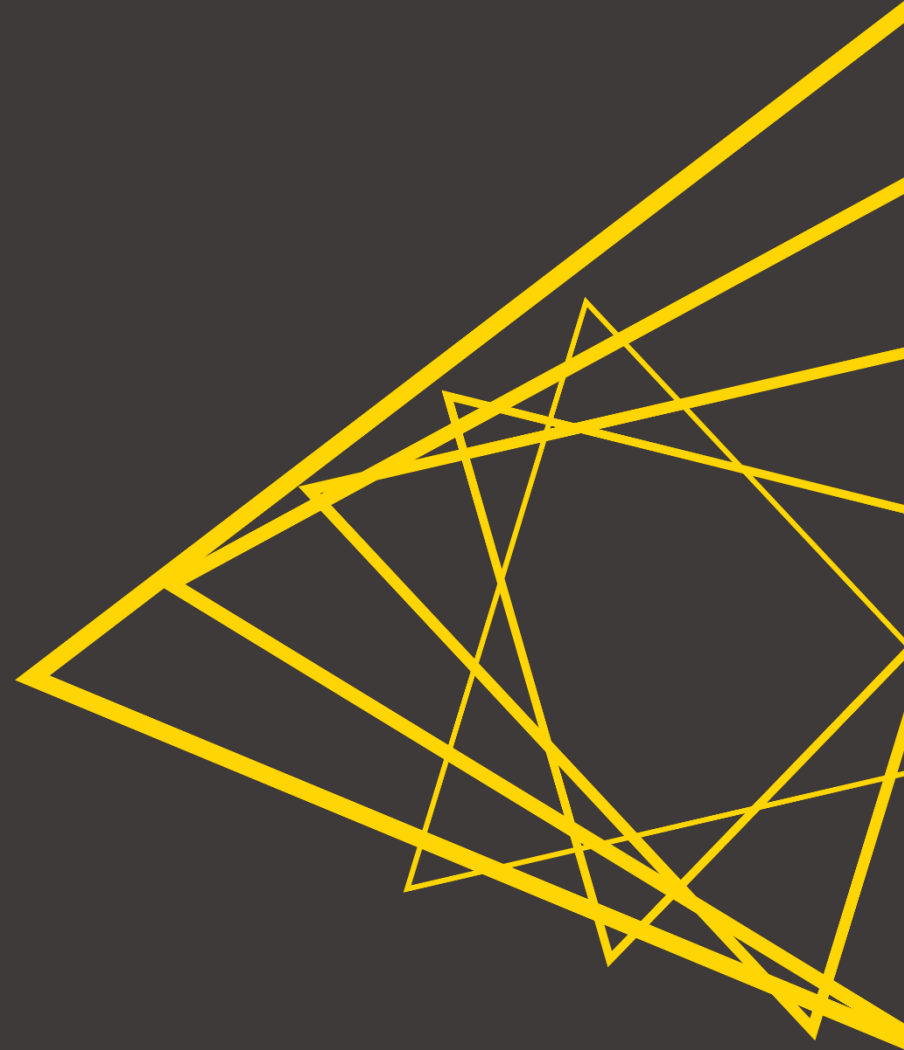# Boston Workshop

Monday, June 12 2023

**Dr. Alice Krebs**
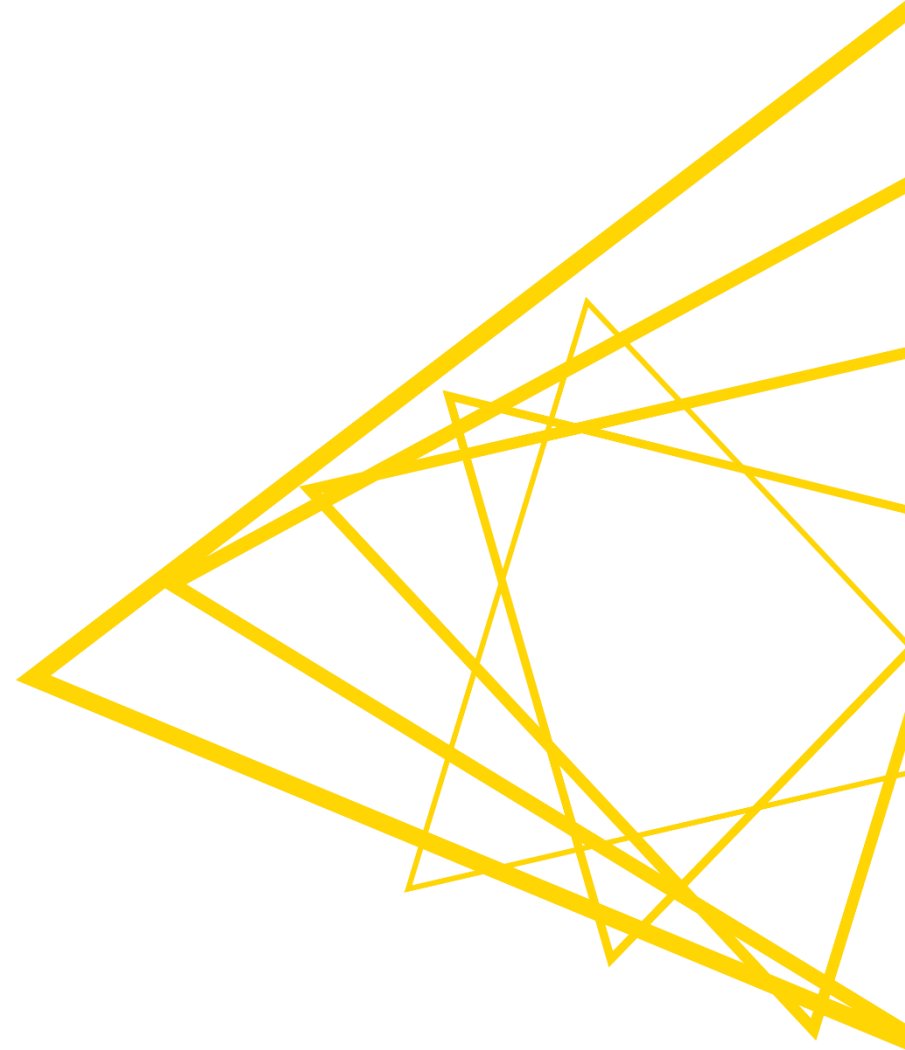
# Before we get started…

**… some questions for you:**

- What tools are you using for data analysis?
- What do you expect/hope to learn today?
- Did you manage to open the "00_Setup" workflow already?
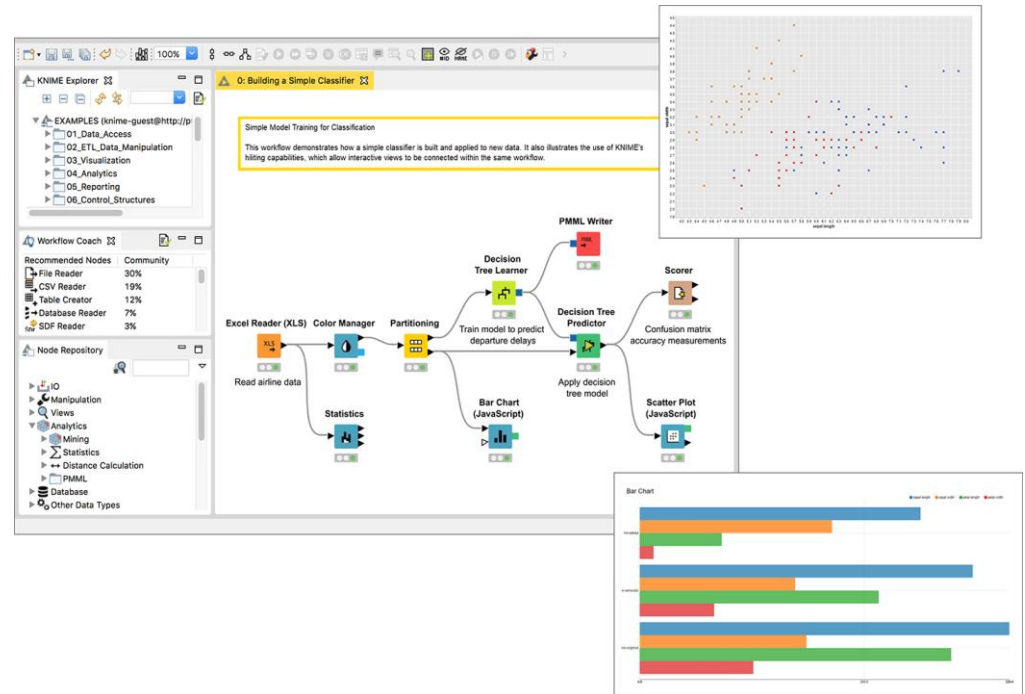
# Introduction to KNIME Analytics Platform

# What is KNIME Analytics Platform?

- A tool for data analysis, manipulation, visualization, and reporting

- Based on the graphical programming paradigm
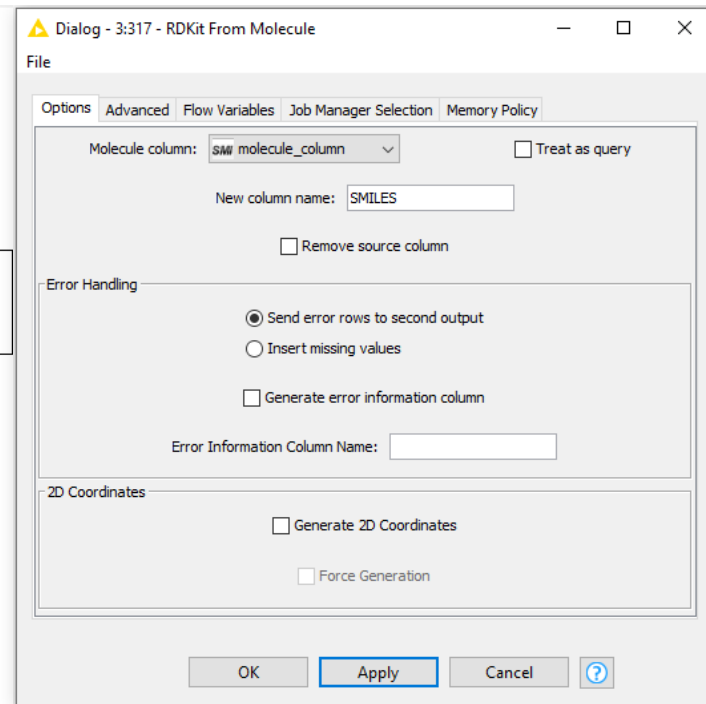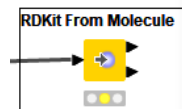
