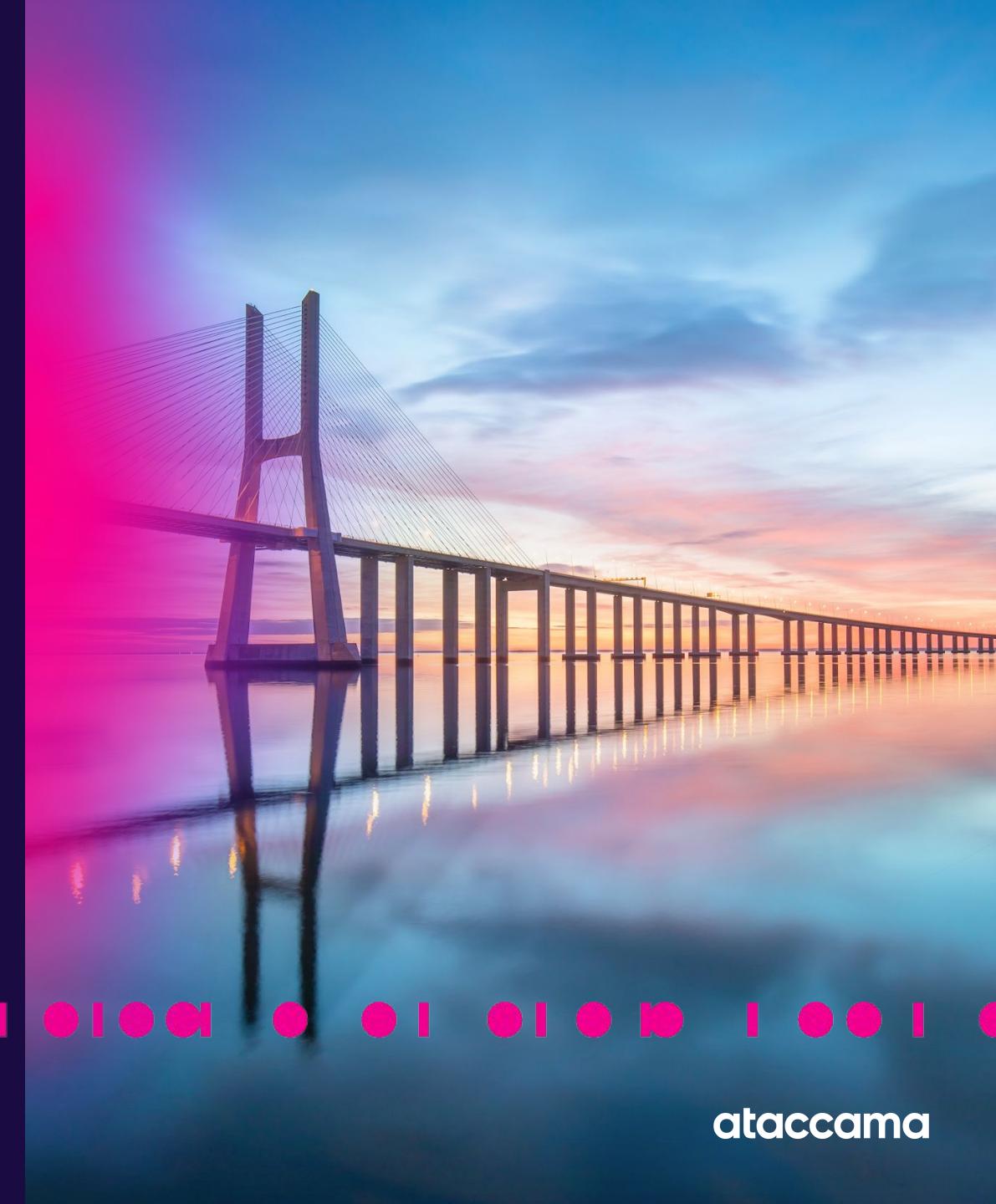


ONE Desktop Core

v15.4.x



ataccama



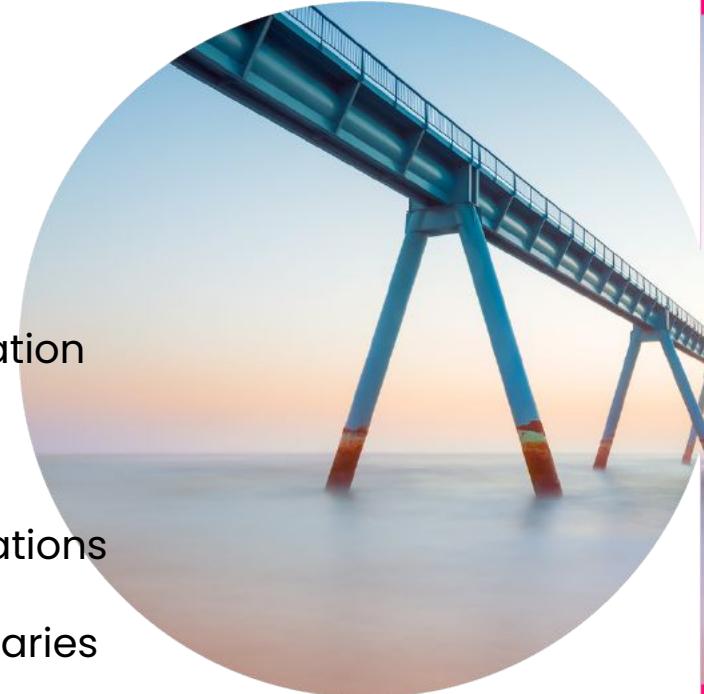
Agenda

Day 1

- Introduction to ONE Desktop
- Quick Demo
- Read / Write Operations
- Data Profiling
- Flow Control

Day 2

- Best Practices & Conventions
- Expressions
- Data Transformation
- Step Debugging
- Scores & Explanations
- Lookup & Dictionaries



There are 8 practical lab exercises on this course.

Introduction to ONE Desktop

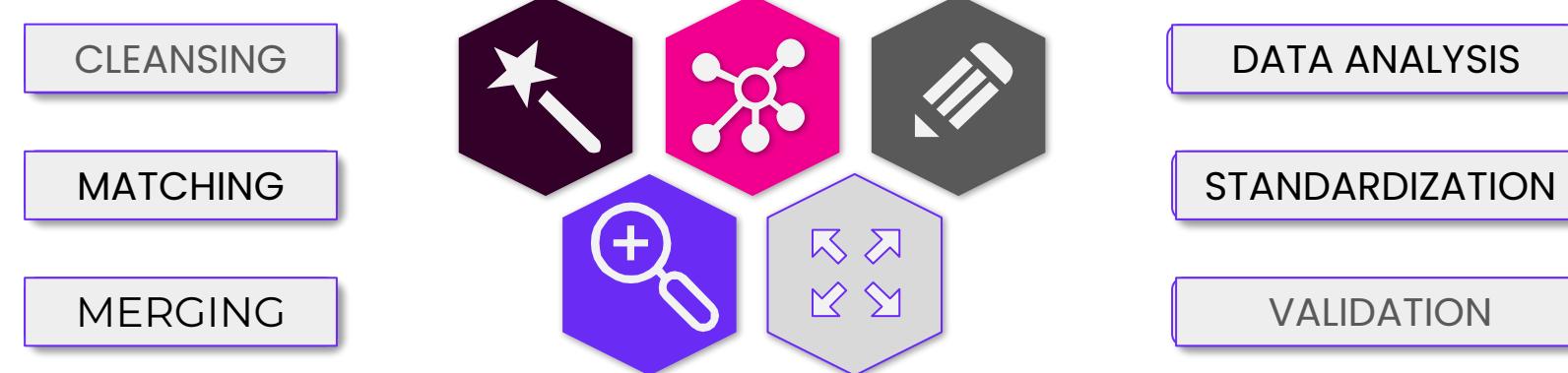


ataccama

ONE Desktop – Introduction

ONE Desktop is a powerful engine for improving data quality coming from various data sources.

Key features:



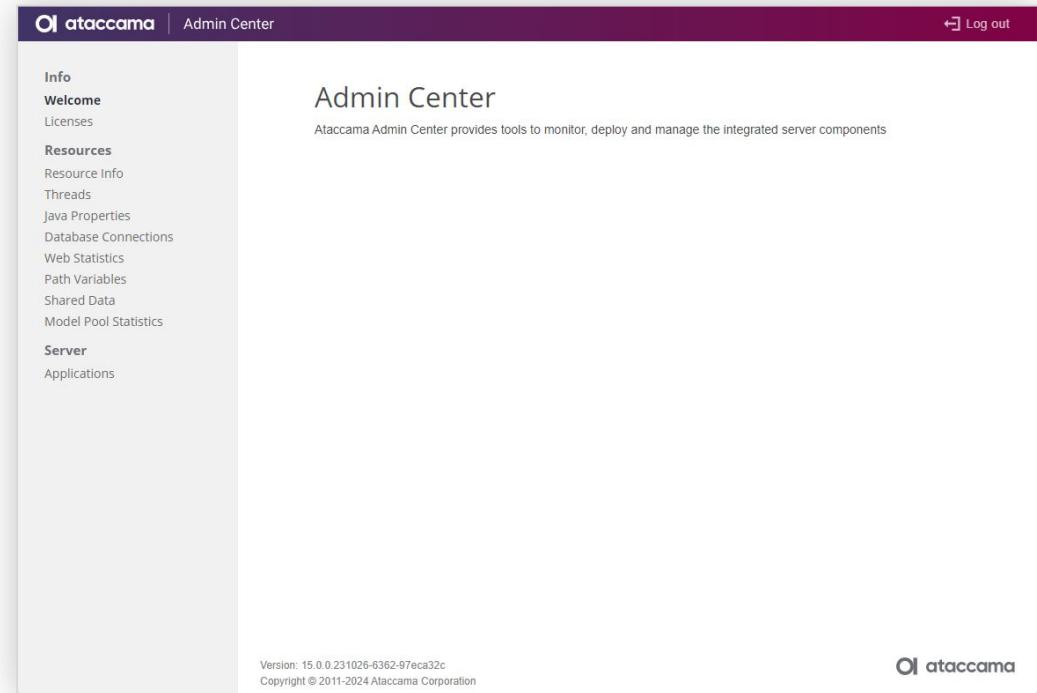
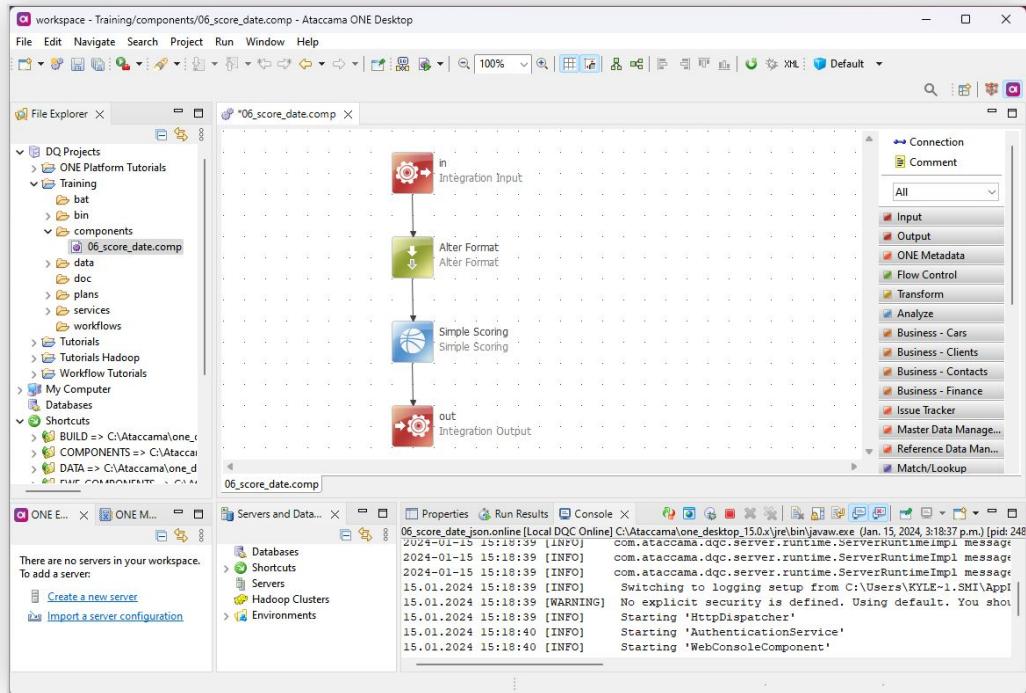
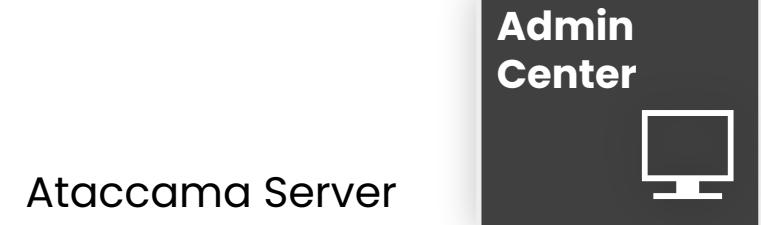
Highlights:

- Available ready-to-use DQ modules
- Flexible and customizable component architecture
- Works in both online and batch modes

Typical use case scenarios:

- Assessment: Data Profiling
- Data Cleansing
- Data Matching and Enrichment
- Prevention: Online Services & DQ Firewall
- Address Cleansing and Identification
- Rule-based engine for complex Data Quality Monitoring and Reporting

ONE Desktop – Architecture



ataccama

ONE Desktop – Architecture

ONE Desktop components and requirements:

Server engine (core)

+ System requirements

- › Server engine (runtime core):
 - Executes plans and services
 - Contains various tools and documentation
- › Platform-independent
 - Can be run on almost any platform
(combination of OS and architecture)
- › Requires a suitable Java Runtime Environment (JRE) on the platform
 - Server engine requires JRE 1.8 or higher

GUI / ONE Desktop

- › Eclipse-based application
 - Bundled with an updated JRE
 - No additional packages need to be preinstalled
- › ONE Desktop GUI
 - Interprets the XML code from plans and displays it in an easy-to-use form
 - Calls the runtime whenever a plan or service is executed
 - Shows console and monitoring, as well as additional information
 - Allows the user to configure sets of runtime variables for the execution of each specific plan
 - Includes debugging, content assists, and searchable help for ease of development

ONE Desktop: Other Requirements

License Keys

License keys have **PLF** extension.

- Created and managed by Ataccama.
- Version-dependent.
- Can be restricted to only allow the use of a certain subset of modules and components.
- Can be generic or detail-specific (e.g. user, architecture, validity).

How to install a license?

- Place the **PLF** license file in the **USER** folder or in **{ONE/Desktop}/runtime/license_keys**.
- Install the license directly from the GUI/IDE using the menu **Help > License Information > Install License**.

Database Connections

For connectivity with a DB, a JDBC driver must exist.

- Most DBs – and all the popular ones – have JDBC drivers available.

ONE Desktop bundles come with popular JDBC drivers that can be found in **{ONE/Desktop}/runtime/lib/jdbc**.

How to Install a driver?

- Server: Place the JDBC driver as a JAR file in the **{ONE/Desktop}/runtime/lib** folder.
- GUI (Desktop): Navigate to the menu **Window > Preferences** and choose **Ataccama ONE/Desktop > DB drivers**.

Creating a new project

Before you begin doing any transformations or analysis of data, you need to create a project. A project is a set of multiple plans, components, data files, workflows, and custom scripts, which are usually focused on solving a specific task and are logically organized into folders.

Empty project

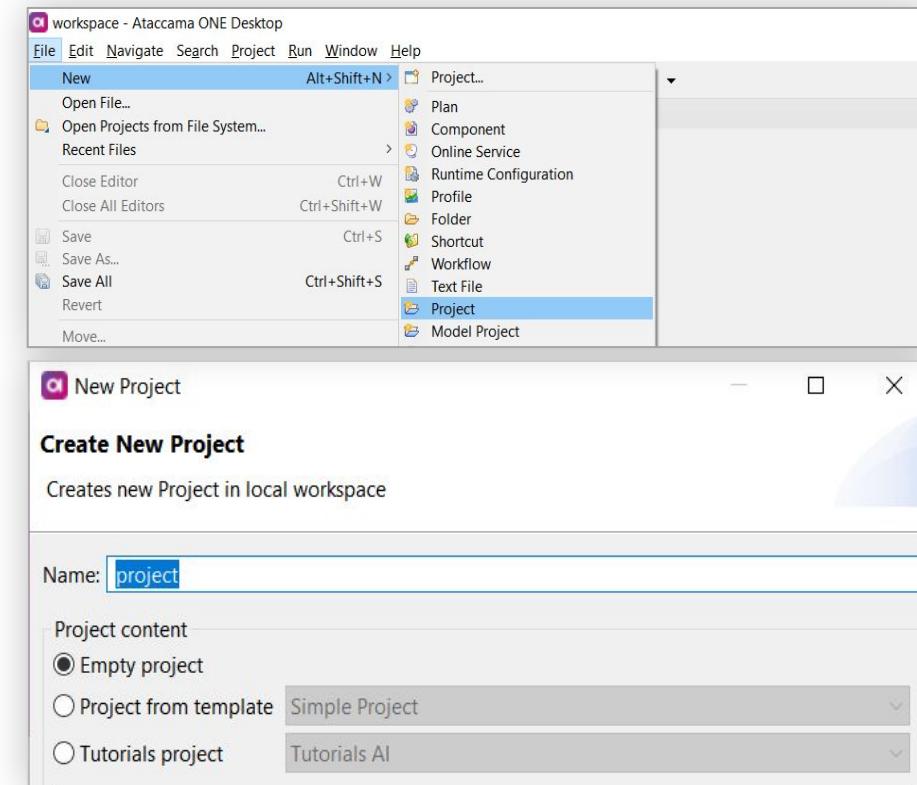
- No predefined folder structure

Project from template

- **Simple Project** – begin with two default directories
- **Complex Project** – begin with a full folder structure

Tutorials Project

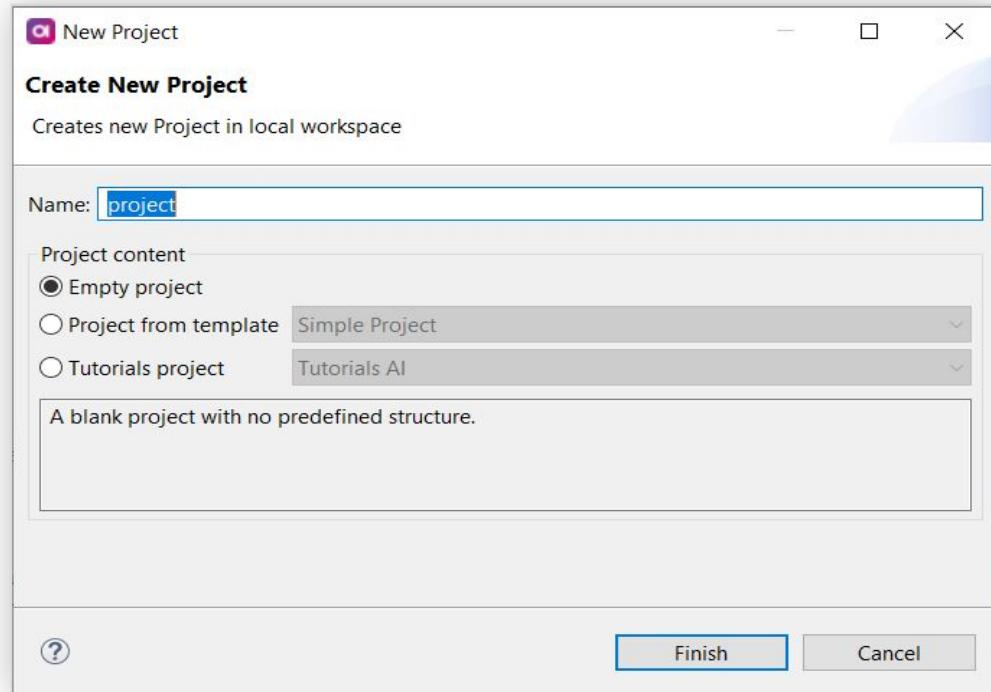
- Tutorials
- Tutorials Hadoop
- Workflow Tutorials
- ONE DQM Server Example



Project Types: Projects

Projects

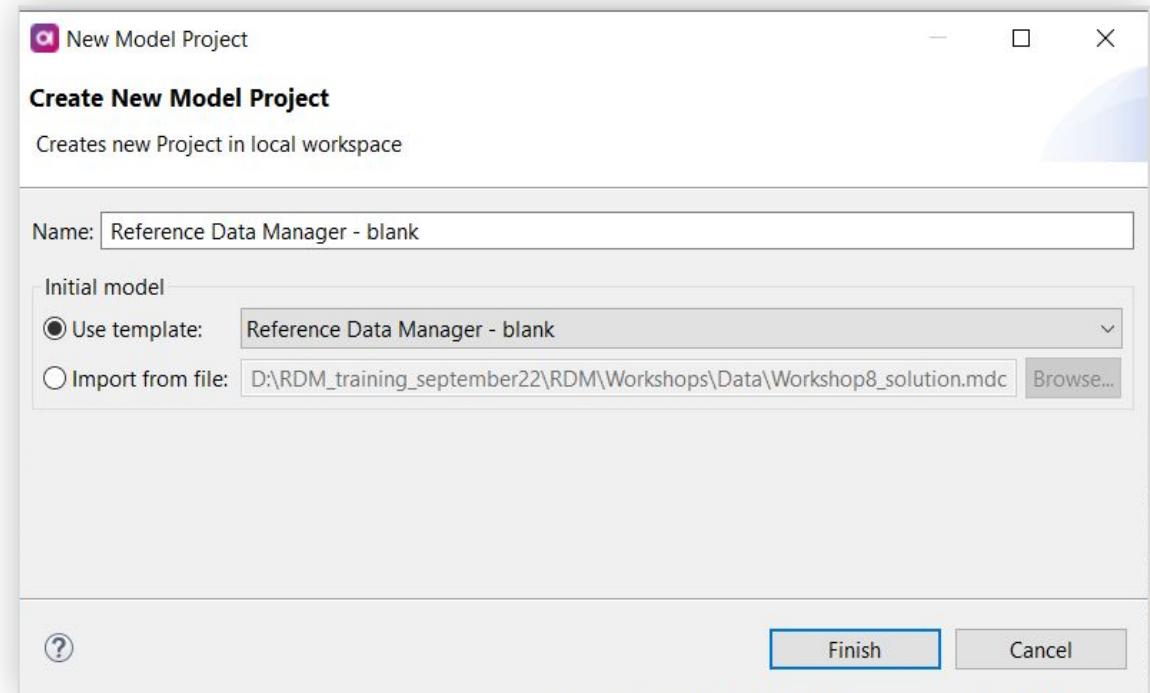
- Empty Project
 - No predefined folder structure
- Project from template
 - Simple Project – two default directories
 - Complex Project – full folder structure
- Tutorials Project
 - Tutorials
 - Tutorials Hadoop
 - Workflow Tutorials
 - ONE Platform Tutorial



