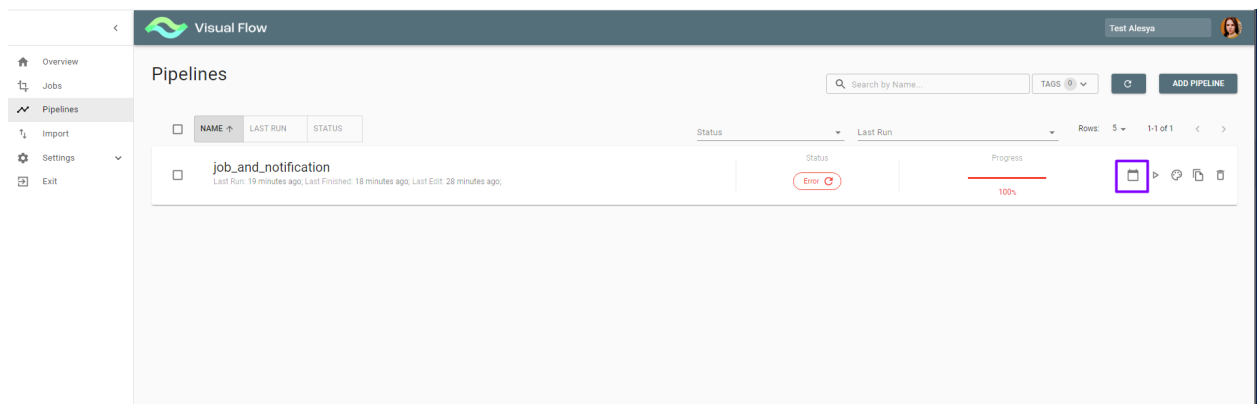
Failed pipeline can be re-run from the point of failure with button  located on the Pipelines Overview Screen.

Important:

Job stage has a *Logs* button  for analyzing logs of a certain job.

5.5. Scheduling a pipeline.

You can schedule a pipeline to run with cron scheduling button:



It opens cron scheduling window:

Scheduling

☐ On/Off

MINUTE (0-59) HOUR (0-23) DAY OF MONTH (1-31) MONTH (1-12) DAY OF WEEK (0-6)

*	any value
,	separator or value list
-	range of values
/	step values

SAVE

CANCEL

Switch cron toggle on and enter values of minute/day/day of month/month/day of week through space according to the tips:

Scheduling

☒ On/Off

"At 01:01"

Next at: 2022-09-23 21:01:00 UTC

MINUTE (0-59) HOUR (0-23) DAY OF MONTH (1-31) MONTH (1-12) DAY OF WEEK (0-6)

*	any value
,	separator or value list
-	range of values
/	step values

SAVE

CANCEL