- Provides a diverse array of extensions:

  - Text Mining
  - Network Mining
  - Cheminformatics
  - Many integrations,
    such as Java, R, Python,
    Weka, Keras, Plotly, H2O, etc.

4

# Graphical programming paradigm?

- Lets users create programs by manipulating program elements graphically rather than by specifying them textually



```
mols = []
for smi in input_1_pandas[self.molecule_column]:
    mols.append(Chem.MolFromSmiles(smi))
```

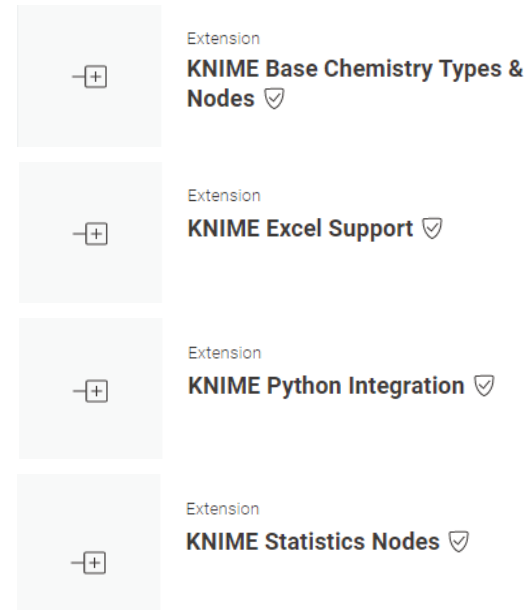rdkit.Chem.rdmolfiles.**MolFromSmiles**(*(AtomPairsParameters)SMILES, (SmilesParserParams)params*) → Mol :

Construct a molecule from a SMILES string.