Project Types: Model Projects

Model Projects

- File for import
 - Reuse or modify existing projects
- Templates Configured to a specific Ataccama module:
 - MDM
 - Profiling
 - DQIT
 - RDM



ONE Desktop: Interface

 **Model Explorer**

 **File Explorer**

- DQ Projects
- Databases
- Shortcuts
- Servers

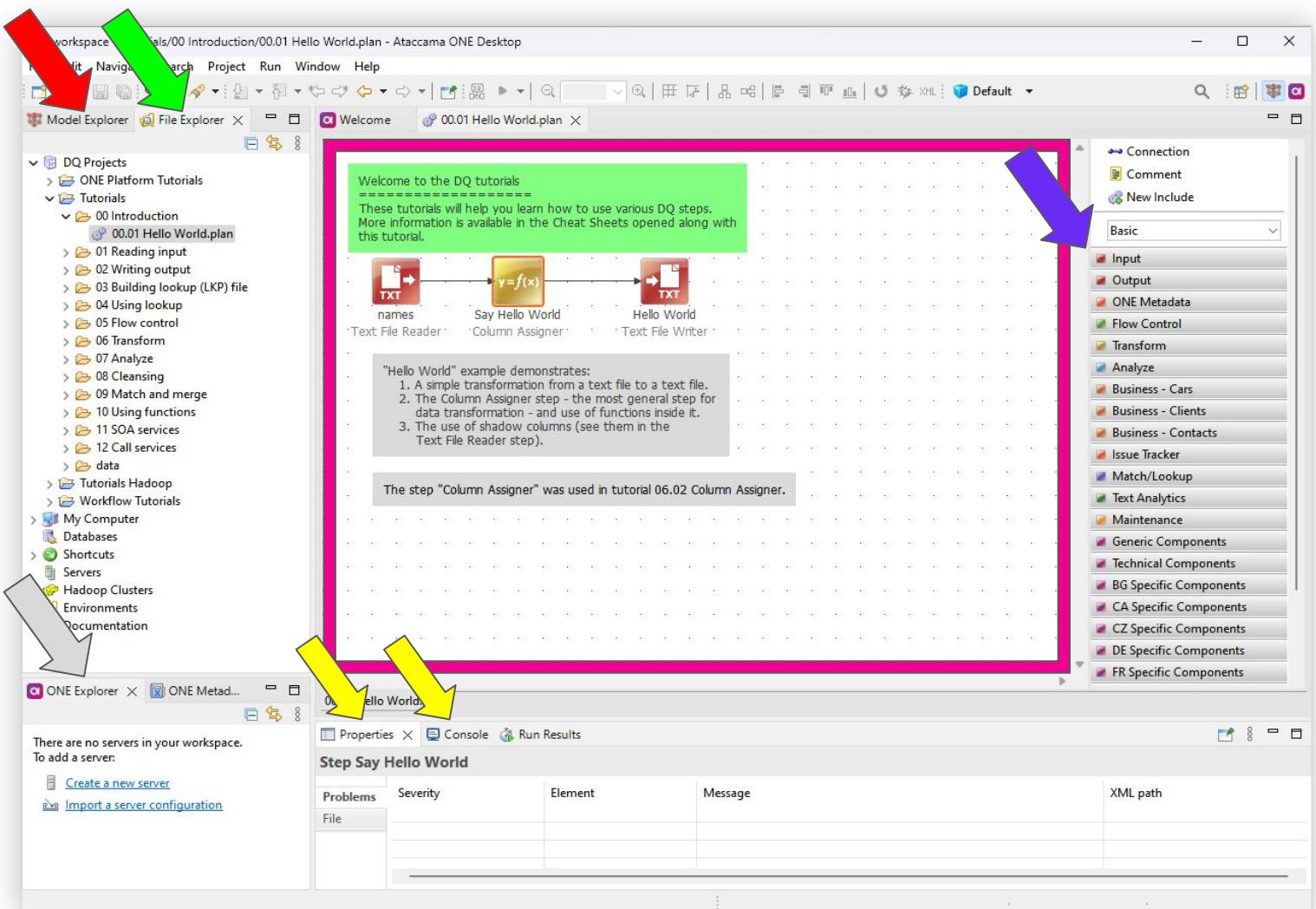
 **Canvas**

 **Palette**

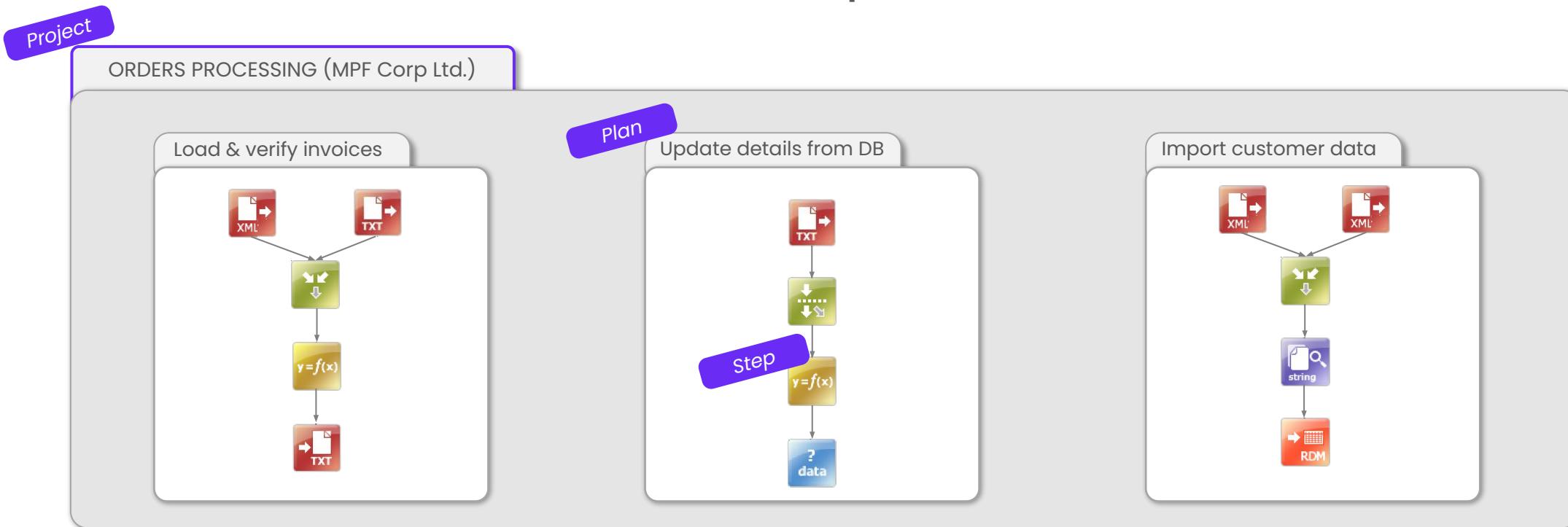
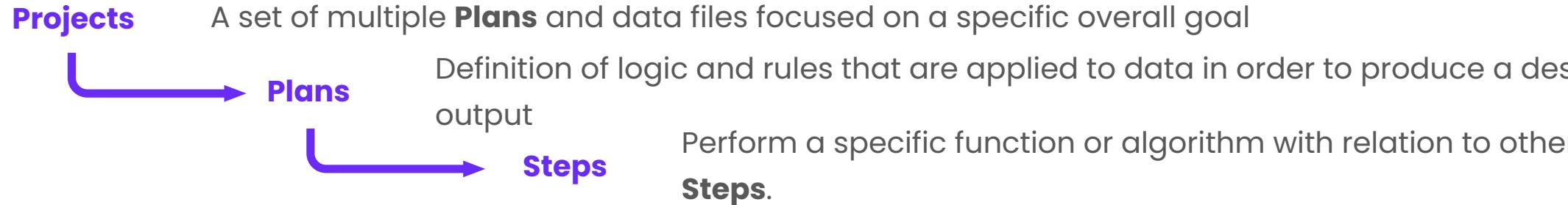
 **Console**

 **Properties**

 **ONE Integration**



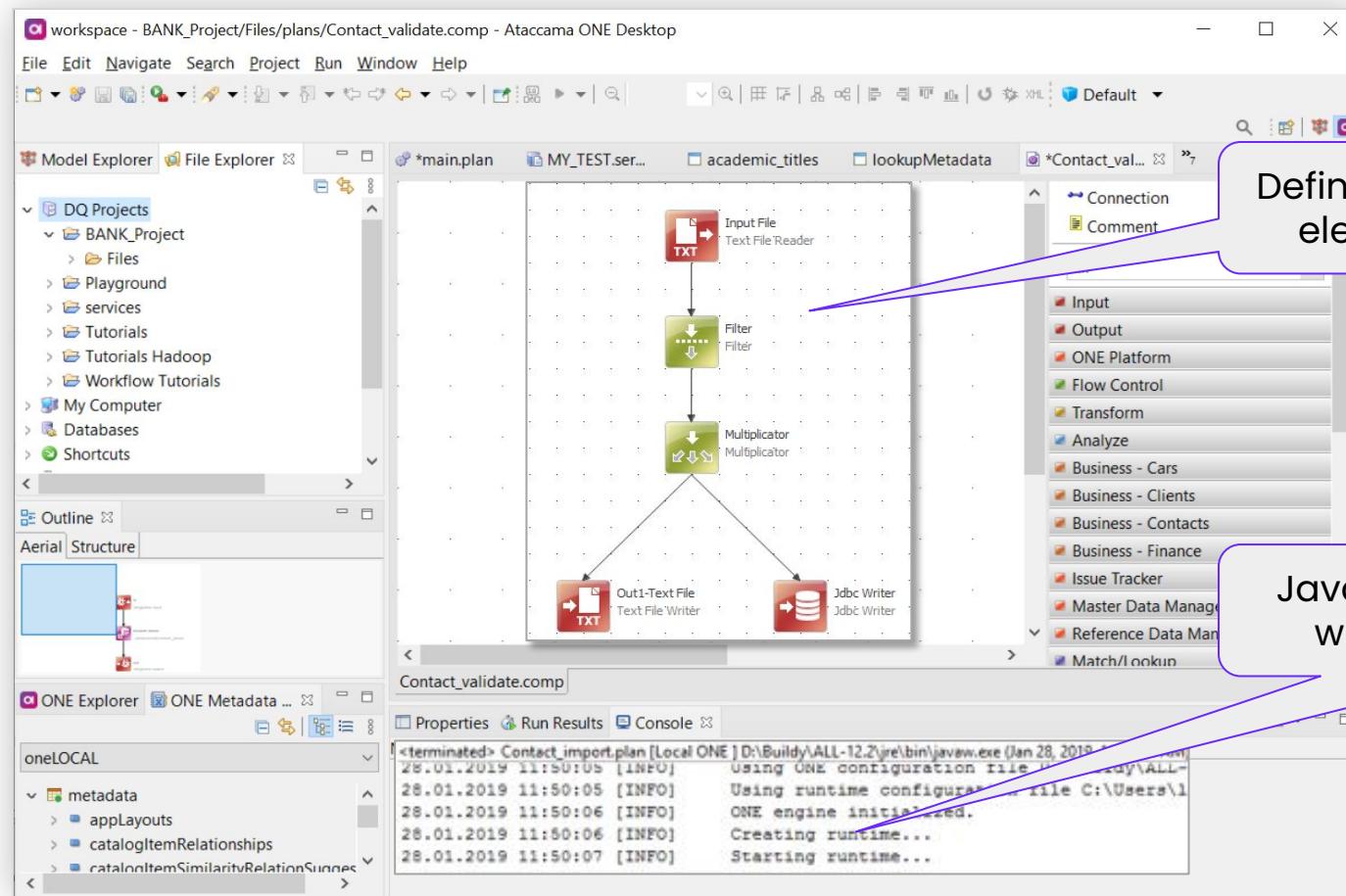
ONE Desktop: How it works



ONE Desktop: Plans

Plans can be described as a canvas where all Steps are placed and linked together.

main.plan



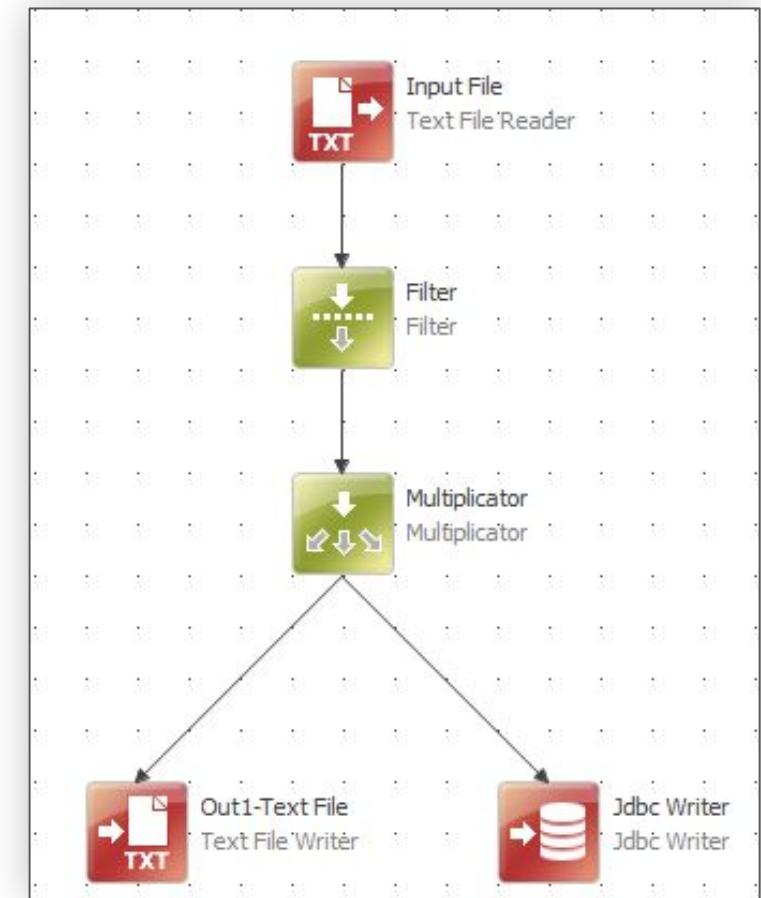
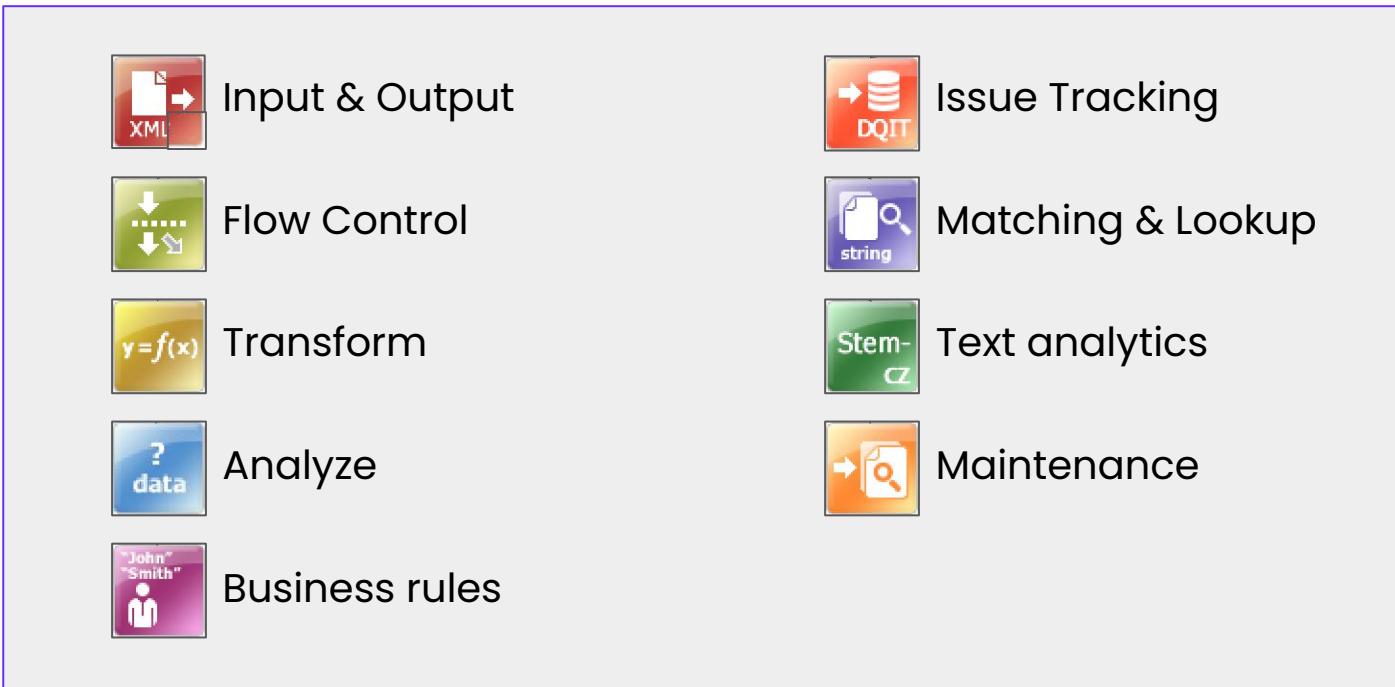
Definition of the plan's logical elements and their order

Java runtime execution log with processing details

ataccama

ONE Desktop: Step Types

Steps are divided into the following categories:



ONE Desktop: Commonly Used Steps

The most useful steps to begin with:



Input / Output

- Text File Reader/Writer
- JDBC Reader/Writer
- Excel File Reader/Writer
- Integration Input/Output
- XML Reader/Writer
- Random Record Generator



Flow Control

- Alter Format
- Multiplicator
- Filter
- Extract Filter
- Condition
- Union & Union Same
- Join



Transformations

- Column Assigner
- SQL Select / Execute
- Matching Values
- Apply Replacements
- Regex Matching
- Splitter
- JSON & XML Parser
- Pattern Parser



Analysis

- Profiling
- Simple Scoring
- Data Quality Indicator
- Group Aggregator



Schema Discovery

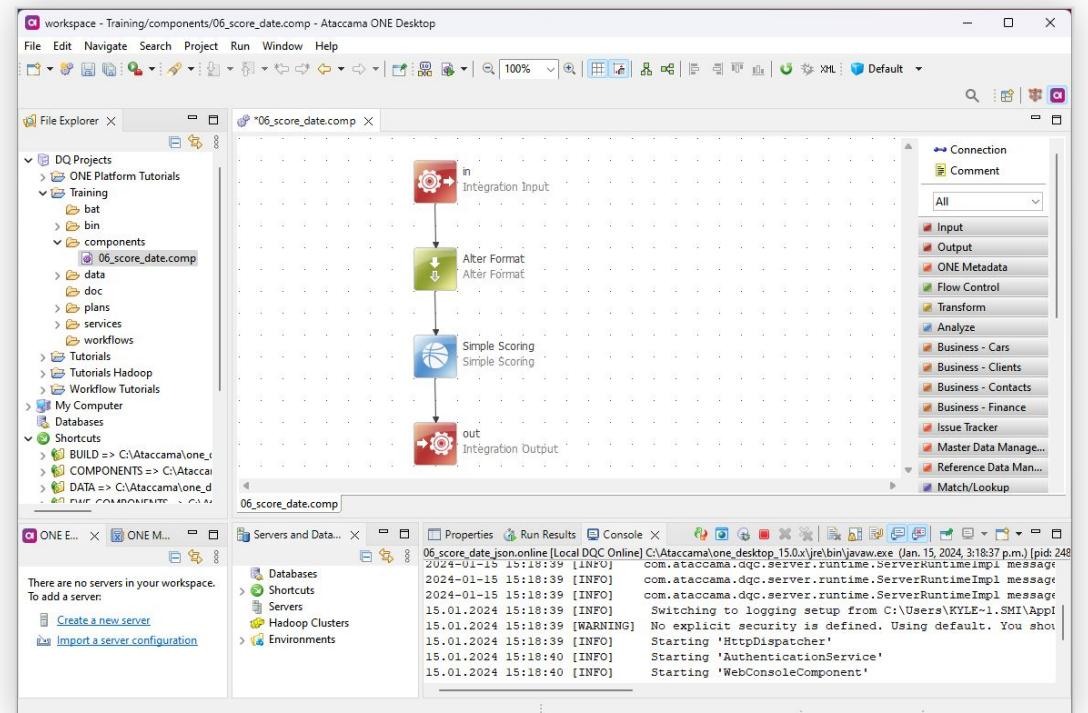


Matching/Lookups

- Lookup & Lookup Builder
- Representative Creator
- Matching
- SOAP / JSON call
- WEB Lookup
- Canopy Clustering

Topic Highlights

- **ONE Desktop** is an **Eclipse-based IDE** designed for creating DQG solutions based on Ataccama products.
- **Projects** are sets of plans, components, workflows, and data, which are logically organized into folders.
- **Plans** define the logic and rules that are applied to data in order to produce the desired output.
- **Steps** perform specific functions or algorithms to data with relation to other **Steps** within plans, components, or workflows.