# Extensions?

- KNIME AP comes by default only with basic functionality



- Extensions (and integrations) add specific functionality

- Not default, need to be added

Extension
**KNIME Base Chemistry Types & Nodes**

Extension
**KNIME Excel Support**

Extension
**KNIME Python Integration**

Extension
**KNIME Statistics Nodes**

# The KNIME Workbench

# The KNIME Workbench – most important windows

# Visual KNIME Workflows

**NODES** perform tasks on data



**Not Configured**

**Configured**

**Executed**

**Error**

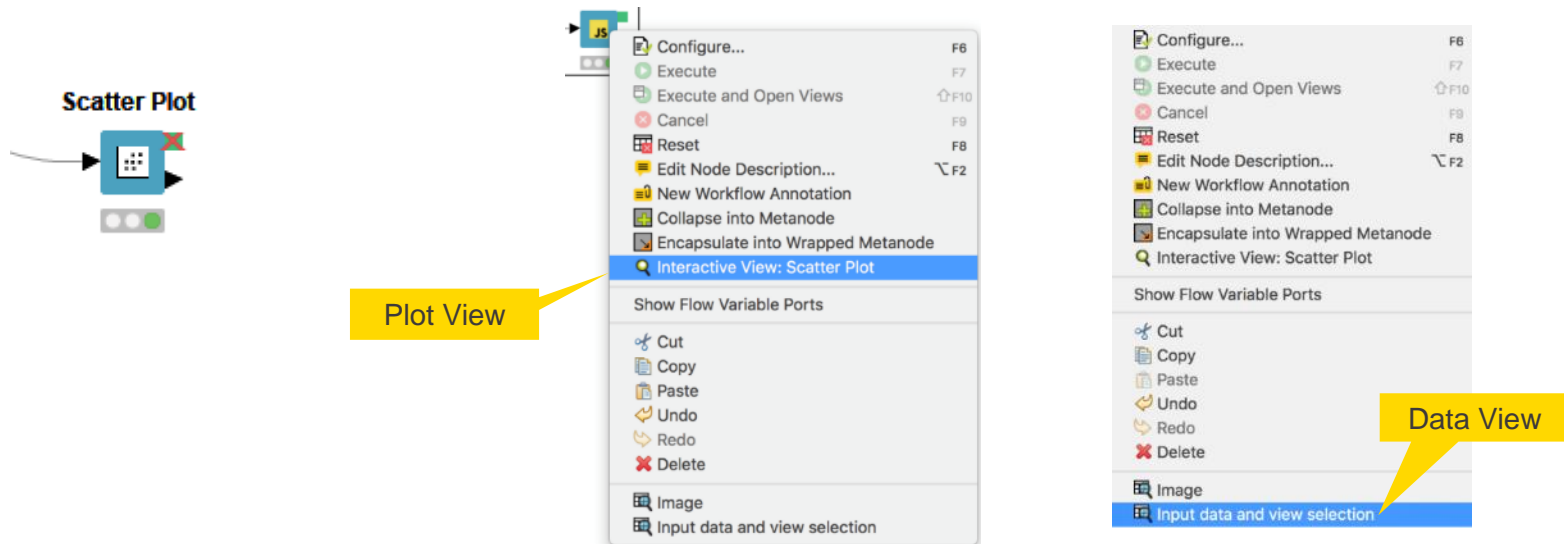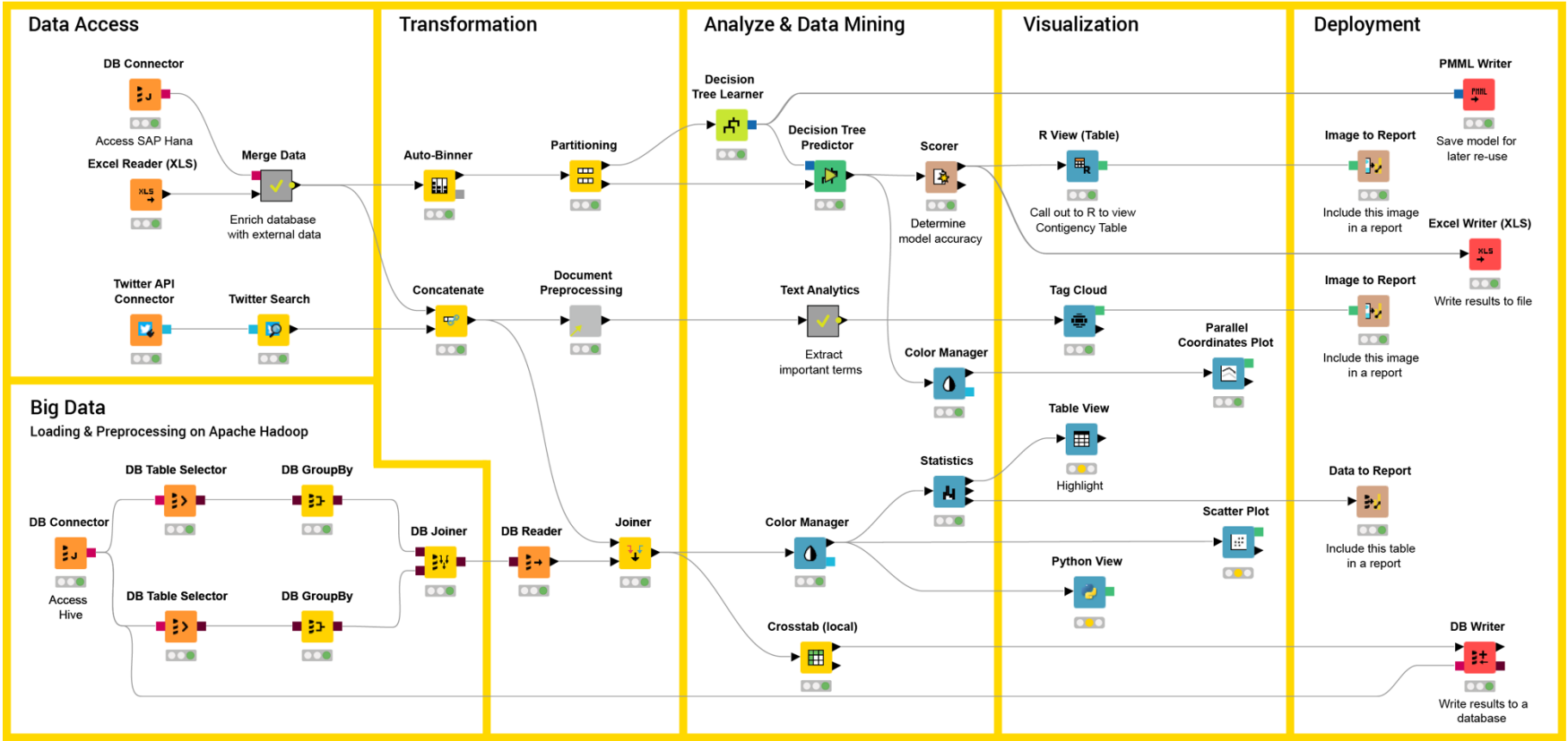Nodes are combined to create
**WORKFLOWS**

# Node Outputs and Views

- Right-click executed node
- Select View option in context menu

OR

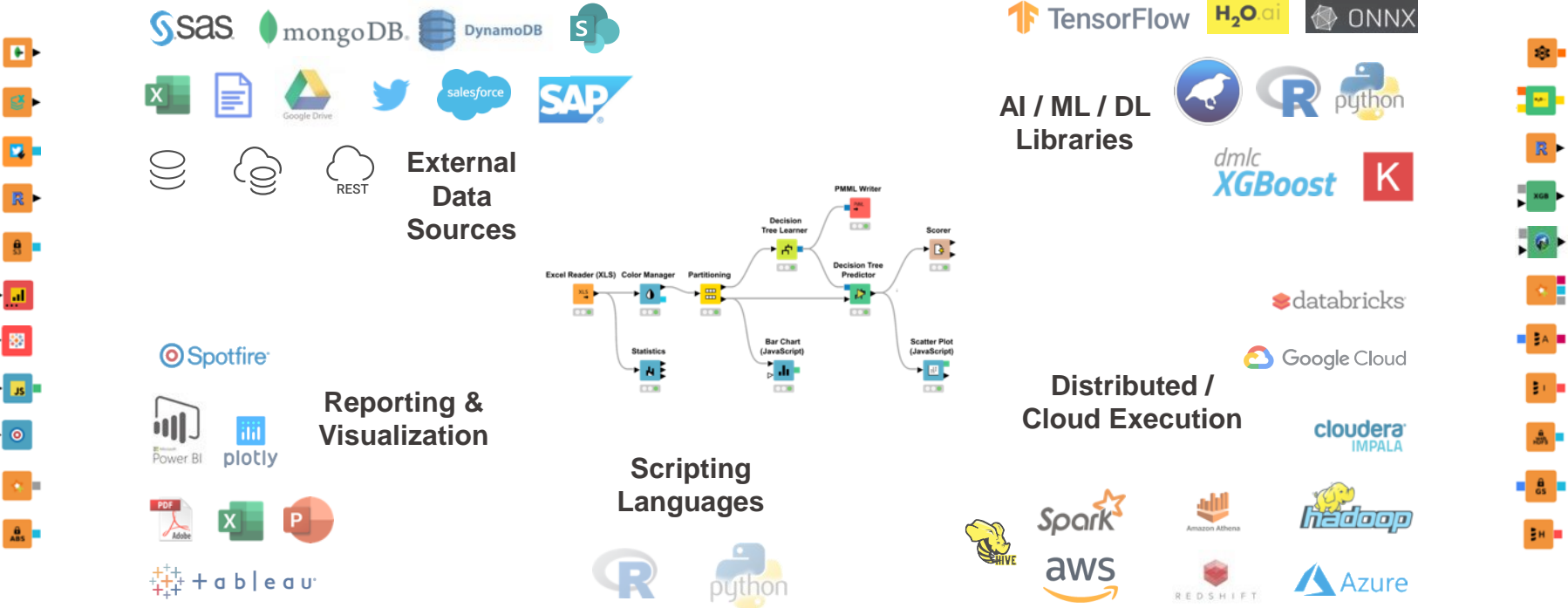- Select output port (last item) to inspect execution results



**Scatter Plot**

Plot View

Data View

# 4000+ Nodes for all Steps of End-To-End Data Science

# Mix & Match Data Sources, Technologies, and Execution



**External Data Sources**

**AI / ML / DL Libraries**

**Reporting & Visualization**

**Scripting Languages**

**Distributed / Cloud Execution**

KNIME — Open for Innovation

# Data Visualization
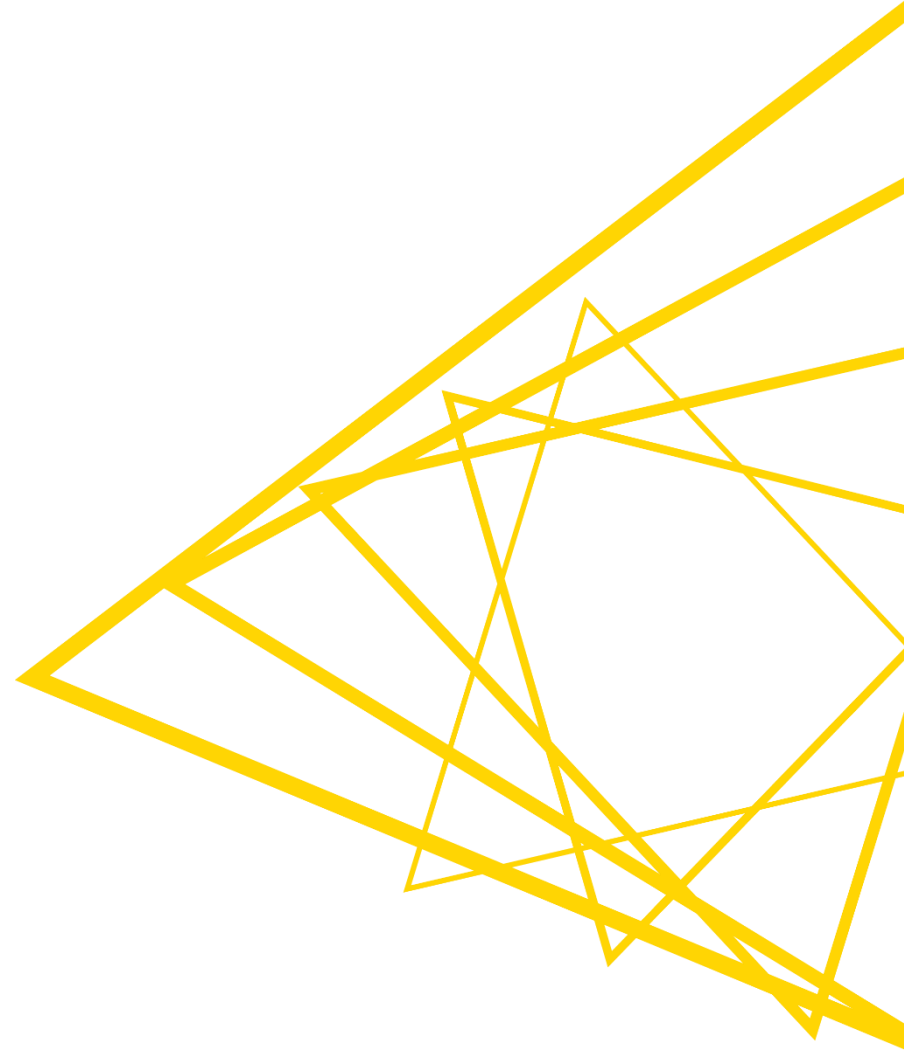
# Why KNIME?

- Self-documenting
  - Every workflow is a track of what you did to your data

- Reproducible
  - Once configured, it will run the same way every time and have the same results

- Data is available at every manipulation step
  - Trace and control how the data table is changing throughout the workflow

- Interactive data visualization

- Chemistry capabilities

# Setup KNIME AP for today

# Set up KNIME Analytics Platform

- Download workflows from the KNIME Community Hub
  - https://kni.me/s/5Goi-yKrwMLzAU4o