Memory Refresher #1

Introduction to ONE Desktop



ataccama

Which of the following statements about ONE Desktop is True?

- 1.** Simple, Complex, and Empty projects are the main project types.
- 2.** There are two types of projects: project and model project.
- 3.** ONE Desktop can not be used without any web-applications.

Match the terms with their meanings:

A. Plans

B. Steps

C. Components

C. Components

1. Algorithms that can be reused in different plans.

B. Steps

2. Algorithms for reading, transforming, and analyzing data.

A. Plans

3. Defines the logic by placing elements and connecting them together.

ONE Desktop Demonstration



ataccama

Lab Exercise #1

Environment Preparations



Read / Write Operations



ataccama

Data Sources: Traditional



Steps to read data from several sources:

- **Text File Reader**
- **JDBC Reader**
- **Excel File Reader**
- **XML File Reader**
- **JSON File Reader**
- ...and more



csv and txt (incl. compressed files)

common DB systems, anything with a JDBC Driver

.xls and .xlsx files

flat or hierarchical XML files

flat or hierarchical JSON files



Steps to interact with the ONE Web Application:

- **ONE Metadata**

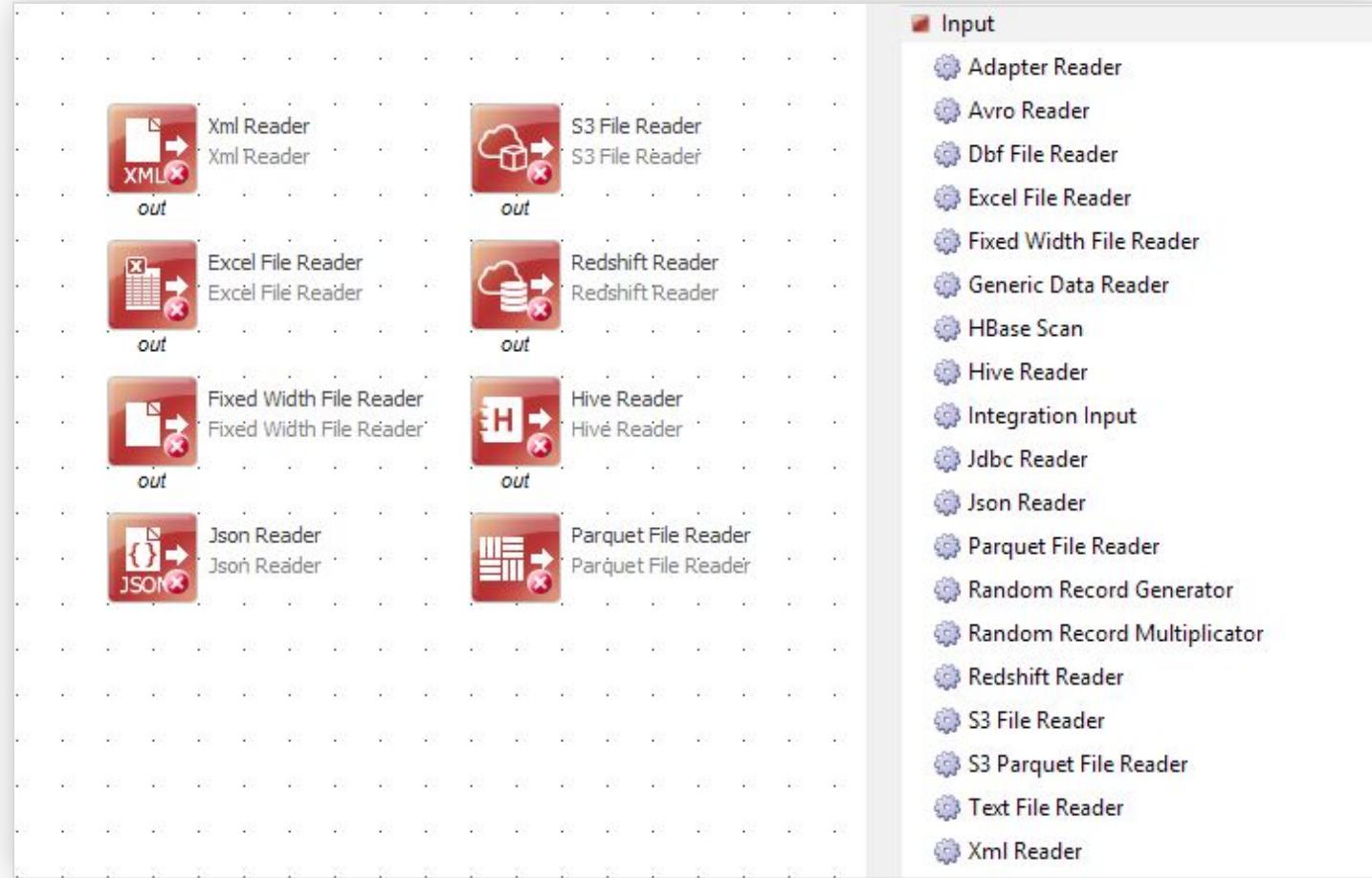


Steps to use in Components and Web Services:

- **Integration Input** and **Integration Output** can be used instead of other, more specialized, input steps to create a plan file that is usable as either a component or an online service.

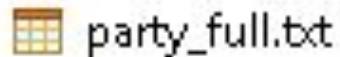


Other Reader Steps



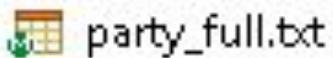
Flat Files: Metadata

In ONE Desktop, a flat file (csv, txt) has no metadata by default.



party_full.txt

Once the metadata is created, the file gets a green „M“ (metadata) icon.



party_full.txt
party_full.txt.metaData

Metadata is stored by default as a hidden XML file.

To create or edit metadata, **double-click the file or use right-click > Edit Metadata** in the File Explorer.

Flat Files: Metadata Editor

Format Advanced

File type: Delimited Encoding: UTF 8 Preferences...

Line separator: CRLF String qualifier: none
Field separator: ; String qualifier escape: none
Data starts at line: 2 Escape character: none

	Use	Name	Type	Format
4	<input checked="" type="checkbox"/>	src_sin	STRING	
5	<input checked="" type="checkbox"/>	src_card	STRING	
6	<input checked="" type="checkbox"/>	src_address	STRING	
7	<input checked="" type="checkbox"/>	src_primary_key	INTEGER	thousands separator=","
8	<input checked="" type="checkbox"/>	meta_last_update	DAY	locale='en_US', format='yyyy/MM/dd'
*				

Treat All as Strings
Reset Columns

Automatic refresh Refresh

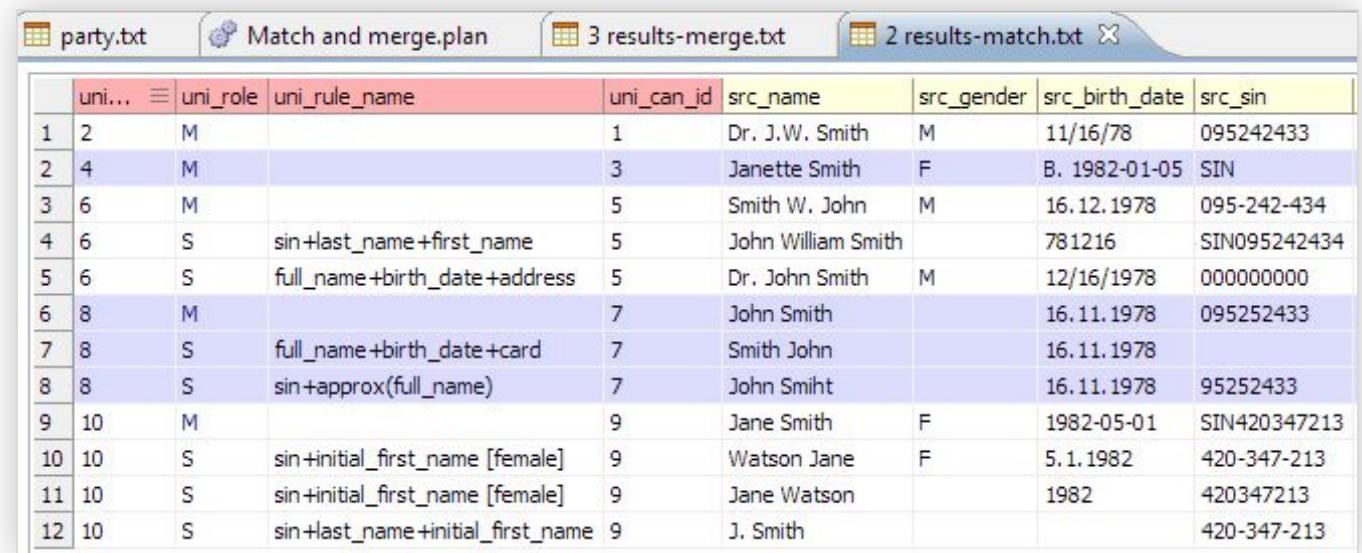
Result data Decoded original data

	src_name	src_gender	src_birth_date	src_sin	src_card
1	Dr. John Smith	M	12/16/1978	000000000	88682239496
2	Smith W. John	M	16.12.1978	095-242-434	266805R079844QR

Flat Files: CSV Viewer

Flat files are opened in the CSV Viewer which allows sorting, filtering, coloring.

See "Getting started with Ataccama DQC" in **Help > Help Contents** for details.



The screenshot shows the Ataccama CSV Viewer interface. At the top, there are four tabs: "party.txt", "Match and merge.plan", "3 results-merge.txt", and "2 results-match.txt". The "2 results-match.txt" tab is currently active. Below the tabs is a table with 12 rows and 9 columns. The columns are labeled: "uni...", "uni_role", "uni_rule_name", "uni_can_id", "src_name", "src_gender", "src_birth_date", and "src_sin". The data in the table is as follows:

	uni...	uni_role	uni_rule_name	uni_can_id	src_name	src_gender	src_birth_date	src_sin
1	2	M		1	Dr. J.W. Smith	M	11/16/78	095242433
2	4	M		3	Janette Smith	F	B. 1982-01-05	SIN
3	6	M		5	Smith W. John	M	16.12.1978	095-242-434
4	6	S	sin+last_name+first_name	5	John William Smith		781216	SIN095242434
5	6	S	full_name+birth_date+address	5	Dr. John Smith	M	12/16/1978	000000000
6	8	M		7	John Smith		16.11.1978	095252433
7	8	S	full_name+birth_date+card	7	Smith John		16.11.1978	
8	8	S	sin+approx(full_name)	7	John Smiht		16.11.1978	95252433
9	10	M		9	Jane Smith	F	1982-05-01	SIN420347213
10	10	S	sin+initial_first_name [female]	9	Watson Jane	F	5.1.1982	420-347-213
11	10	S	sin+initial_first_name [female]	9	Jane Watson		1982	420347213
12	10	S	sin+last_name+initial_first_name	9	J. Smith			420-347-213

Excel Files: Metadata

To create metadata for an Excel file:

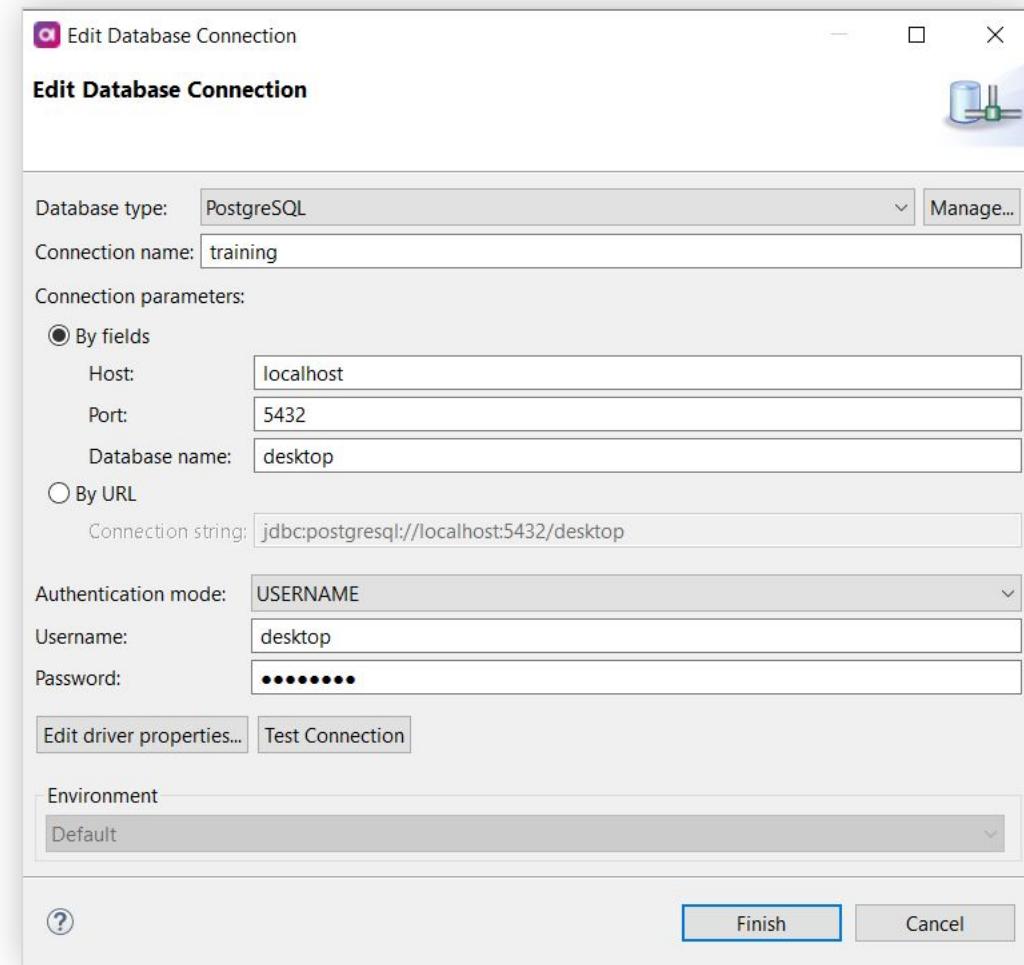
- Right-click the file and choose '**Edit Metadata**'.
- Drag & drop the file onto a plan.

The screenshot shows the Ataccama interface with the 'Edit Metadata' dialog open. The top section displays the selected data range from 'Sheet1 (default)' starting at row 2. The data consists of five columns: Name, Gender, Date, SIN, and Address. The bottom section shows the 'Column names and types' mapping, where each column is identified by its name and assigned a type of 'STRING'.

Name	Type
1 Name	STRING
2 Gender	STRING
3 Date	STRING
4 SIN	STRING

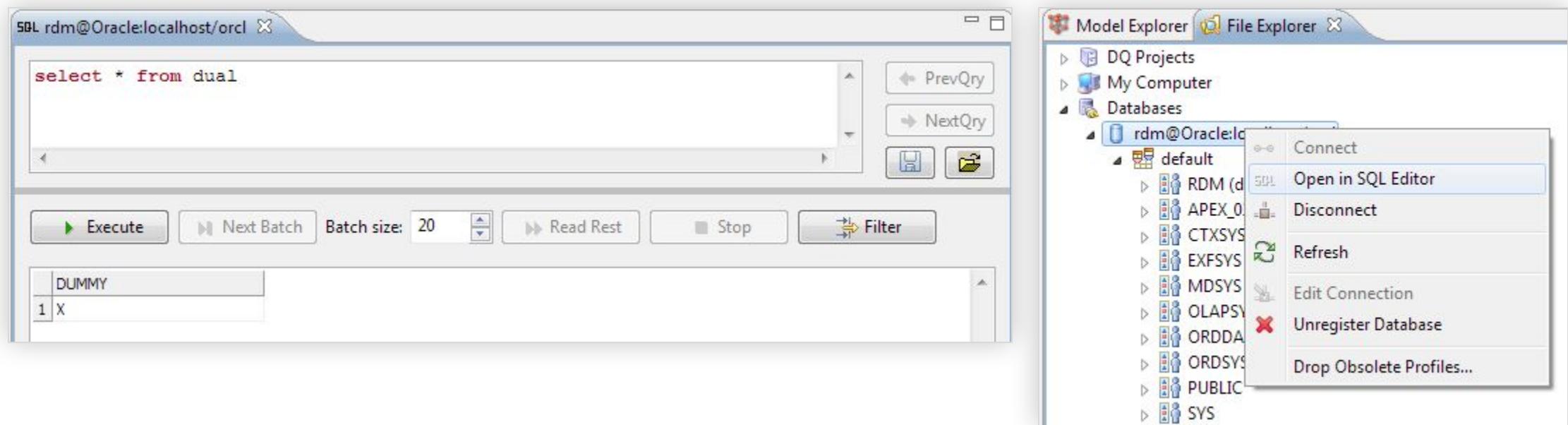
Database Connections: Creation

ONE Desktop supports direct connection to supported DBs via the File Explorer.



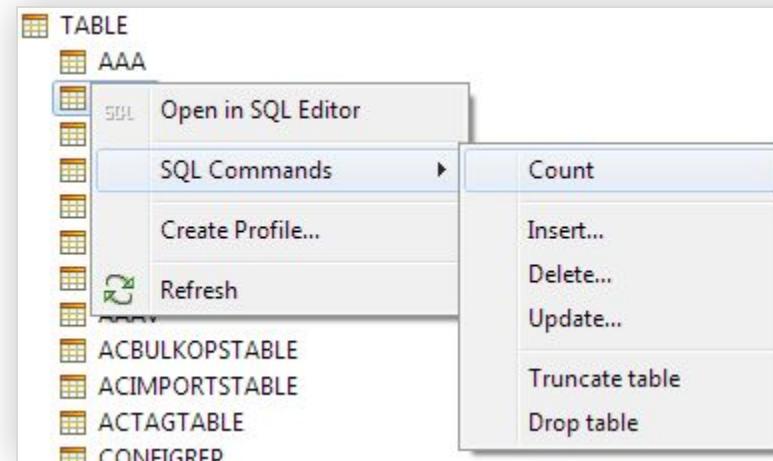
Database Connections: SQL Editor

- Simple SQL queries can be executed directly within the ONE Desktop.
- Double-click on a table to open it in SQL Editor (or right-click > Open in SQL Editor).
- The SQL query field offers content assist (CTRL + Space).



Database Connections: Custom Commands

- Right-click a table and choose **SQL Commands** to execute common queries.
- Custom SQL Commands can be configured in **Preferences > ONE Desktop > Database > Table Commands**.



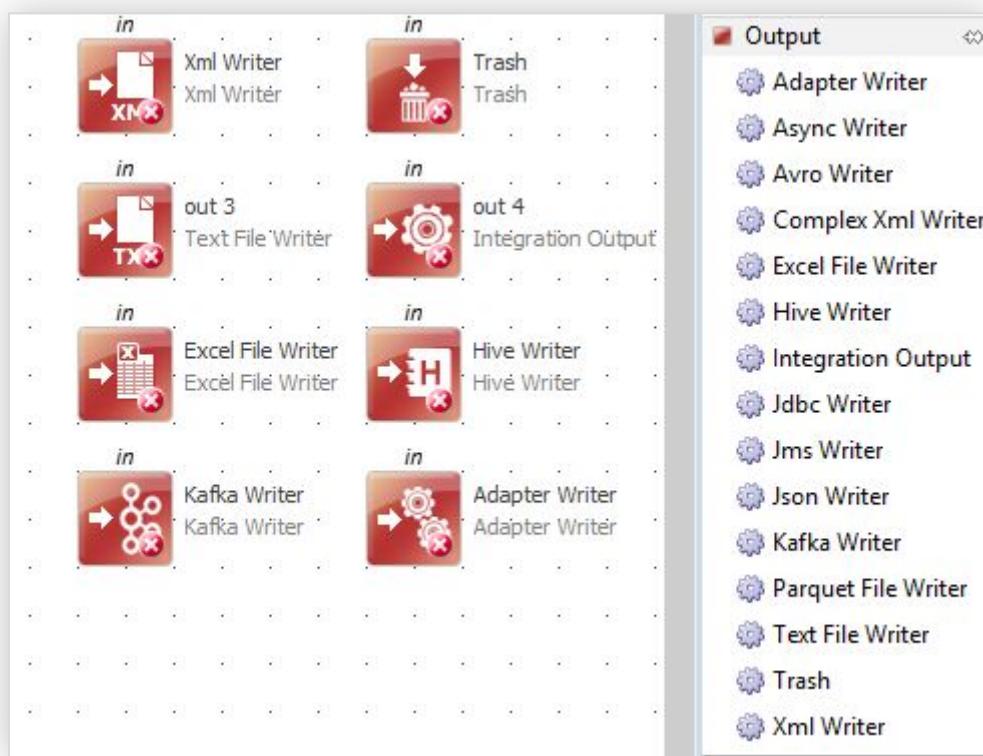
Data Sources: Creating Reader Steps

Drag & drop a file (text or Excel) or a database table into a plan to create the appropriate **Reader Step** (Text File/Excel/Jdbc).

- For **files**, the reader step is configured based on the metadata.
- For **DB tables**, a configuration wizard is shown to help set up the reader step.

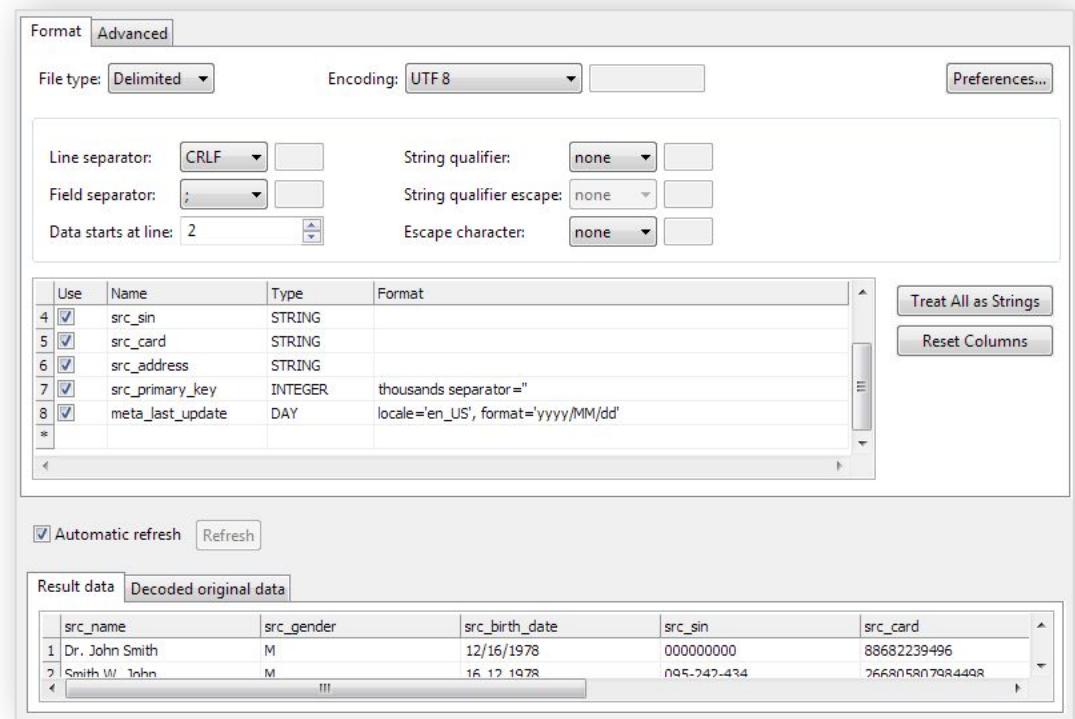
Data Output

ONE Desktop can write data in several ways:



Topic Highlights

- **Metadata** describes the format of the data, primarily for **CSV**, **TXT**, and **XLSX** files, for properly using in plans.
- Text files can be opened in the **CSV Viewer** by double-clicking. Clicking the name of the column in the header row will sort that column.
- A data file can be directly added to a plan by **dragging & dropping** it into the plan's canvas, where the corresponding reader step will be generated.



Memory Refresher #2 Read / Write Operations



Metadata must be created before a data file can be added to a plan.

- TRUE
- FALSE

There are Reader steps that will allow input from and output to the ONE Web Application.

- TRUE
- FALSE

Data Profiling



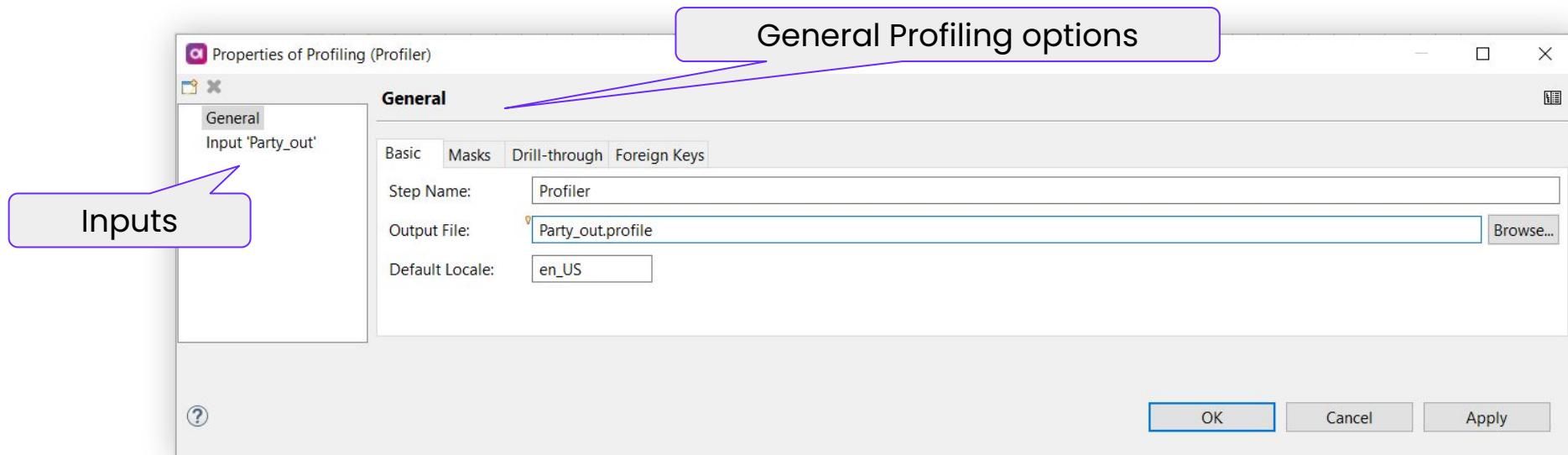
ataccama



Profiling Step

The Profiling step provides its outputs directly to the ONE Desktop Environment. It serves primarily for data analysis.

- Generates a special file '*filename.profile*' which contains the profiling results.
- Besides common statistical analysis, the **Profiling** step also allows highly advanced analysis.
- Columns can be automatically mapped with the '**Fill Columns**' button.
- **Drill-through** enables quick access to sliced data; limits can be set for working with large volumes of data.



Statistical Analysis

The ONE Desktop IDE provides common statistical analysis as a simple, one-click operation.

- Profile files can be created for inputs from a text file or a DB table.
- Right-click the source file/table and choose Create Profile in the contextual menu.

The screenshot displays two windows from the Ataccama ONE Desktop IDE. The left window, titled 'Input: "party"', shows 'Column Analyses' for various columns like src_name, src_gender, and src_card. The right window, titled 'Basic Analyses', shows statistics for the src_name column, including a pie chart of data types and a table of counts. Callouts point to specific features:

- Source-specific analysis**: Points to the 'Column Analyses' section of the left window.
- Field-specific analysis**: Points to the 'src_name' row in the 'Basic Analyses' table.
- Other statistical values**: Points to the pie chart in the 'Basic Analyses' window.
- General statistics**: Points to the 'Counts' table in the 'Basic Analyses' window.

Column Analyses (Left Window)

Expression	Type	Domain	Non-null	Null	Unique	Distinct	Min	Median	Max
src_name	STRING	specval pa...	139.317	0	130.464	132.575	A. Aakre	L. Burglin	Zwiener St...
src_gender	STRING	enum	15.780	123.537	0	4	F	F	MALE
src_birth_d...	DAY		103.535	35.782	2.791	25.007	1900-01-01	1969-12-19	2006-12-31
src_sin	STRING	pattern	134.830	4.487	134.766	134.798	000000000	773843081	SIN945055...
src_card	STRING	long	55.904	83.413	15.800	35.630	266805807...	55103109...	NULL
src_address	STRING	enum patt...	9	139.308	5	7	14618110 ...	8500 Leslei...	Surrey 146...
src_primary...	INTEGER		139.317	0	139.315	139.316	1	69.659	139.317
meta_last...	DAY		139.317	0	12	341	2000-01-01	2008-03-03	2009-12-19

Basic Analyses (Right Window)

Expression	Data type	Domain	Rows
src_name	STRING	specval pattern	139.317

Counts

Type	Count	%
Null	0	0,00%
Non-null	139.317	100,00%
Duplicate	6.742	4,84%
Distinct	132.575	95,16%
Non-unique	2.111	1,52%
Unique	130.464	93,65%

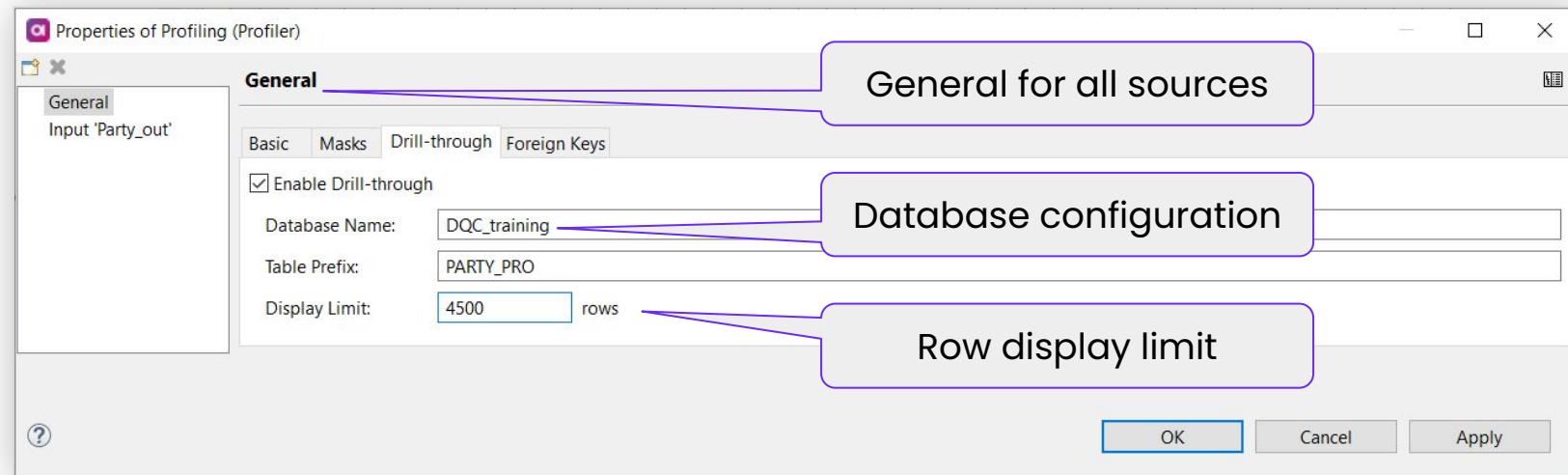
Statistics

Type	Count	%
Unique	130.464	93,65%
Non-unique	2.111	1,52%
Duplicate	6.742	4,84%

Drill-Through Profiling

The drill-through feature can be enabled for any profile.

- Tables will be created in a DB with different prefixes for different profiles.
- Row limit can be set to avoid overloading the DB with huge data volumes.



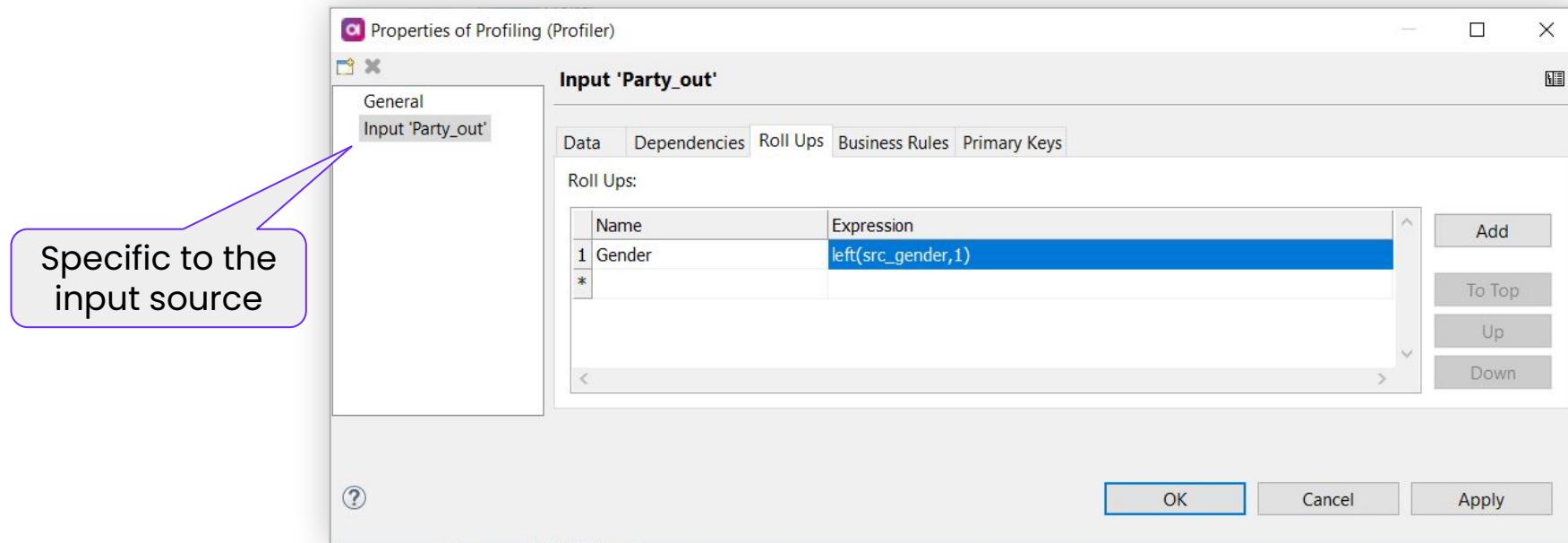
Right-clicking on any value (analysis result) from statistics, business rules, primary key analysis, dependency analysis, etc. will offer you the possibility to drill-through and display the corresponding values.

- EXAMPLE: This feature can be used to find all invalid (FALSE) results of a business rule

Roll-Ups Profiling

Roll-ups allow you to create a kind of sub-profile within a profile.

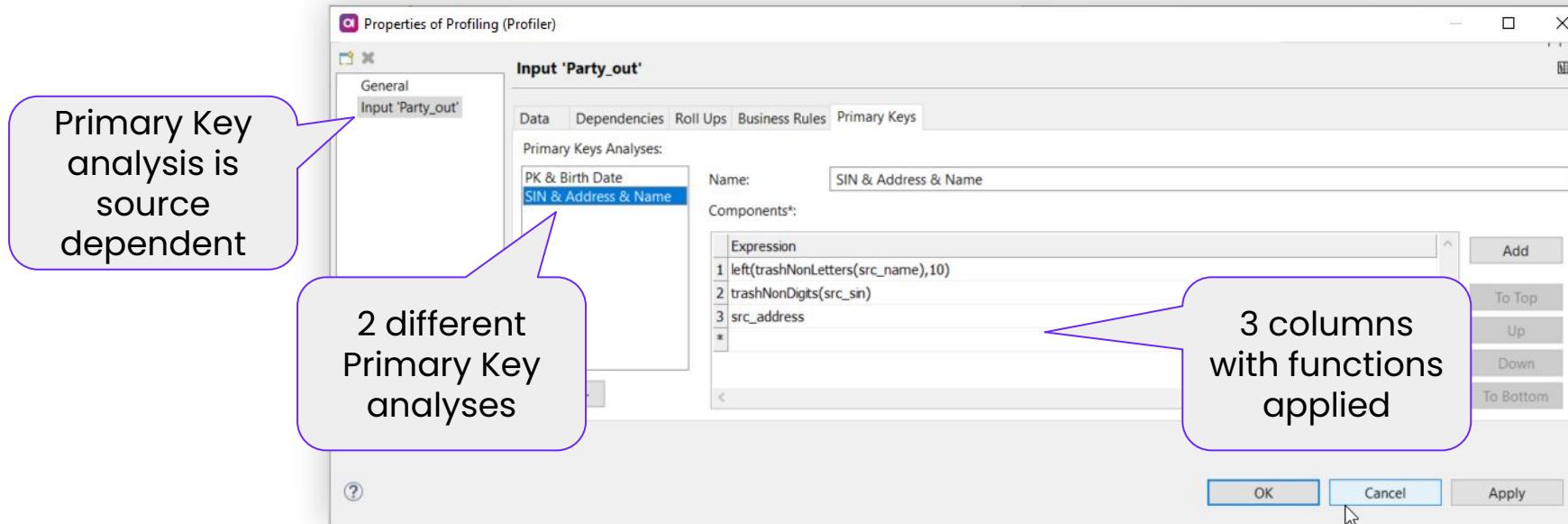
- Example: for a group of people, you can create a roll-up by gender (if this value is available). A sub-profile for males and females will then be created within the common profile.
- A roll-up can be a particular column value or any expression, such as a function applied to a column value
- Hierarchies of roll-ups can be created



Primary Key Analysis

Various primary key analyses can be created:

- Primary key can be a column, set of columns, or any expression or function applied to columns.

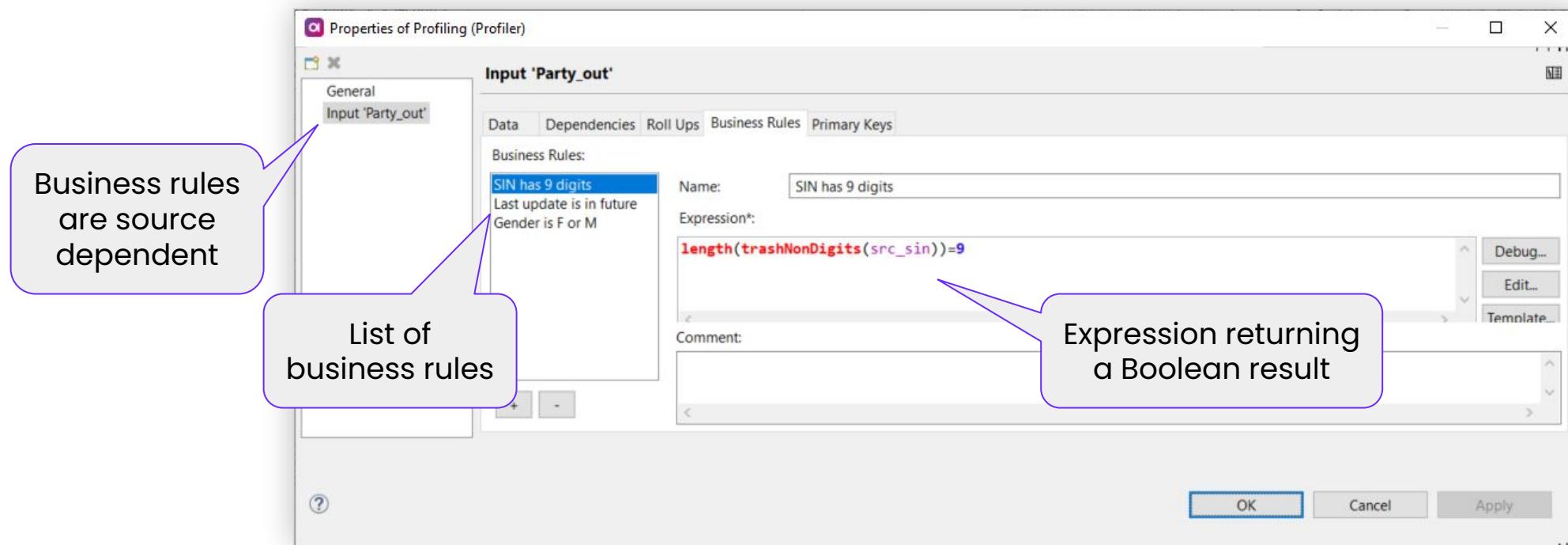


Primary key analysis ensures the uniqueness of identifiers that might be primary keys.

- i.e. reading from a DB without constraints.

Business Rules

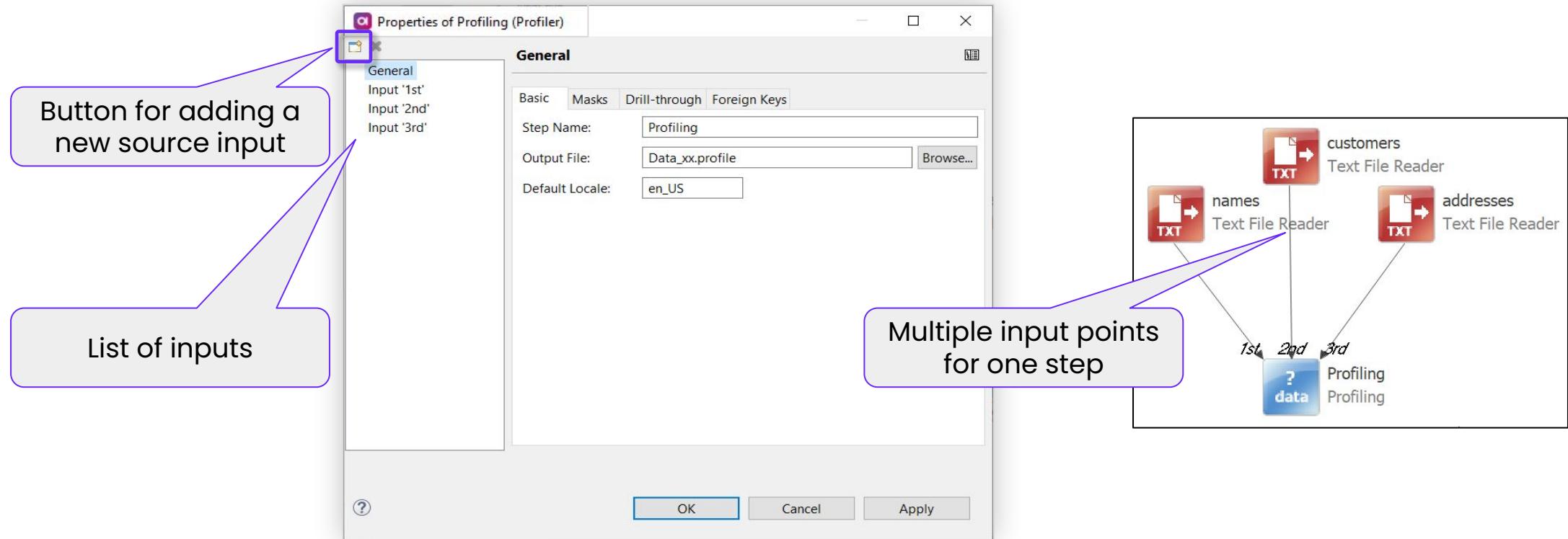
- Business rules are statements which can be evaluated as either **true** or **false**.
- All business rules are evaluated for each record.
- There is no limit to the complexity of business rules expressions, as long as they return a Boolean value at the end.



Multi-Table Profiling

The Profiling step can read from multiple input sources.