# Import the workflow group
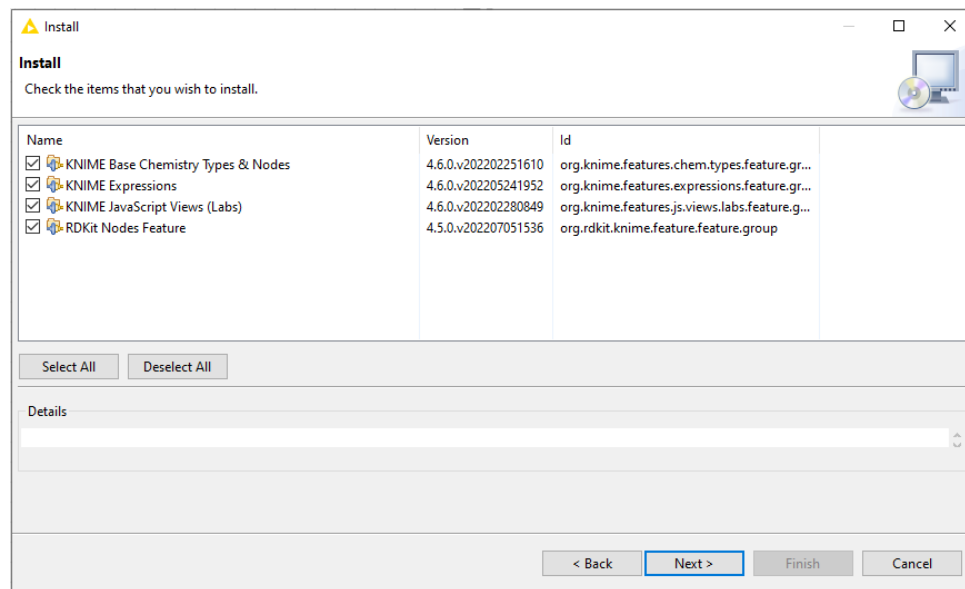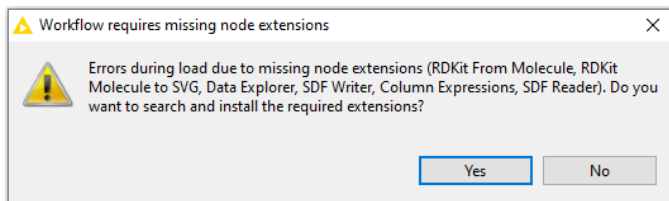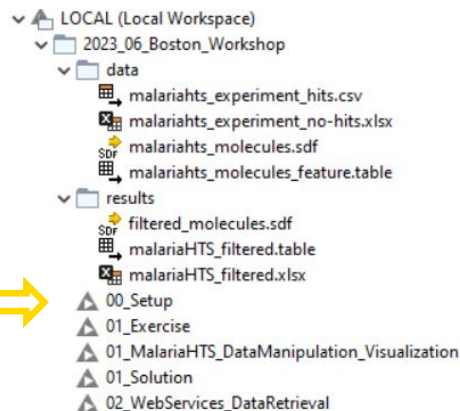
- Navigate to your download folder and choose the "2023_06_Boston_Workshop.knar" file
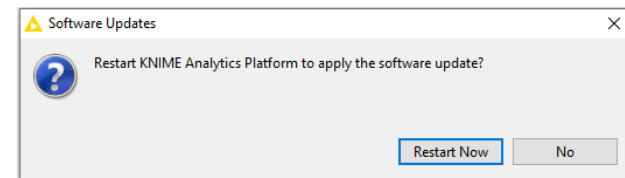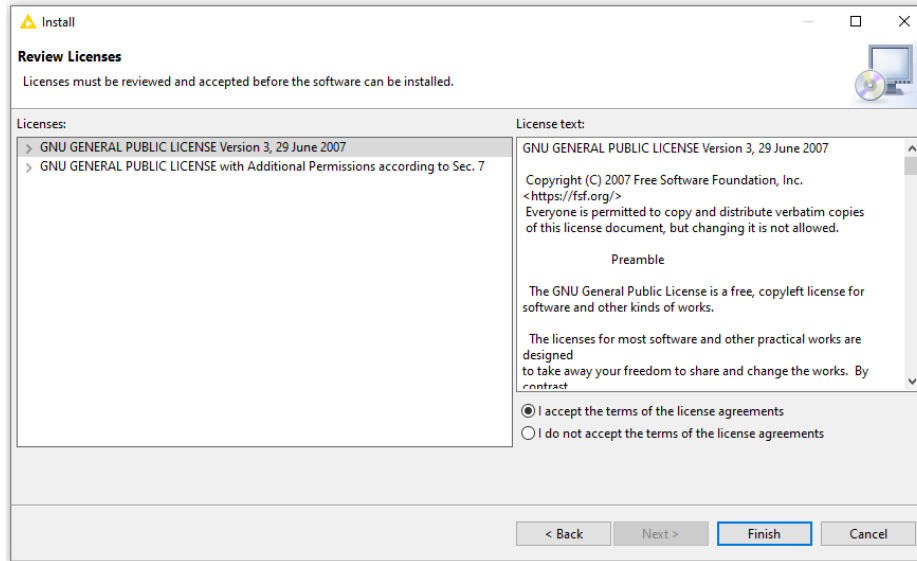
.knwf → *single workflow*
.knar → *group of workflows*