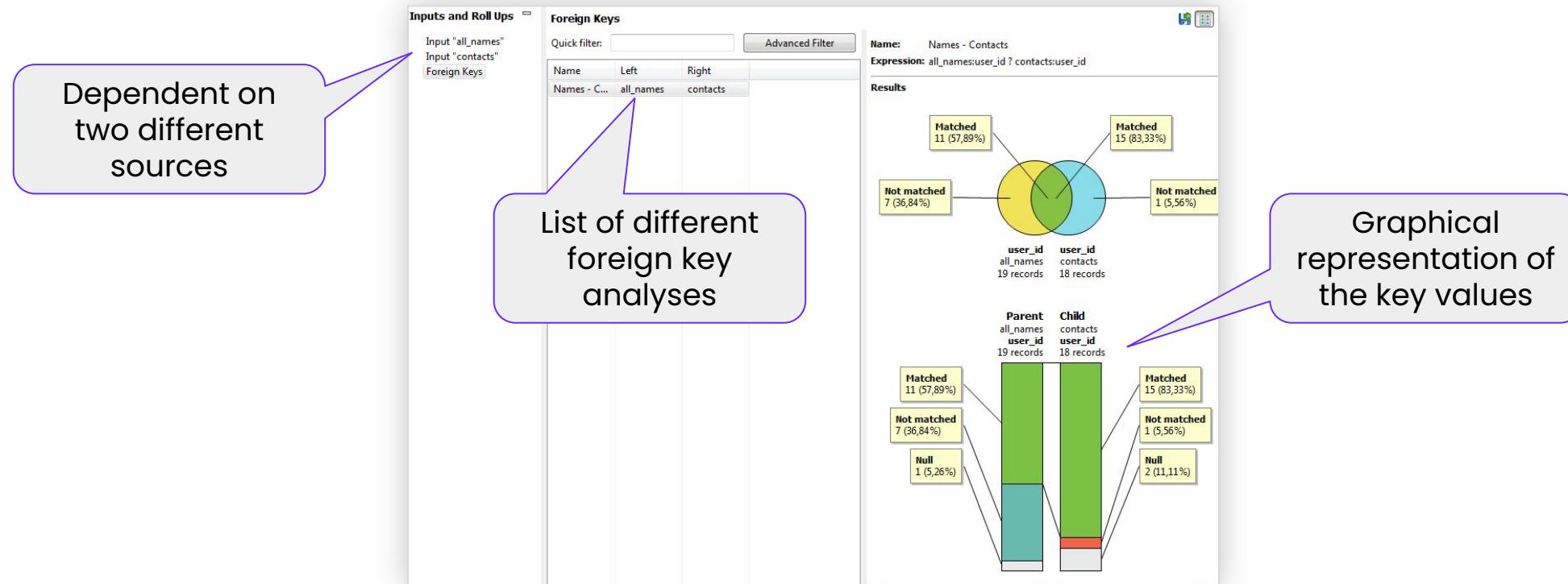
- Each source gets analyzed independently, but cross-source analysis is possible as well (e.g. foreign key analysis)



Foreign Key Analysis

Analysis between two of the sources going into the Profiling step.

- Requires definition of both the parent key and the child key.
- Graphic representation of relationship statistics (matched, non-matched, null in the parent, etc.).



Dependency Analysis

Searches for dependencies between determinant and dependent values.

- **Determinant** and **Dependant** can be column values or any expression.
- You can set thresholds for what is accepted as a dependency and what isn't.

The image shows two windows from the Ataccama Profiler tool. The left window is titled 'Properties of Profiling (Profiler)' and contains a tab labeled 'Input 'in''. It displays a 'Dependencies' section with a single entry: 'ID defines NAME'. This entry has a 'Determinant*' field containing the expression '1 squeezeSpaces(toString(src_customer_id))' and a 'Dependants*' field containing 'src_customer_name' with a 'Threshold' of 90. A purple callout points to this section with the labels 'Dependencies list', 'Determinants', 'Dependants', and 'Acceptability threshold'. The right window shows the results of the analysis. It includes a pie chart with the following data: Non-trivial dependencies (80,66%), Trivial dependencies (19,32%), and Violations (0,02%). Below the pie chart is a table of 'All Records' and another table for 'Non-trivial Records'. A purple callout points to the pie chart with the label 'Pie chart with results'. Another purple callout points to the bottom tables with the label 'Counts on: • Trivial • Non-trivial • Null • Violations'.

Pie chart with results

Counts on:

- Trivial
- Non-trivial
- Null
- Violations

Determinant V...	Rows Count	%	Distinct Count
Total	36.084	100,00%	14.551
Null	0	0,00%	0
Violations	8	0,02%	2
Dependencies	36.076	99,98%	14.549
Trivial	6.970	19,32%	6.970
Non-trivial	29.106	80,66%	7.579
Non-trivial null	0	0,00%	0

Determinant V...	Rows Count	%	Distinct Count
Total	29.114	100,00%	7.581
Dependencies	29.106	99,97%	7.579
Null dependen...	0	0,00%	0
Violations	8	0,03%	2

Sharing Results

Results can be simply shared via the profile files.

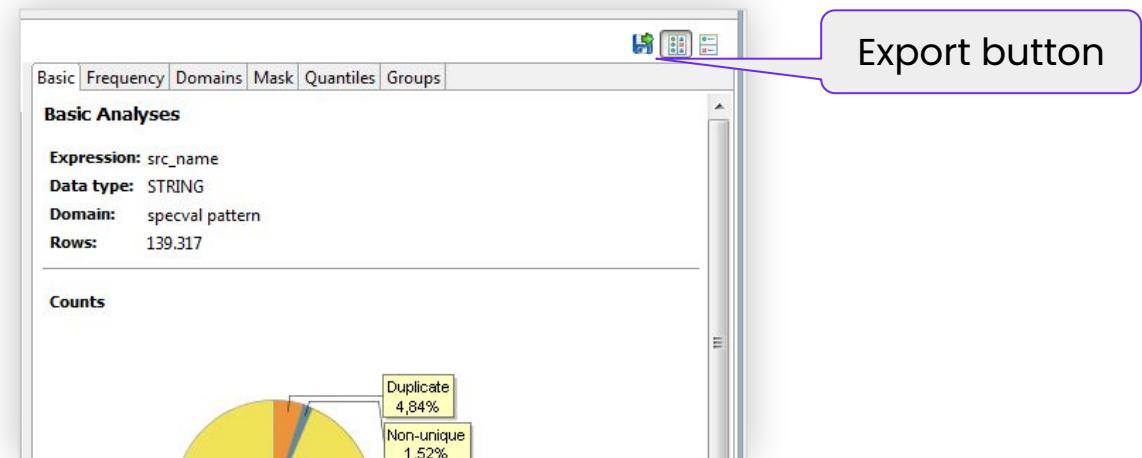
- Profile files have a very small footprint as they only store aggregations and statistical information, not the data itself.

ONE Desktop plans can also be shared.

- Plans are abstractions of the ONE Desktop GUI constructed from an XML file.
- XML files can be stored in a VCS and allow access to all parties.
- Because it is a plain text file, it is easy to keep track of changes and document them.

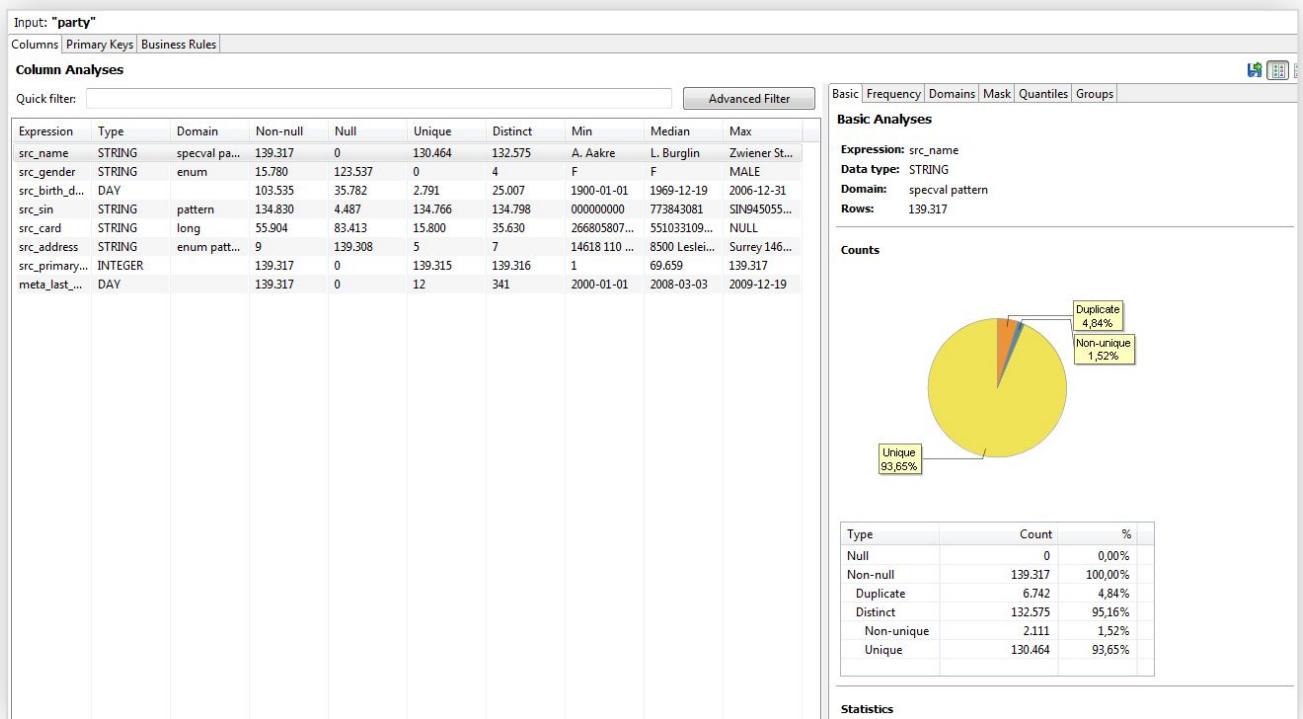
Profiles can also be exported to HTML via the Export button

- Allows for partial export
- Export to HTML or XML
- Can be auto-zipped



Topic Highlights

- Profiling provides the user with a **summary of basic statistical information**: unique values, duplicates, dependencies, frequency analysis, masks, and patterns.
- Allows the evaluation of business rules and observation of patterns in data, **to determine future steps for data cleansing**.
- Provides a **drill-through** feature, allowing quick access to sliced data samplings.



Memory Refresher #3 Data Profiling



ataccama

Match the terms with their meanings:

A. Frequency Analysis

B. Mask Analysis

C. Quantiles

D. Groups

D. Groups

1. Shows the number of times that each non-null frequency count is repeated in the selected column.

C. Quantiles

2. Displays the data values that occur at designated intervals in the ordered data set.

B. Mask Analysis

3. Shows the syntactic patterns of the data.

A. Frequency Analysis

4. Shows the number of times each value occurs in the data.

Lab Exercise #2

Data Profiling



Flow Control

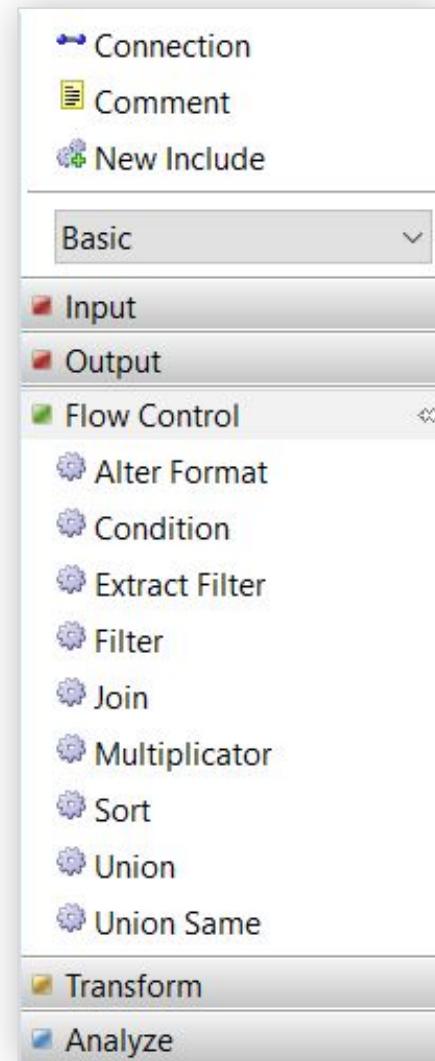


ataccama

Flow Control Steps

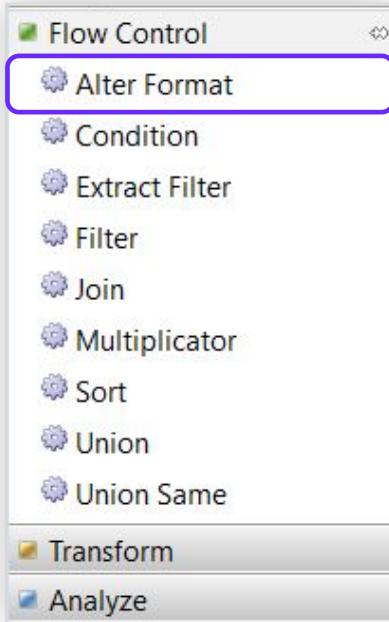
Flow control steps can be used to:

- Add or remove columns from the data flow.
- Split the data flow.
- Extract or segregate certain data sets.
- Join two or more data flows.





Alter Format Step



Alter Format

Adds or removes columns from the data flow.

This step is used to include additional or temporary columns that are not present in the original source input data. Common usage is to have two **Alter Formats** in a plan – one at the beginning to create new columns and one at the end to remove the unwanted ones not designed for output.

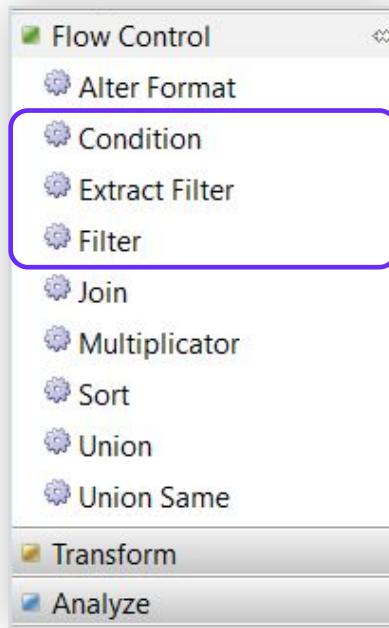
Typically used in every plan.



Alter Format Step



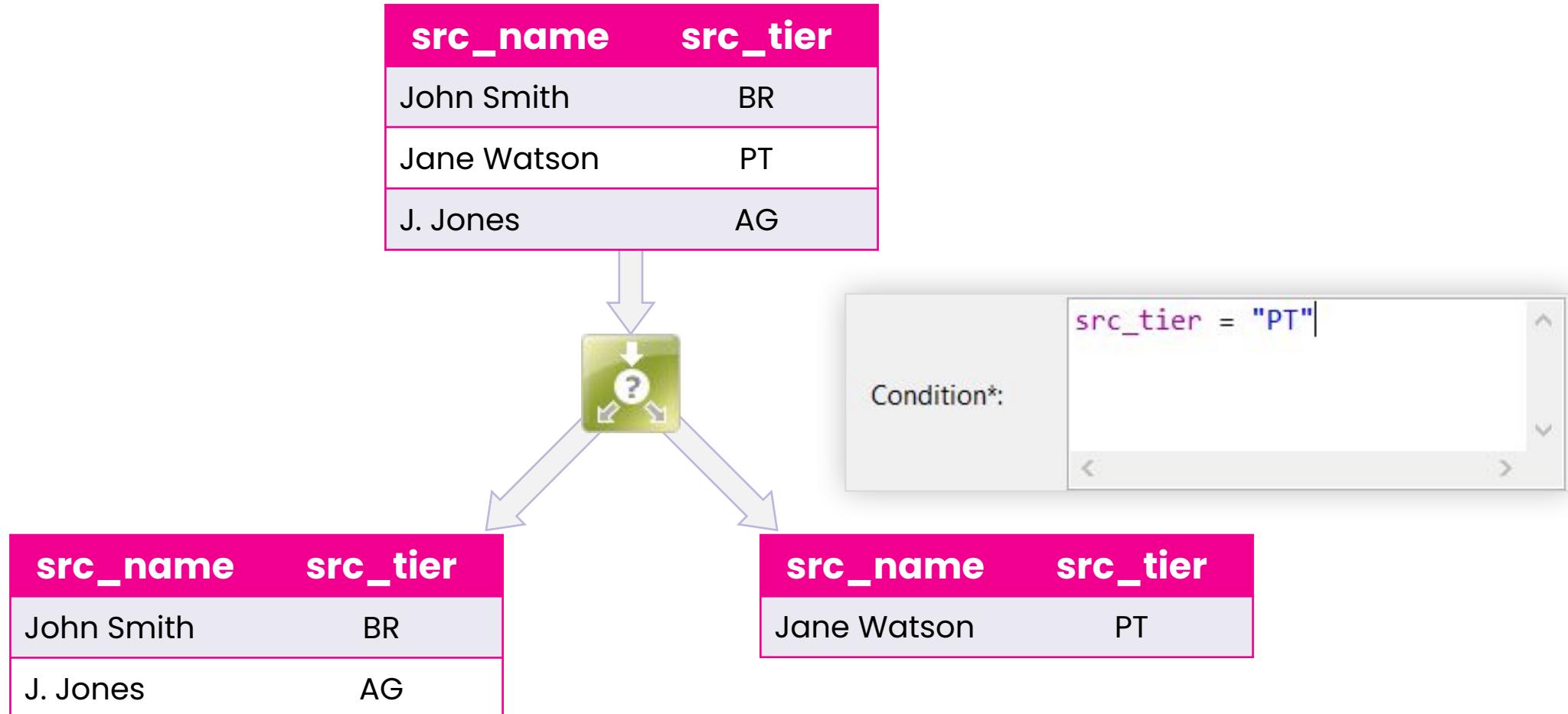
More Flow Control Steps



Condition	
	Allows branching of the data flow – divides into 'true' and 'false' outputs based on the result of specified condition. Usually used to extract certain data set (such as bad records) for manual verification.
Filter	
	Filters out records which do not satisfy a given condition. The remaining records are lost.
Extract Filter	
	Conditional branching of the data flow (same as condition), but the original flow is left untouched (all records are still there). Useful when you need to leave the original flow intact.

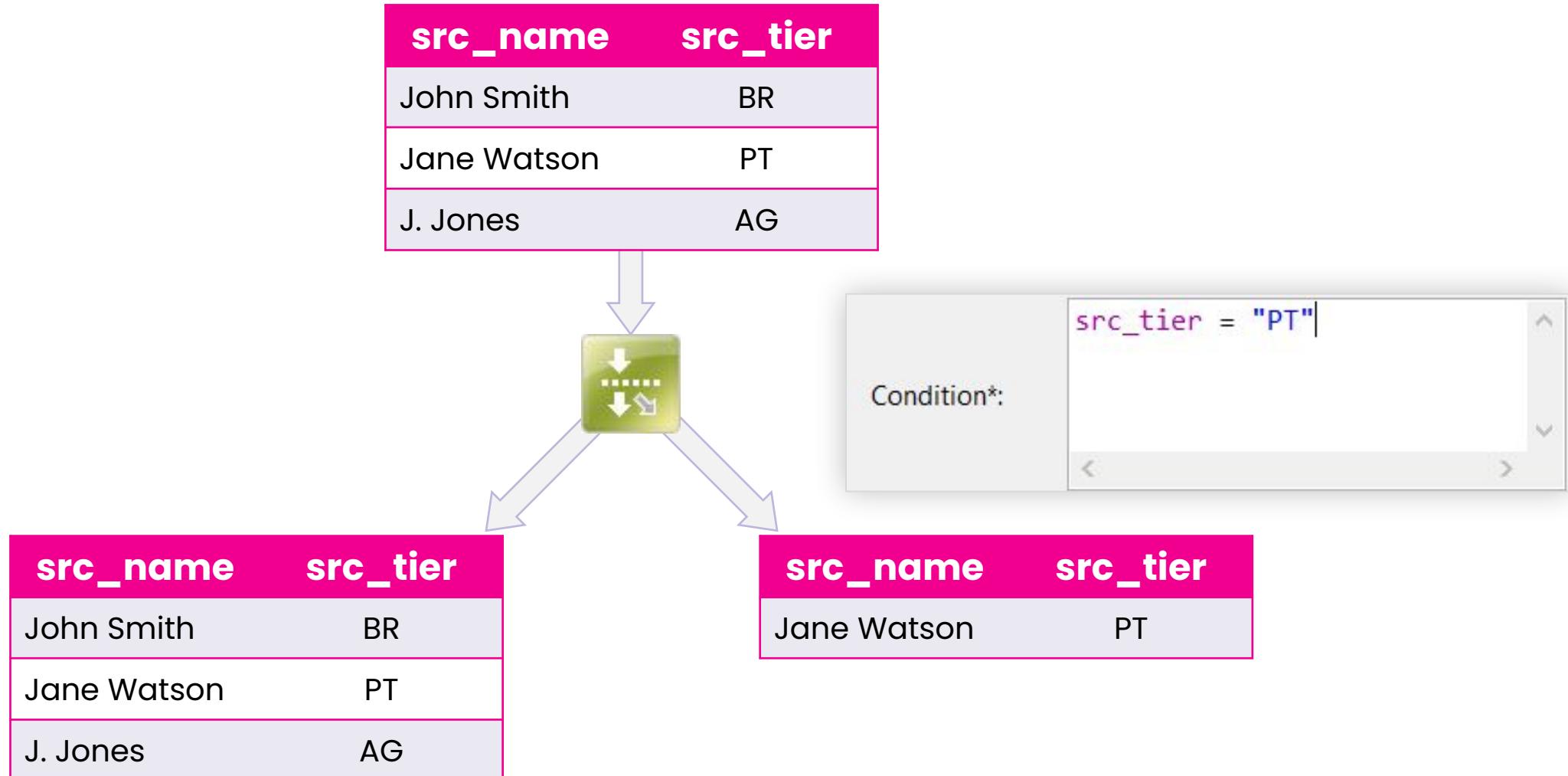


Condition Step

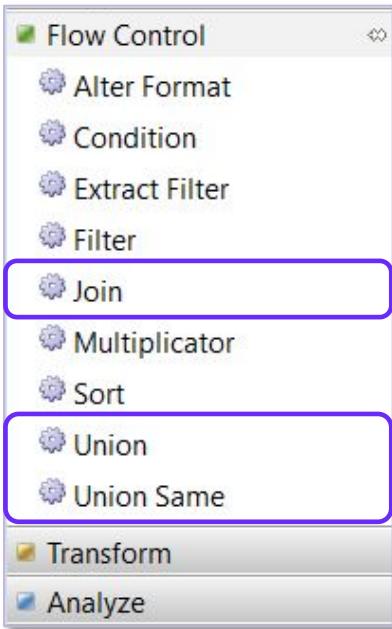




Extract Filter Step



Other Flow Control Steps



Join	
	Works in the same way as SQL JOIN. You must specify a key column and type of join (LEFT, RIGHT, INNER or OUTER).
Union	
	Works in the same way as SQL UNION. Apply a vertical merge of selected columns.
Union Same	
	Special type of Union step for working with identical data streams – all columns must have the same name and type.