# Getting started

- Open the "00_Setup" workflow to install all the extensions you need today
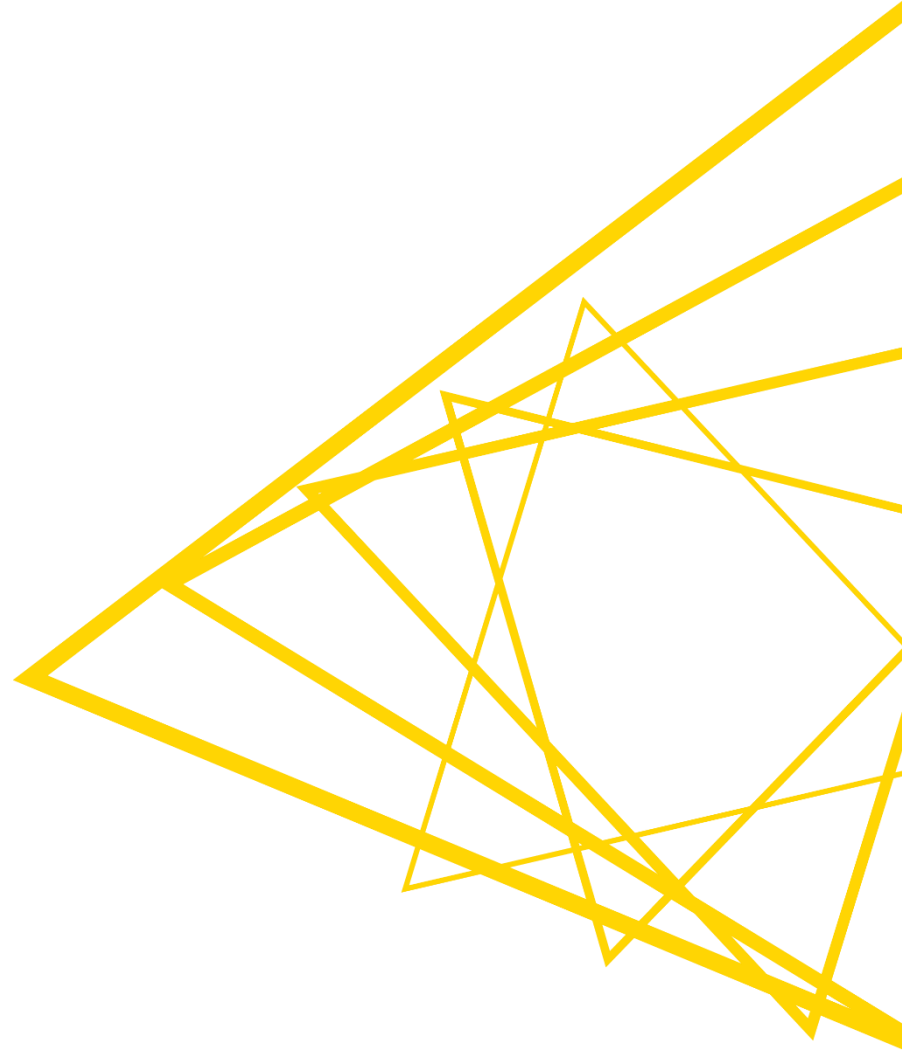
# Getting started
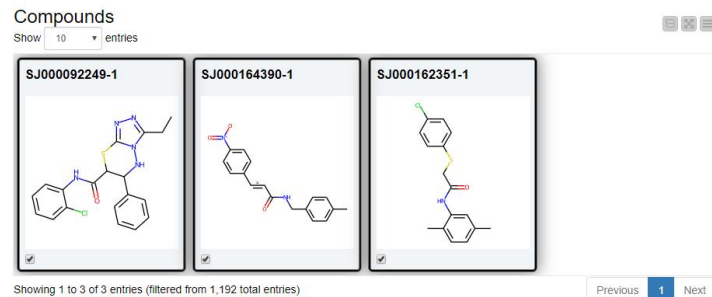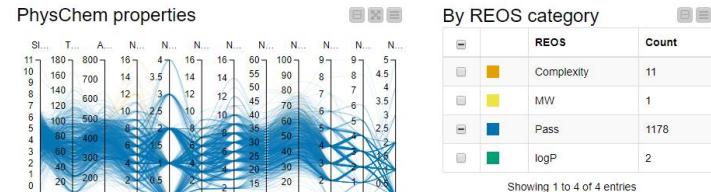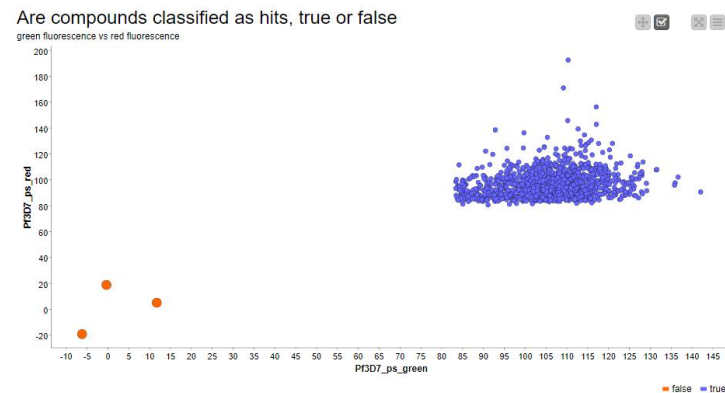
- Open the "00_Setup" workflow to install all the extensions you need today

# Ready to go!

# Today's workflows I

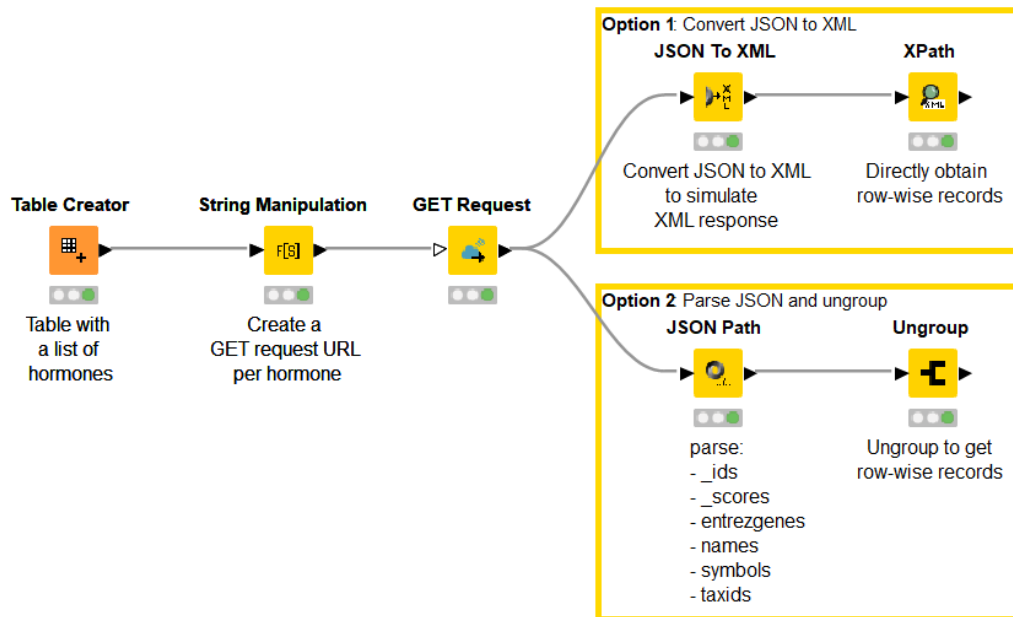- 01_MalariaHTS_DataManipulation_Visualization

# Today's workflows II
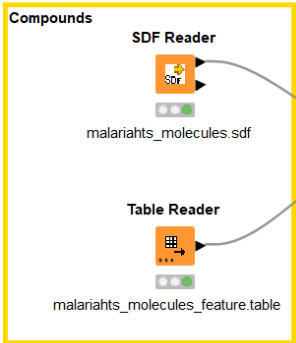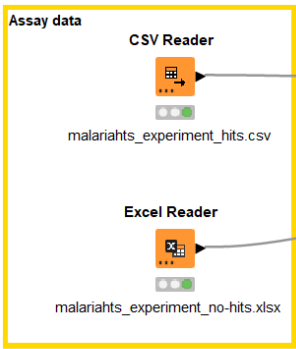
- 02_WebServices_DataRetrieval



In this workflow snippet, we will use the REST service provided by MyGene.info to obtain a list of human genes related to specific hormones. Then, we parse the JSON response into a table that is easy to read.