Join Step

customer_id	src_name
1	John Smith
2	George Kent
3	Jane Watson

customer_id	src_address
1	8500 Leslei St
3	20 St John Ave
3	1221 Queens Blvd
4	49 River Ave



Left Key*:

Right Key*:

Join Type*:

Match Strategy*:

customer_id	src_name	src_address
1	John Smith	8500 Leslei St
3	Jane Watson	20 St John Ave
3	Jane Watson	1221 Queens Blvd



Union Step

src_name	src_address
John Smith	8500 Leslei St
George Kent	27 Town Square

src_name	residence
Jane Watson	20 St John Ave

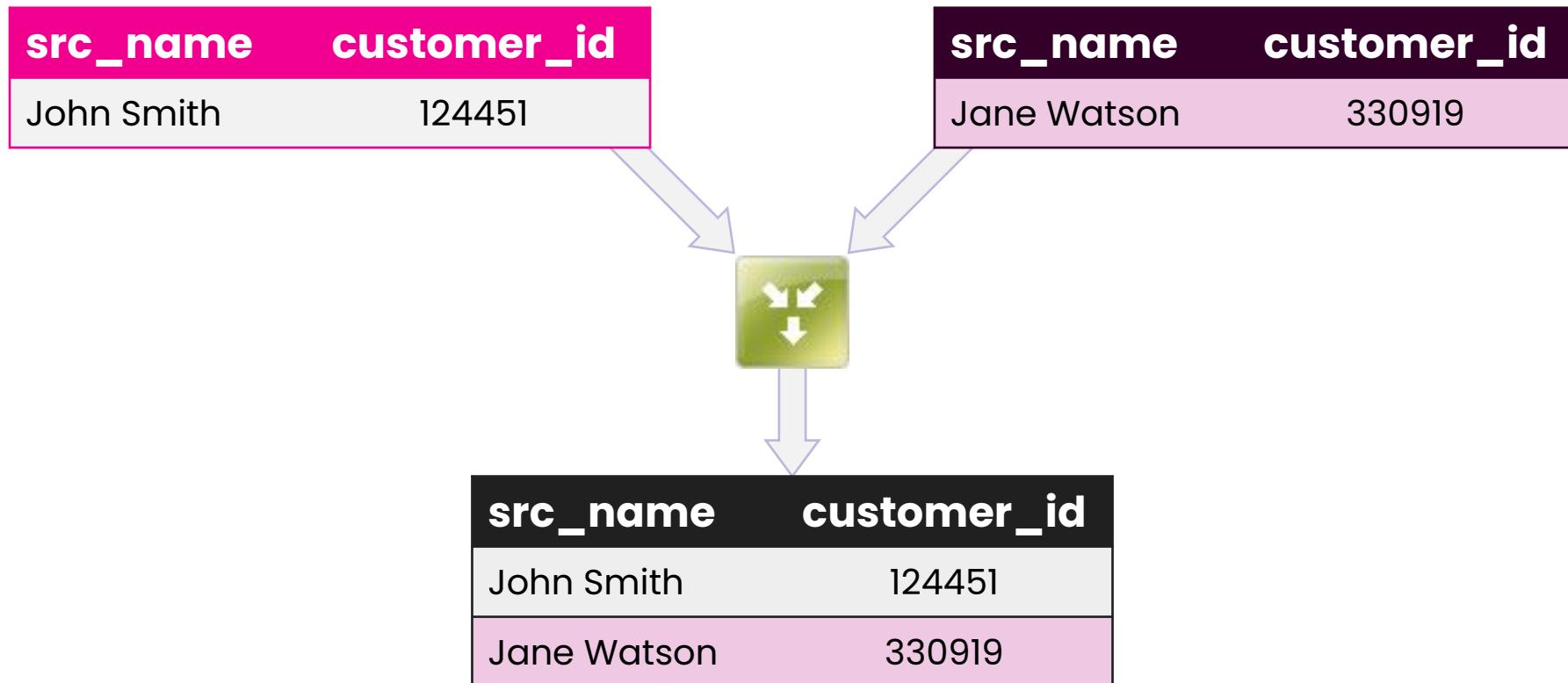
Column Mappings:

Dest	In_a	In_b
1 src_name	src_name	src_name
2 address	src_address	residence
*		

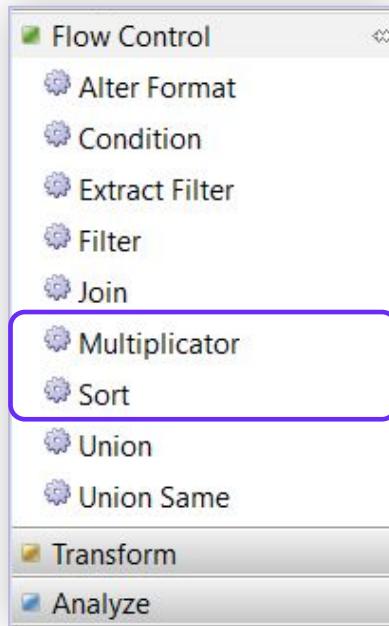
src_name	address
John Smith	8500 Leslei St
George Kent	27 Town Square
Jane Watson	20 St John Ave



Union Same Step



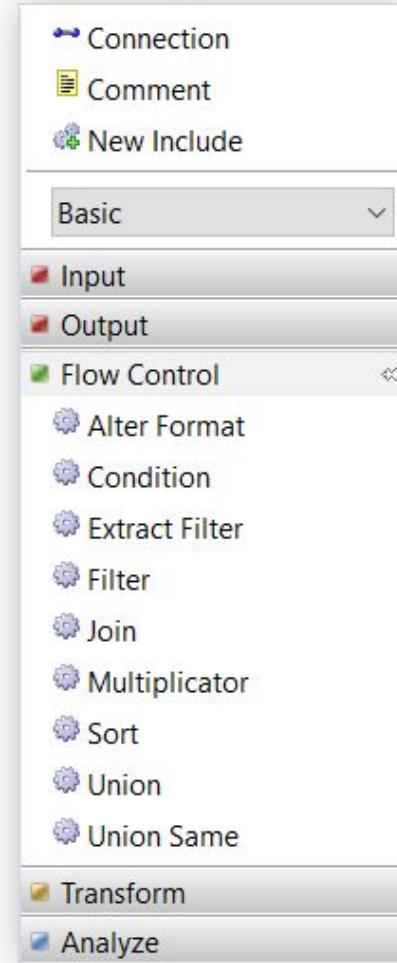
Multiplicator & Sort Steps



Multiplicator	
	Multiplies the data stream into more identical streams. Several connections can be connected from the 'out' endpoint.
Sort	
	Sorts the data in the stream. Sometimes used before writing the output to a file.

Topic Highlights

- Add or remove columns from the data flow (**Alter Format**).
- Split the data flow (**Multiplicator**).
- Extract or segregate certain data sets (**Condition, Filter, Extract Filter**).
- Join two or more data flows (**Join, Union, Union Same**).



Memory Refresher #4 Flow Control



ataccama

Extract Filter retrieves records which pass a specific condition.

- TRUE
- FALSE

Filter splits the data flow into two branches based on a logic expression.

- TRUE
- FALSE

Union Same merges multiple data flows with the same format into a single data flow.

- TRUE
- FALSE

Lab Exercise #3

Flow Control



Best Practices & Conventions



ataccama

Data Quality Conventions

- **Understand your data!** The first step is always to analyze the data you have.
- **Never trust any data or any assertion.**
 - The data that was supposed to be perfect is often far from perfect.
 - Analyzing the data from all angles could shed light on new problems.
 - Assumptions are dangerous.
- Use the **Profiling**  step to perform an initial analysis.
- Names of plans and components **shouldn't have spaces** – the best word separator to use is the underscore.
- Column name lengths should not exceed 30 characters (including prefix).

Naming Conventions

General Data Quality conventions dictate that:

- Source data should never be modified or transformed.
- A new set of columns (with prefixes – see below) should be created for each meaningful value in a data record .
- These columns will store transformed values, standardized values, data quality indicators, and audit information.

Basic column prefixes:

- src_ The source value
- pur_ A value used internally for transformations
- std_ A value that has been standardized against a dictionary
- cio_ The best possible value after all processing (STD or PUR)
- mat_ A simplified version of data used for matching purposes only
- dic_ A value that was verified in a dictionary

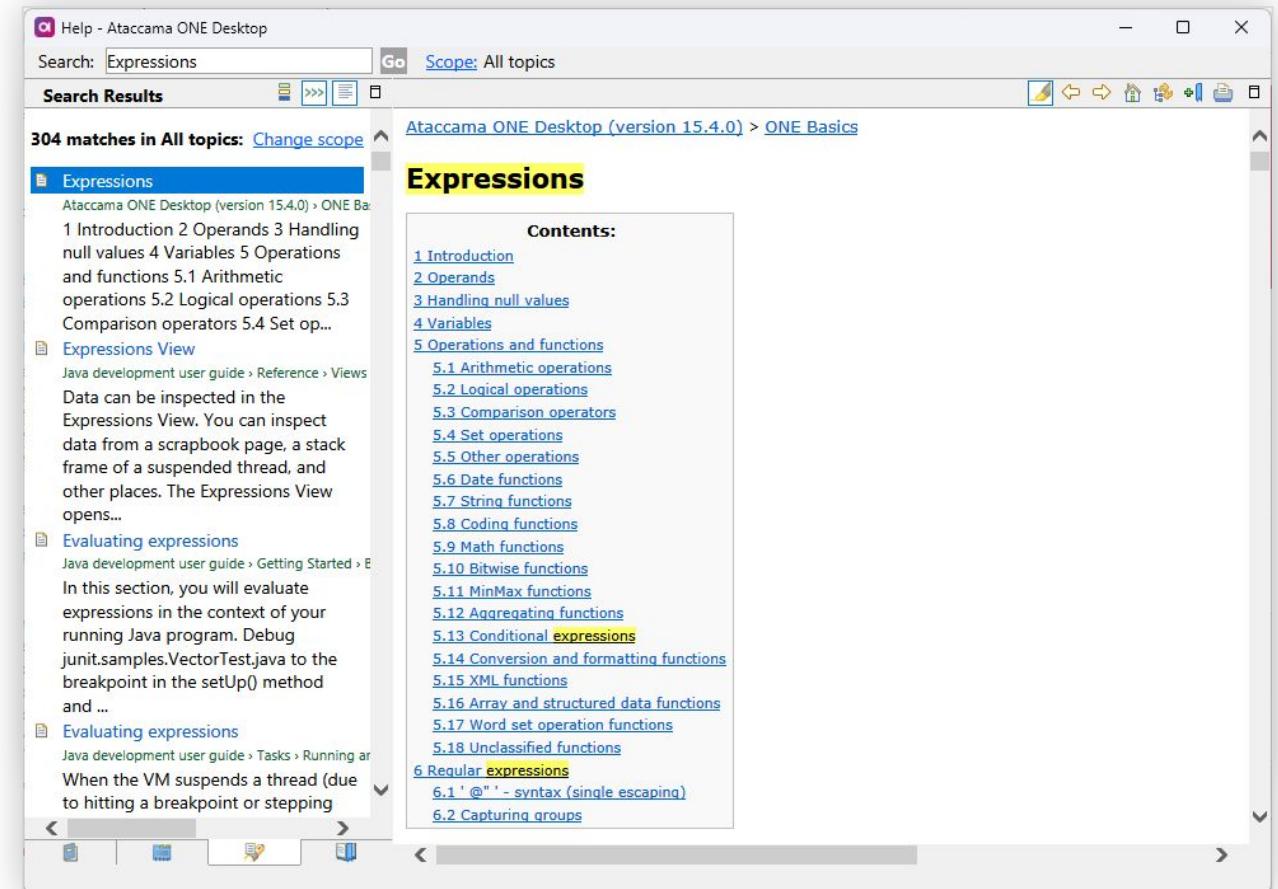
Expressions



ataccama

Ataccama Expression Language

- Write your own transformations.
- Full reference in the built-in help.
- Tutorials (Using functions).
- Expression templates.

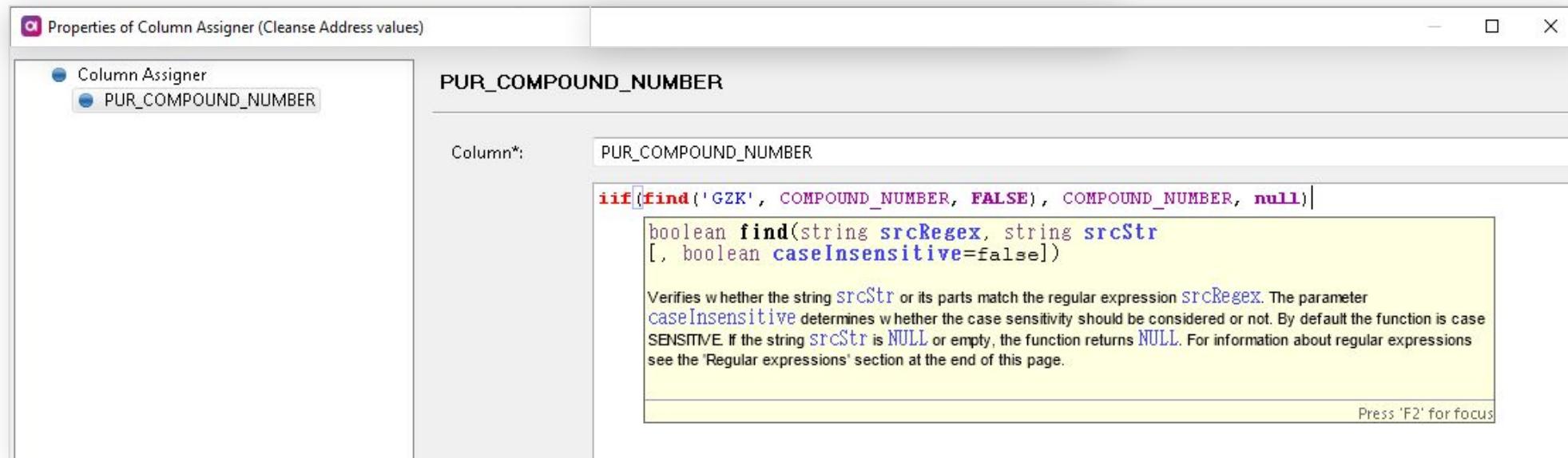


Expressions Categories

- String functions.
- Date functions.
- Conditional expressions.
- Conversion and formatting functions.
- Math functions.
- Word set functions.
- And more; a complete list of expressions can be found in ONE Desktop's Help > Help Contents menu – enter the keyword "Expressions".

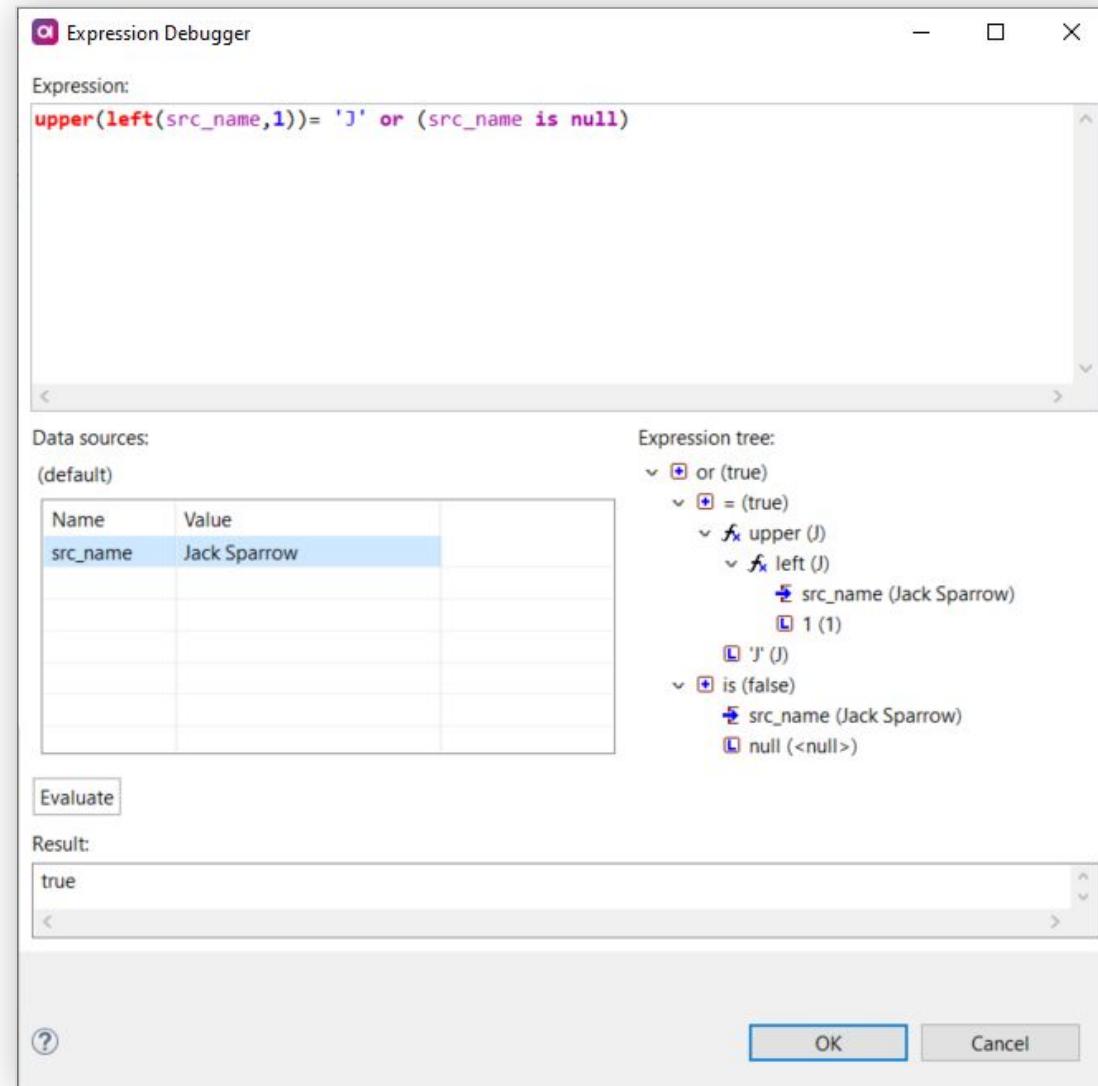
Expression Content Assist

- Every time you write a function/expression in ONE Desktop, you can invoke content assist by pressing.
 - **CTRL + Spacebar**
- You can also place your cursor over the function name to see help for that function.



Expression combinations

- Nesting expressions
 - **upper (left (src_string, 2)**
- Combining
 - **AND, OR, XOR,...**
- Evaluation priorities
 - visualized in Debug Expression Tree
- Conditionals - **If, Case, NVL**



Expressions syntax

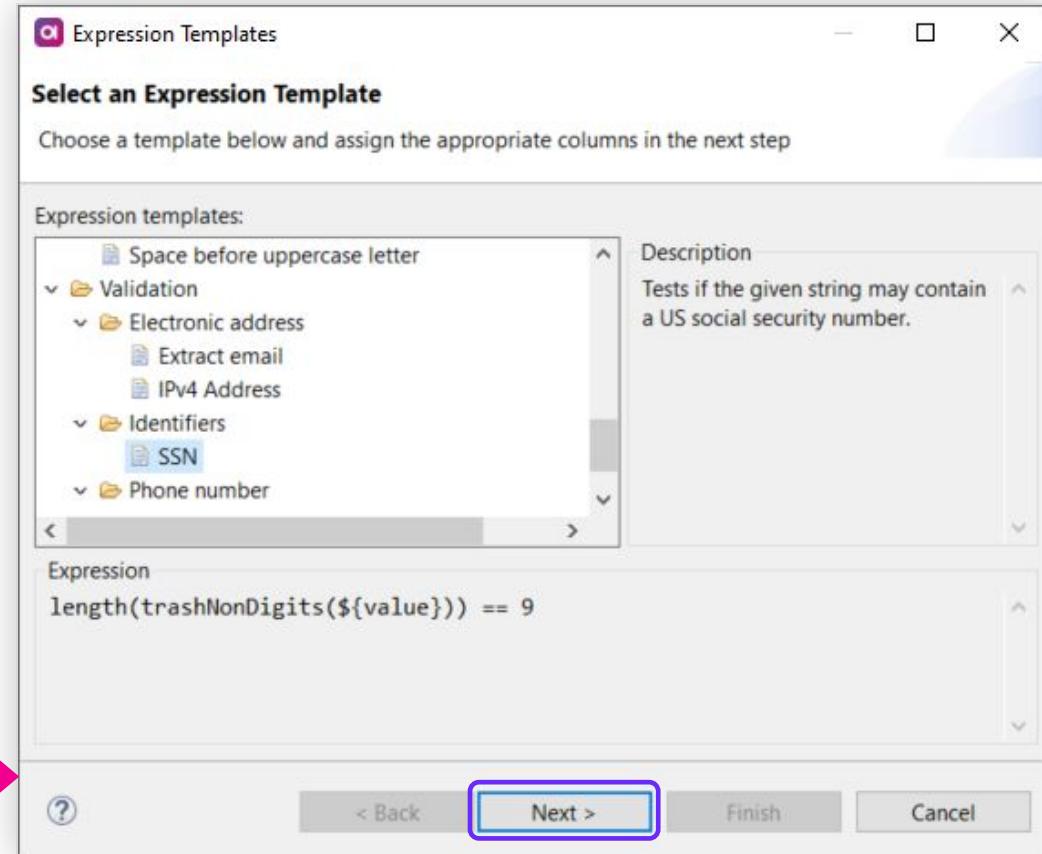
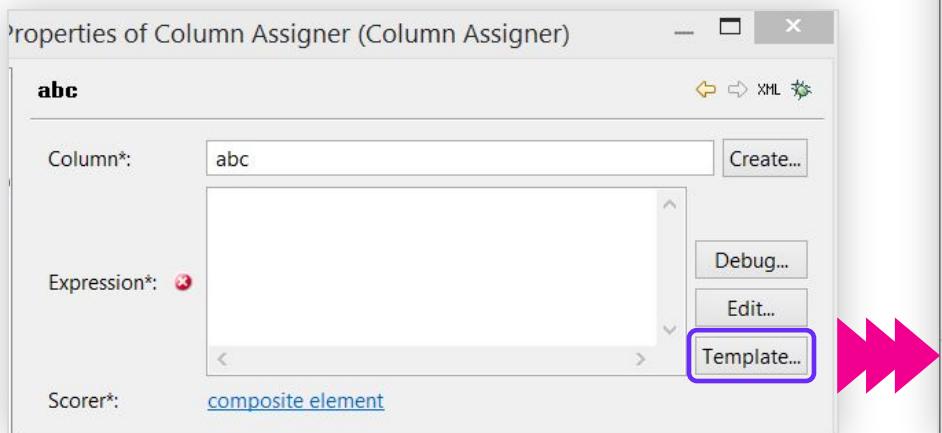
- In-built functions (length, right,..)
- Variables:
 - Simplify complex expressions.
 - Automatically typed (String, Integer).
- Evaluation:
 - The value of the last expression is returned.
 - Expression should produce a value matching the output column type.

```
number_of_chars := 10;  
  
src := src_num_as_string;  
  
filler := '0';  
  
iif(      //boolean predicate  
      length(src)<number_of_chars and src is not  
      null,  
      //true  
      right(replicate(fillers, number_of_chars) +  
            src, number_of_chars),  
      //false  
      Value  
    )
```

Expression Templates

- Ready to use expressions and their combinations.

Note: \${value} is replaced by a column name after pressing the 'Next >' button.



Aggregations in Expressions

- **Context specific expressions (Aggregations)**

- Count Functions

- [count\(\)](#)
 - [countDistinct\(\)](#)
 - [countDistinctif\(\)](#)
 - [countif\(\)](#)
 - [countNonAsciiLetters\(\)](#)
 - [countUnique\(\)](#)
 - [countUniqueif\(\)](#)

- Concatenate

- [concatenate\(\)](#)
 - [concatenateif\(\)](#)

- Averages

- [avg\(\)](#)
 - [avgif\(\)](#)

- Selectors

- [first\(\)](#)
 - [firstif\(\)](#)
 - [last\(\)](#)
 - [lastif\(\)](#)

- Sum

- [sum\(\)](#)
 - [sumif\(\)](#)

- Min/Max

- [maximum\(\)](#)
 - [maximumif\(\)](#)
 - [minimum\(\)](#)
 - [minimumif\(\)](#)

- Modus

- [modus\(\)](#)
 - [modusif\(\)](#)

Expressions on strings

Edit distance on string:

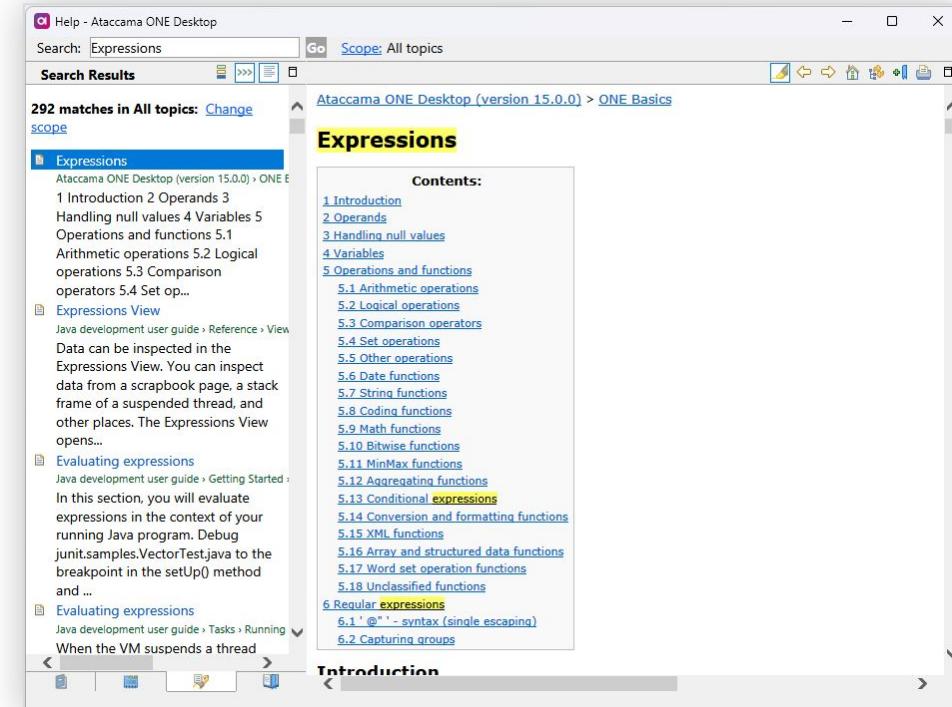
- Levenshtein (string difference - how many "typos")

String set operations:

- Set functions accept string(s) and interpret them as a set of items:
 - `set.sort("Orange Banana Kiwi", " ") --> "Banana Kiwi Orange"`
- Union, difference, intersection of 2 strings
- Size, Sort
- MapExp, SumExp

Topic Highlights

- **Expression Content Assist** can be invoked assist by pressing **CTRL + Spacebar**.
- **Expression combinations**; expressions can be nested, combined with **AND, OR, XOR**, etc., and contain conditionals such as **IF, CASE, NVL**.
- **Expression Templates**, ready to use expressions and their combinations.



Memory Refresher #5 Expressions



Can Expression Templates be used as-is, without modification?

- 1.** Yes.
- 2.** No.

Expression Content Assist can invoked by:

- 1.** Hovering your mouse over the expression input field.
- 2.** Pressing CTRL SPACEBAR on your keyboard.
- 3.** Shouting loudly until a colleague comes to your assistance.

Can a logical expression contain inputs of different types, such as string or integer?

1. Yes.
2. No, all inputs must be of the same type
3. Yes, if the inputs are converted to the same type.

Lab Exercise #4

Data Quality Indicator

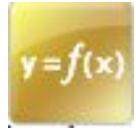


Data Transformation



ataccama

Basic Steps for Direct Data Transformation



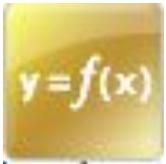
- **Column assigner**
 - Allows data mapping (maps values from input columns to working columns, etc.).
 - Allows data value transformations with the use of IDE expressions.



- **Transliterate**
 - Translates any character or string into another defined character or string.



- **Regex matching**
 - Parses and extracts values from input strings via defined patterns and regular expressions.



Column Assigner Step

Properties of Column Assigner (Column Assigner)

Column Assigner

General Advanced

Id: Column Assigner

Assignments:

Column	Expression	Score
1 pur_name	upper(src_name)	Score
*		

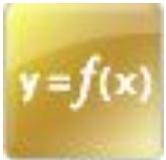
Add To Top Up Down To Bottom Fill Columns...

Comments [\(Hide\)](#)

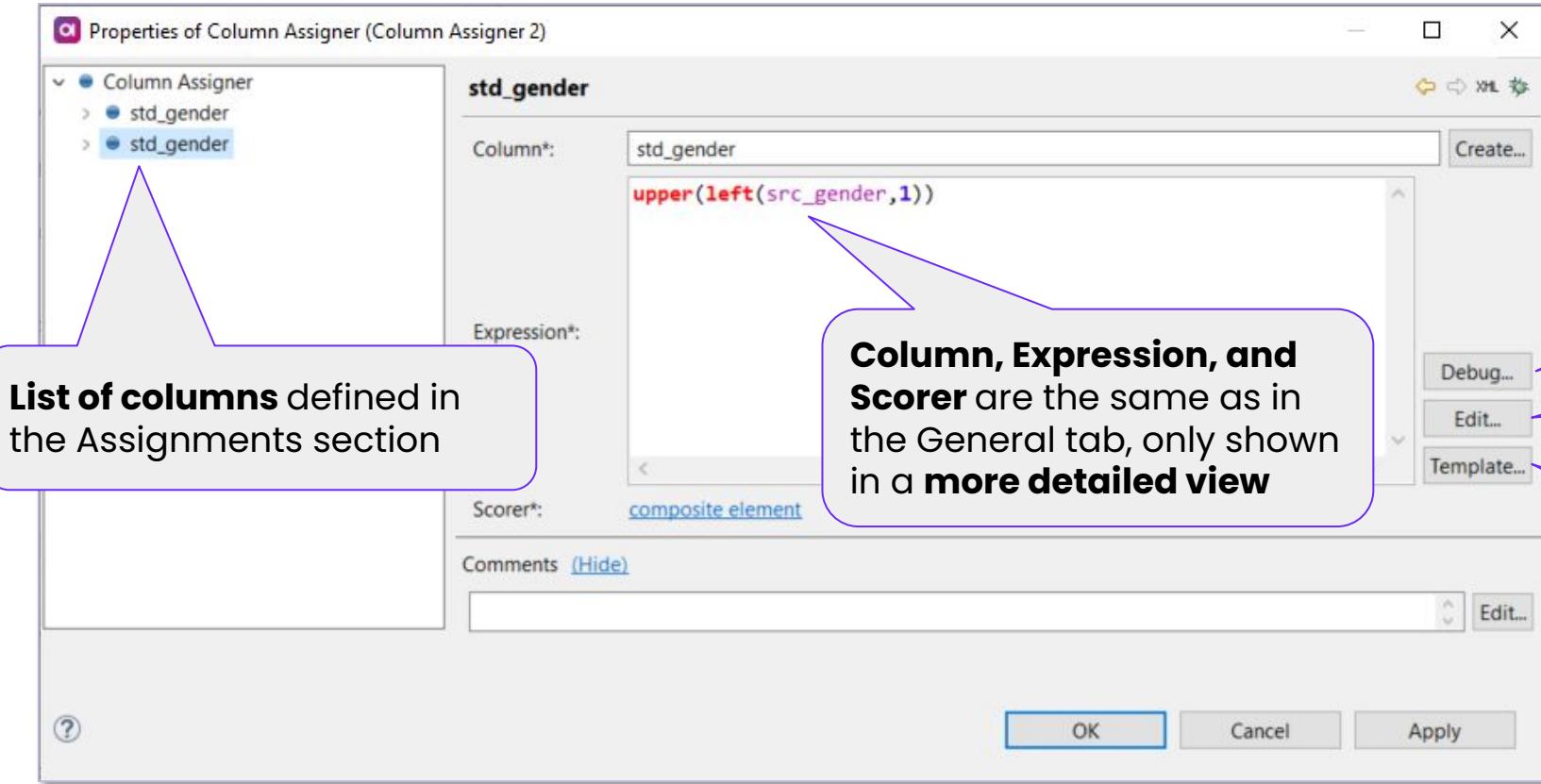
OK Cancel Apply

Assignments:

- **Mapping or transformation** operations
- **Scoring** of transformations
- **Comments** for each performed transformation



Column Assigner Step



List of **columns** defined in the Assignments section

Column, Expression, and Scorer are the same as in the General tab, only shown in a **more detailed view**

Debug allows you to debug your defined expression

Edit opens the Expression text in a separate window

Template opens a list of common expressions and operations and allows you to use them as a template



Transliterate Step

Properties of Transliterate (Transliterate)

Transliterate

General Advanced

General

Id: Transliterate

Other

Rules*:

From	To	Comment
1 GSK	SSK	
2 ABC	DEF	
3 Smith	Smyth	
*		

Add To Top Up Down To Bottom

Columns*:

From	Scorer	To
1 src_name	Scorer	std_name
*		

Add To Top

Comments [\(Hide\)](#)

Edit...

OK Cancel Apply

Defined **transliteration rules**

Defined **input column From** and
result column To
(both may be the same column)



Regex Matching Step

Properties of Regex Matching (Regex Matching)

Email

Name*: Email

Pattern*: `^([\d\p{L}!#$%^&_=?^-_|{~-}{1,127})@([\d\p{L}!\.-]{1,127})$`

Case Insensitive:

Partial Match:

Multiplicative:

Result Columns:

Name	Substitution	Comment
1 user	\$1	
2 domain	\$2	
*		

Comments [\(Hide\)](#)

OK Cancel Apply

Pattern definition

Mapping of matched parts of the input string



Regex Matching Step

Regular expressions (Example usage):

- Parsing values in various formats (dates, phone numbers, identification numbers).
- Extracting values from fields containing multiple types of information (phone number mixed with address).
- Transforming input values.

The matched and unmatched parts of the input text can be accessed via special variables:

- **\$& or \$0** string matching the whole regular expression.
- **\$1, \$2, ...** capturing groups (parts of matched string).
- **\$`** part of input string before the match.
- **\$'** part of input string after the match.



Regex Matching Debugger

The screenshot shows the "Regular Expression Debugger" window. The input text is `!/?AB34578T**`. The regular expression is `([A-Z]{2}\w+)(\d+)`. The "Partial match" checkbox is checked, while "Ignore case" is unchecked. The "Evaluate" button is visible. The results table displays one match with elements \$:, \$& (\$0), \$', \$1, and \$2, corresponding to the tokens !?, AB34578, T**, AB3457, and 8 respectively. Navigation buttons for "Prev Match" and "Next Match" are shown at the bottom of the results table.

Element	Value
\$:	!/?
\$& (\$0)	AB34578
\$'	T**
\$1	AB3457
\$2	8

Input string

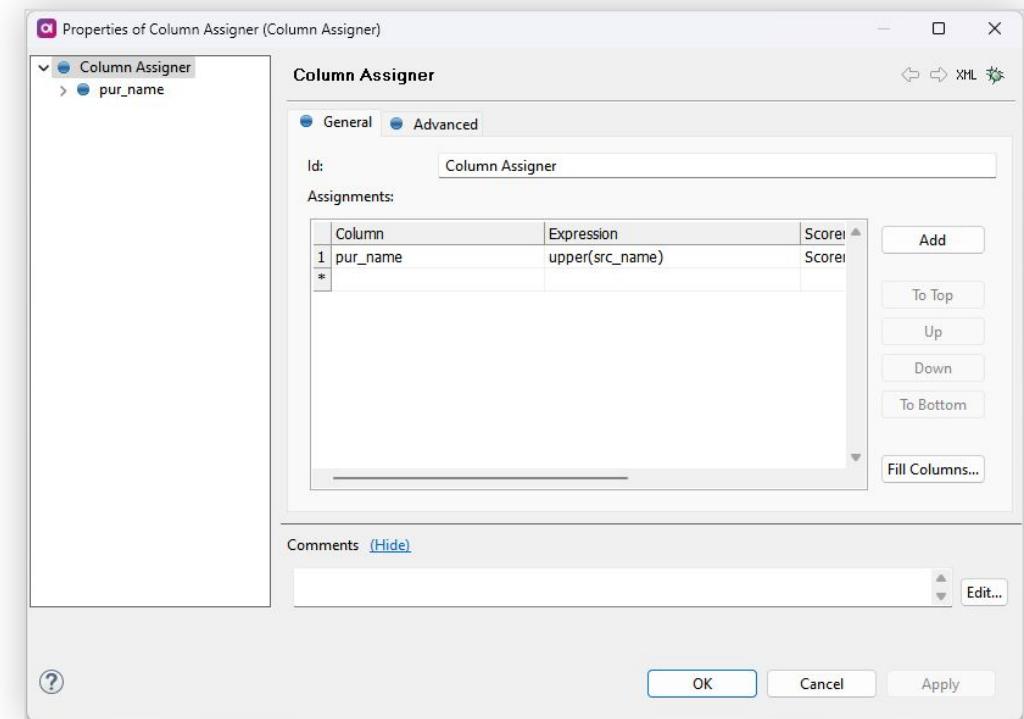
Pattern for matching

Additional settings

Matched parts of the input

Topic Highlights

- **Column Assigner** inputs the result of a defined expression to a specified column.
- **Transliterate** transforms a specified string into another string.
- **Regex matching** parses data and extracts values in the input string using defined patterns of Regular expressions.



Memory Refresher #6 Data Transformation



Attribute names can be dynamically changed in the Column Assigner.

- TRUE
- FALSE

The Column Assigner can be used to add a new column.

- TRUE
- FALSE

Lab Exercise #5

Data Transformation



Step Debugging



ataccama



Debugging Overview

Debugging expressions:

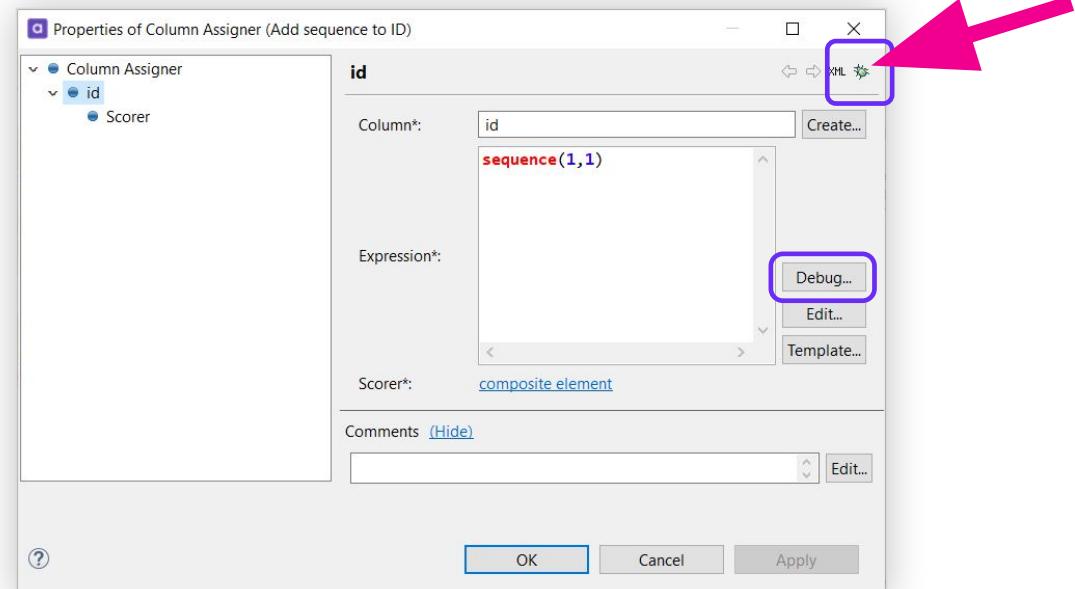
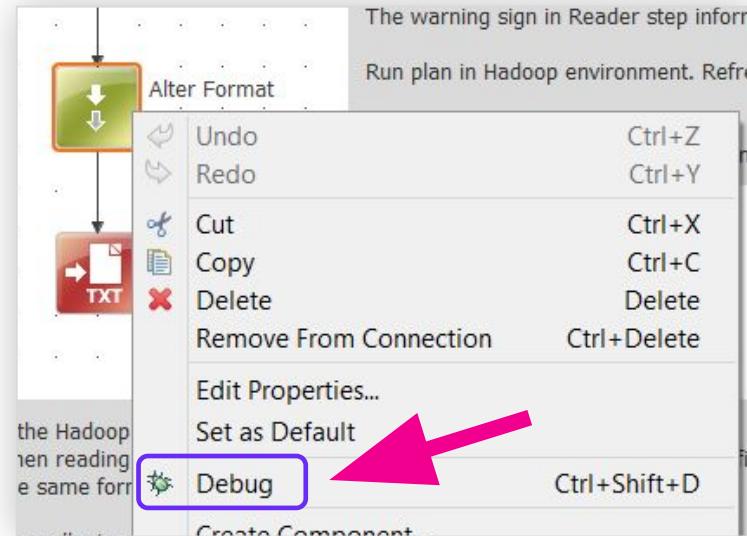
- Performed in a special form.

Debugging steps:

- Right-click on a step or click the bug icon in the upper right corner.
- Debugging a component step is also possible.