Sometimes, responses come in XML format. We have also included a way to parse XML responses by first converting the JSON response directly to XML.

**Option 1**: Convert JSON to XML

**JSON To XML** — Convert JSON to XML to simulate XML response

**XPath** — Directly obtain row-wise records

**Table Creator** — Table with a list of hormones

**String Manipulation** — Create a GET request URL per hormone

**GET Request**

**Option 2** Parse JSON and ungroup

**JSON Path** — parse:
- _ids
- _scores
- entrezgenes
- names
- symbols
- taxids

**Ungroup** — Ungroup to get row-wise records

# 01_MalariaHTS_DataManipulation_Visualization

- **Step 1:** import data from 4 different sources

# 01_MalariaHTS_DataManipulation_Visualization

- **Step 1:** import data from 4 different sources

# 01_MalariaHTS_DataManipulation_Visualization

- **Step 2:** data cleaning, merging datasets, manipulate data, create classification

Open for Innovation
KNIME

# 01_MalariaHTS_DataManipulation_Visualization

- **Step 3:** visualize and interactively select data
  - build a component to combine different plots in one view

**Interactive data visualization**

INTERACTION REQUIRED

# 01_MalariaHTS_DataManipulation_Visualization

- **Step 4:** Export the data
  - SD file
  - KNIME-native table (only re-usable in KNIME)
  - Excel

**SDF Writer**

**Table Writer**

KNIME-native format

**Excel Writer**

Open for Innovation

KNIME

# Content, theory, background on…

- Import data
  - Local file system
  - Workflow relative path
- Merging data sets
  - Concatenate
  - Join
- Exclude data
  - Filter
  - Splitter
- Handle missing values and duplicates

- Change data
  - Change assigned data type
  - Capitalize the letters in the sample ID
  - Rename columns
  - Create classification of REOS rules
- Visualize data in components

KNIME
Open for Innovation

# Import data

- Different Reader nodes for different file formats

- "Table" is a KNIME native format

- Drag and drop the files from the data folder in the KNIME Explorer

**Excel Reader**

malariahts_experiment_no-hits.xlsx

**CSV Reader**

malariahts_experiment_hits.csv

**SDF Reader**

malariahts_molecules.sdf

**Table Reader**

malariahts_molecules_feature.table

# Common Settings: File Path

- A path consists of three parts:
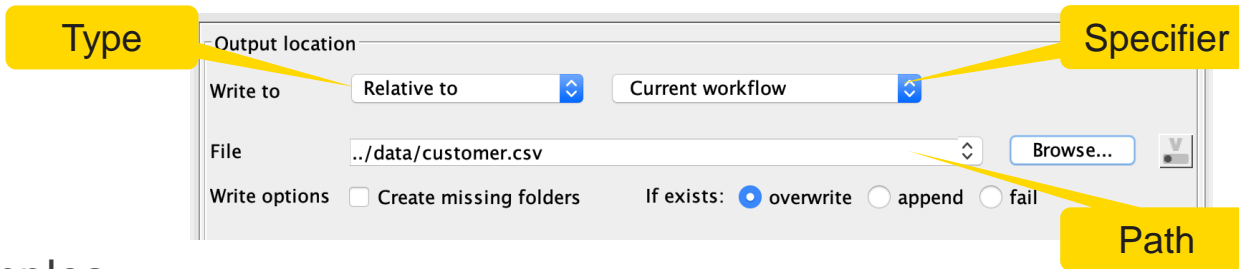  - **Type**: Specifies the file system type e.g. local, relative, mountpoint, custome_url or connected.
  - **Specifier**: Optional string with additional file system specific information e.g. relative to which location (knime.workflow)
  - **Path**: Specifies the location within the file system



- Examples:
  - (LOCAL, , C:\Users\username\Desktop)
  - (RELATIVE, knime.workflow, file1.csv)
  - (MOUNTPOINT, MOUNTPOINT_NAME, /path/to/file1.csv)
  - (CONNECTED, amazon-s3:eu-west-1, /mybucket/file1.csv)

# Common Settings: 4 Default File Systems

- **Local File System**



- **Relative to …**



- **Mountpoint**



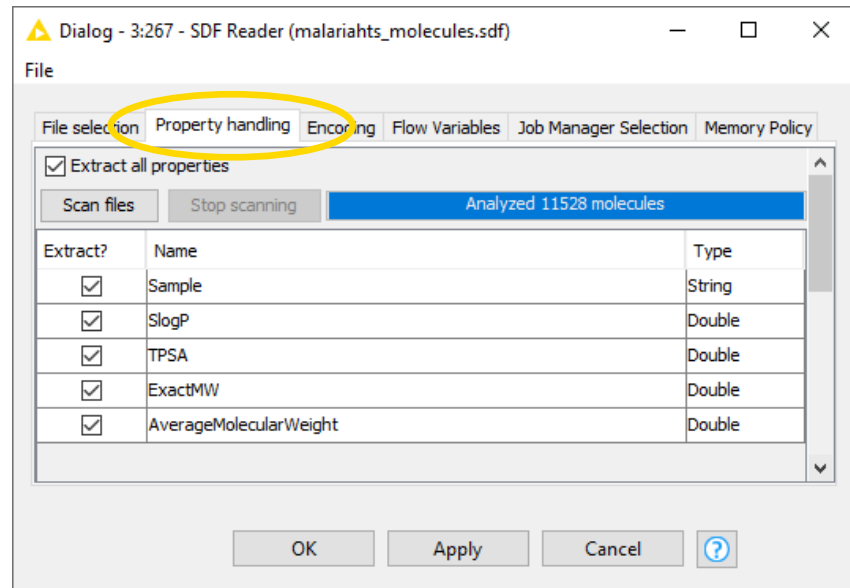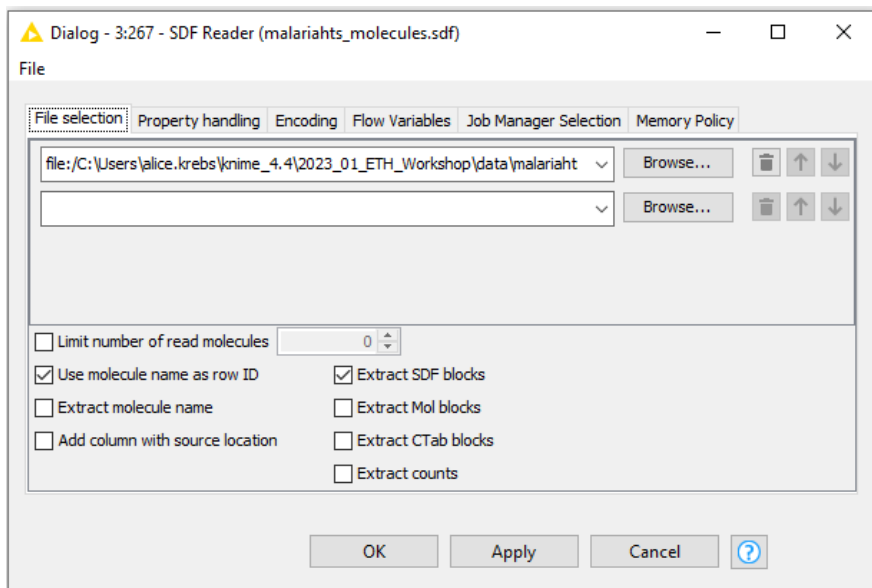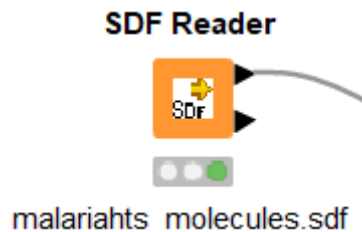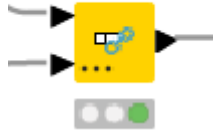- **Custom URL**