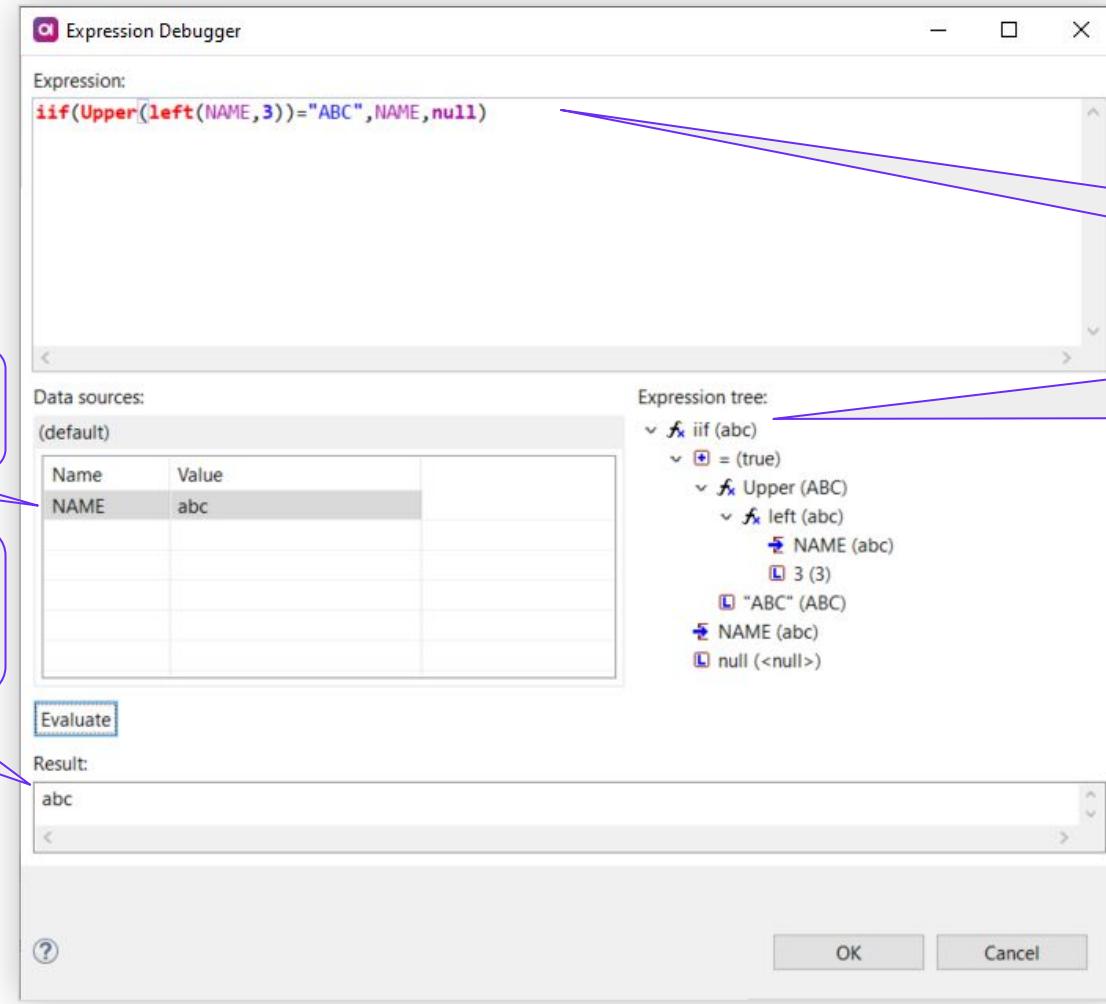
Testing some parts/whole part of the flow:

- Add Multiplicator + Text File Writer steps.





Debugging Expressions



Input value used for debugging

Result shows the input value transformed by the defined expression

Expression definition comment format: `/**/` or `//`

Expression tree shows individual operations in the expression step by step



Debugging Steps

You can debug the entire step for a given sample input:

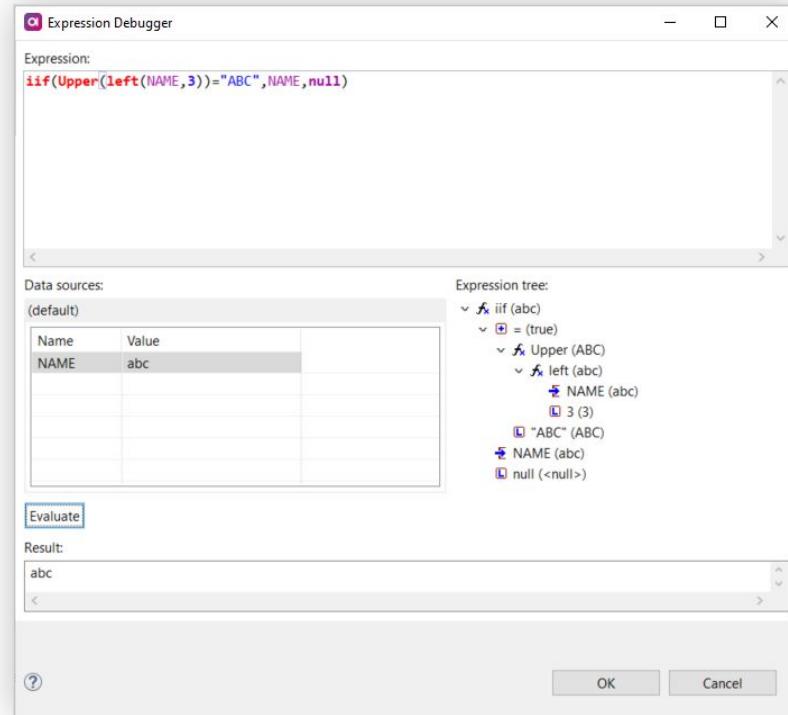
The screenshot illustrates the debugging process for a Column Assigner step. On the left, the data flow interface shows a sequence of steps: Alter Format, Column Assigner (highlighted with a yellow box), Column Assigner (highlighted with a yellow box), and Transliterate. A context menu is open over the second Column Assigner step, with the 'Debug' option highlighted by a pink arrow. This leads to the 'Step Debugger' window on the right, which displays sample input data:

	src_name	src_gender	src_birth_date	src_sin	src_card	src_address	src_email	src_
1	Kyle Smith		1963-11-22	555-555-555	4500111100001111	55 Neverland Dr	kyle.smith@ataccama.com	
*								

Below the input table is an 'Evaluate Input' button. The 'Properties of Column Assigner (Column Assigner)' dialog is also visible, showing the configuration for the 'pur_name' column. The 'Expression*' field contains the expression `upper(src_name)`, and the 'Scorer*' field is set to `composite element`. A pink arrow points from the 'Properties' dialog to the 'Debug...' button in the 'Properties of Column Assigner' dialog.

Topic Highlights

- Debugging can be performed on expressions, steps, components, and parts of or entire plans.
- Functions or regular expressions can be debugged by clicking on **Debug** inside the step.
- Likewise, steps can be debugged by right clicking on the step then selecting **Debug** from the context menu.
- Plans can be debugged by inserting a **Multiplicator step** and **Text Writer** step at the desired point to view the plan's results to that point.



Memory Refresher #7 Step Debugging



Regular expressions can be debugged
in the Regex Matching step.

- TRUE
- FALSE

You can run a testing process for the whole plan or component.

- TRUE
- FALSE

Is it possible to see intermediate results of an expression?

- 1.** No.
- 2.** Yes, during expression debug.
- 3.** Yes, during debug for a step.

Lab Exercise #6

Debugging



Scores & Explanations



ataccama

Scoring And Explanation Codes

- During scoring, every record is evaluated – scored – based on the quality of each value.
- The **higher the score, the bigger the difference between input and cleansed values**, and the lower the credibility of the input value.
- The **Cleansing Score** of any record is simply a sum of the cleansing scores of its individual attributes.
- Score is not self-explanatory by default and should be explained via **Explanation Codes**.

Best Practices For Scoring

Four scoring (quantification) levels:

Scoring level	Scoring description	Scoring result
0	The input value was not modified .	VALID
< 10,000	The input value needed small modifications (trim, squeeze spaces, uppercase, remove unsupported characters, safe replacements, etc.).	VALID
< 10,000,000	The input value needed complex modifications (not found in dictionary, does not meet official rules, major error, unsafe replacements, etc.).	UNSAFE / UNKNOWN
> 10,000,000	The input value (or pre-cleanse value after small modifications) is null .	NOT VALID

- The scoring value must be set in the specified range but does not have to exceed the maximum value in case the issue is repeated.
- It is recommended to use multiplications of 10 if you want to emphasize/differentiate the importance within one scoring level.

Topic Highlights

- During scoring, every record is evaluated – scored – based on the quality of each value.
- The higher the score, the bigger the difference between input and cleansed values, and the lower the credibility of the input value.
- It is recommended to use multiplications of 10 if you want to emphasize/differentiate the importance within one scoring level.
- Explanation should be without spaces.

Lab Exercise #7

Scores & Explanations



Lookup & Dictionaries



ataccama

Building and using lookups

1 Build or obtain a Lookup File

Reference data (source)

- List of correct values (names, addresses)
- List of replacements (wrong ☐ correct)
- In a text file or database table

Lookup file

- A DQC binary (dictionary) file prepared for high-performance lookup operations
- Created by a plan with a **Lookup Builder** step



2 Use the Lookup file in a plan

Steps in your plan with the Lookup file:



- **Lookup**



- **Apply Replacements**



- **Strip Titles**



- **Pattern Parser**

Lookups – Steps



- **Lookup Step**
 - Lookup structure: key + **additional columns** (optional)
 - Determines if the key exists in the lookup file
 - Enriches the input record with values from **additional columns** (if defined)



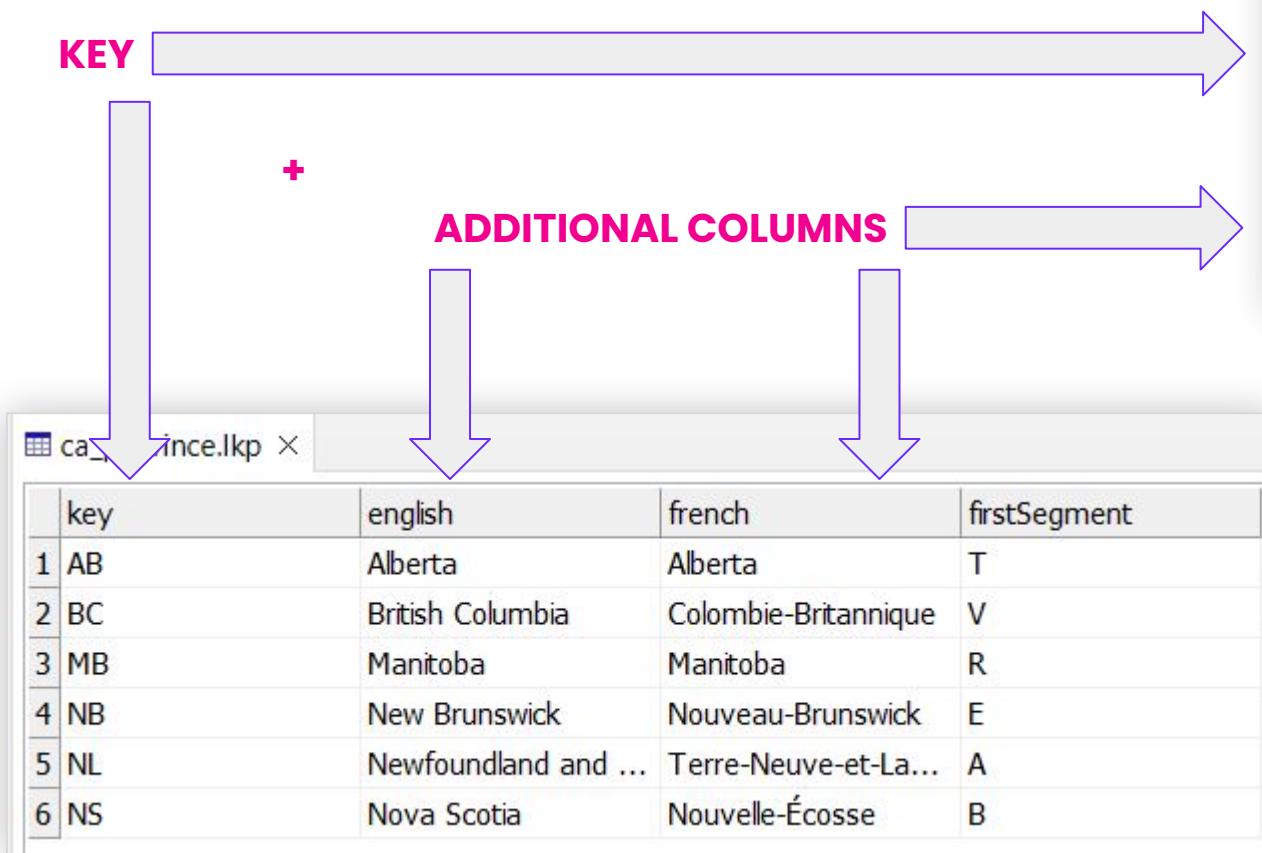
- **Apply Replacements**
 - Lookup structure: key + **replacement**
 - Replaces all occurrences of each key with the corresponding **replacement**



- **Strip Titles**
 - Lookup structure: key
 - **Removes** all occurrences of the key



Lookup Builder: File Structure



Main

File Name*:/data/ext/build/ca_province.lkp

Key*: key

Additional Columns:

Name	Expression	Comment
1 english		
2 french		
3 firstSegment		



Lookup Builder: Parameters

Lookup Builder

General

Duplicities: [REF] dic_asset_01.lkp

Other

Duplicities*: FIRST

Matching Value*: composite element

Approximative Index:

Best Distance Index:

Bidirect Approximative Index:

Compressed:

User Metadata*: composite element

Main

File Name*: ..\lkp\dic_asset_01.lkp

Key*: KEY

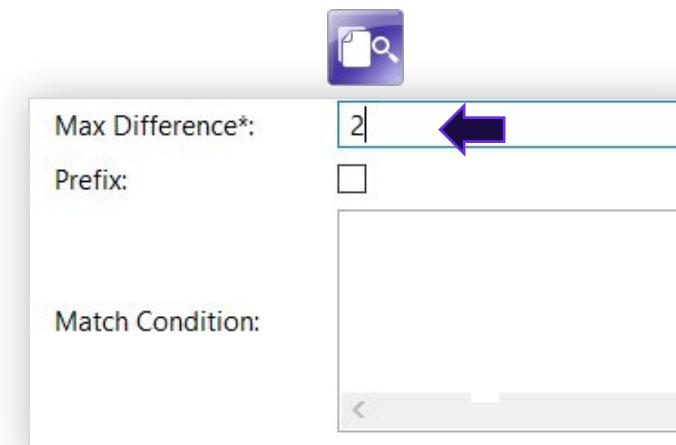
Additional Columns:

Name	Expression	Comment
1 CD_ITEM_LSID		
2 DS_ITEM		
*		

- **Duplicities** (first, accept, omit)
- **Matching Value** (Do upper case, squeeze whitespaces)
- **Approximative Index**
 - Necessary for approximative lookup (only for string keys)
- **Best Distance Index**
 - A more precise lookup stored in the index
- **Bidirect Approximative Index**
 - Improves the performance of approximative lookup
- **Compressed**
 - Compresses rows to improve performance; only useful with very wide datasets
- **User Metadata**
 - User-specific metadata will be included in the output
- **Use the **lkp** folder for storing your lookup files**

Lookups: Approximative Index

- Input strings are matched with lookup keys which have up to „**Max Difference**“ number of differences
- **What is considered a difference:**
 - One different character at any position within the string
 - One missing character at any position within the string
 - One additional character at any position within the string



- Check the **Approximative** Index checkbox in the **Lookup Builder** step

- Set **Max Difference** in the **Lookup** step
- Maximum allowed number of differences between the strings
- No more than 2 is recommended



Lookup Step: General

The screenshot shows the Ataccama Data Integration interface. On the left is the 'Properties of Lookup (Best distance)' dialog, which contains a tree view with 'Lookup' selected. On the right is the 'Lookup' configuration window, specifically the 'General' tab.

Mapping results and their destinations (lookup.key): Points to the 'Key Lookup Value*' field containing 'src_city'.

(Optional) Expression that must be satisfied for this step to be executed: Points to the 'Selection Rules' tab.

The column value used for querying the **lookup**: Points to the 'Id:' field containing 'Asset'.

The **lookup file**: Points to the 'Table File Name*' field containing './data/ext/lkp/ca_city.lkp'.

Scorer used for scoring the processing conditions (null key, no record found, etc.): Points to the 'Scorer*' field containing 'composite element'.

(Optional) Set of rules for picking the best match (Levenshtein**(src_first_name, lookup.fname))**: Points to the 'Select Best Match (0)' item under the 'Lookup' node in the properties tree.



Lookup Step: General

Properties of Lookup (Approximate)

Lookup

General Selection Rules Columns Advanced

Select Best Match:

Expression	Locale
1 levenshtein(src_province,lookup.key)	
*	

Max Difference*: 2

Prefix:

Match Condition:

Comments [\(Hide\)](#)

To Top Up Down To Bottom Fill Columns Debug... Edit... Template...

OK Cancel Apply

Rules to order the results and pick the best match (if multiple matches were found)

Maximum number of differences in approximative search string

Additional condition that has to be met in order for the record to be considered a match

?



Lookup Step: Columns

Name	Expression	Comment
1 query_lookupKey	query.lookupKey	
2 query_matchingKey	query.matchingKey	
3 query_records	query.records	
4 query_difference	query.difference	
5 query_relativeDifference	query.relativeDifference	
6 lookup_key	lookup.key	
*		

Besides the query key, there are many additional attributes which can be used or written to the output, such as:

- **Lookup.key**
- **Query.lookupKey**
- **Query.matchingKey**
- **Query.records**
- **Query.difference**
- **Query.relativeDifference**

The desired, initially specified key from the “Key Lookup Value” parameter.
The lookup key used in a Best Distance Index lookup.
The matching key used in a Best Distance Index lookup.
The number of records matched with a particular row.
The number of differences between the match and the matching key (if they are allowed).
If there were multiple matches, this shows how different each match was on average.



Apply Replacements Step

Scorer for the processing conditions (replacement applied, more replacements, etc.)

When disabled, the algorithm replaces **substrings in the input**. When enabled, replacement only occurs if the **entire input value** is present in the lookup file.

The **lookup file** with the data

Source and target columns

Properties of Apply Replacements (Apply Replacements)

Apply Replacements

General Advanced

General

Id: Apply Replacements

Other

Ignored Separators:

Only Full Replacement:

Preserve Unsupported Chars:

Replacements File Name*: ../data/ext/lkp/ca_street_apply_replacements.lkp

Scorer*: composite element

Tokenizer*: composite element

Column Bindings

In*: src_street

Out*: pur_street

Comments [\(Hide\)](#)

OK Cancel Apply



Strip Titles Step

Properties of Strip Titles (Remove titles)

Strip Titles

General Advanced

General

Id: Remove titles

Other

Ignored Separators: .

Min Word Count*: 2

Preserve Unsupported Chars:

Scorer*: composite element

Separator: |

Title Lookup File Name*:/data/ext/lkp/titles.lkp [Browse...](#)

Tokenizer*: composite element

Column Bindings

In*: pur_full_name [Create...](#) [Debug...](#) [Edit...](#)

Out: pur_full_name [Create...](#)

Titles Out: std_title [Create...](#)

Comments [\(Hide\)](#)

OK Cancel Apply

Source and target columns

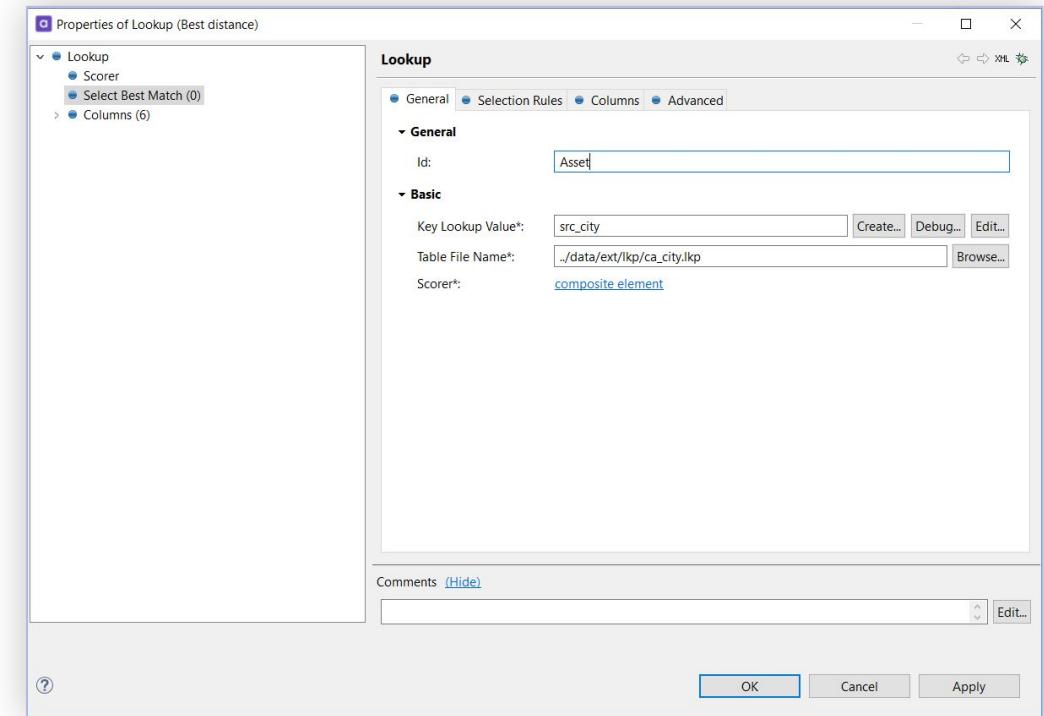
Column for storing removed values (titles)

Removal will not occur if the resulting string would have fewer words than the specified number

Lookup file containing the list of titles to be removed

Topic Highlights

- The **Lookup file** must be built in advance but can be reused as desired.
- **Lookup Builder** is used to create the Lookup file, a generic dictionary file.
- The **Lookup step** determines if the key exists in the lookup file.
- **Apply Replacements** will look for the input string in the dictionary, then perform replacements of all known instances of this string within the given text string.



Memory Refresher #8 Lookup & Dictionaries



What transformations occur with values when searching by lookup?

1. None.
2. You can specify the required transformations in the step's settings.
3. The same as described for the key during lookup building.

How many lines will lookup return if several values fit the key values?

- 1.** One.
- 2.** All.
- 3.** None.

Lab Exercise #8

Lookup & Dictionaries



ONE Desktop Core

v15.4.x



ataccama