# Special case – SDF reader node

- No new file handling
- Extract properties to get all data

# Concatenate vs. Join

**Concatenate**



**Joiner**



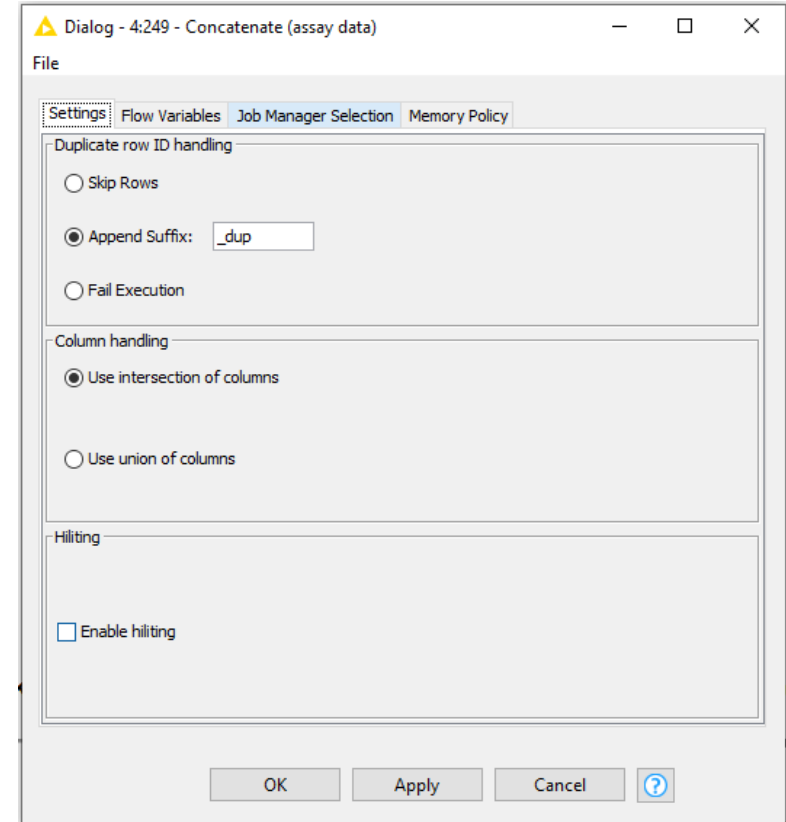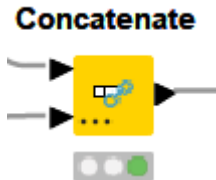- Simply links to tables, appends them underneath each other, like in a chain

| Data table 1 |
|:---:|

$$+$$

| Data table 2 |
|:---:|

$$=$$

| concatenated |
|:---:|
| data table |

- Combines two tables row wise

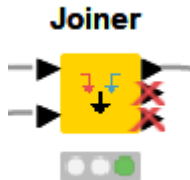| Data table 1 | + | Data table 2 | = | joint data table |
|:---:|:---:|:---:|:---:|:---:|

# Concatenate node

- **Intersection**
  - use only the columns that appear in both input tables

- **Union**
  - use all columns available in the input tables

# Joiner node

# Joining Columns of Data – Inner Join

Left Table

| molregno | chembl_id | SMILES |
|----------|-----------|--------|
| 22 | CHEMBL1794855 | CCCN(CCC) |
| 24 | CHEMBL278751 | CCN(C) |
| 15 | CHEMBL103772 | CCCN1CC |
| 10 | CHEMBL328107 | C1CN(CCN1) |

Right Table

| molregno | Ki_value | Ki_relation | Ki_unit |
|----------|----------|-------------|---------|
| 17 | 76.0 | = | nM |
| 65 | 6.56 | = | nM |
| 35 | 100 | > | nM |
| 15 | 8 | = | nM |
| 10 | 95.8 | = | nM |

Inner Join

| molregno | chembl_id | SMILES | Ki_value | Ki_relation | Ki_unit |
|----------|-----------|--------|----------|-------------|---------|
| 15 | CHEMBL103772 | CCCN1CC | 8 | = | nM |
| 10 | CHEMBL328107 | C1CN(CCN1) | 95.8 | = | nM |

Open for Innovation
KNIME

# Joining Columns of Data – Left Outer Join

Left Table

| molregno | chembl_id | SMILES |
|----------|-----------|--------|
| 22 | CHEMBL1794855 | CCCN(CCC) |
| 24 | CHEMBL278751 | CCN(C) |
| 15 | CHEMBL103772 | CCCN1CC |
| 10 | CHEMBL328107 | C1CN(CCN1) |

Right Table

| molregno | Ki_value | Ki_relation | Ki_unit |
|----------|----------|-------------|---------|
| 17 | 76.0 | = | nM |
| 65 | 6.56 | = | nM |
| 35 | 100 | > | nM |
| 15 | 8 | = | nM |
| 10 | 95.8 | = | nM |

Left Outer Join

| molregno | chembl_id | SMILES | Ki_value | Ki_relation | Ki_unit |
|----------|-----------|--------|----------|-------------|---------|
| 22 | CHEMBL1794855 | CCCN(CCC) | ? | ? | ? |
| 24 | CHEMBL278751 | CCN(C) | ? | ? | ? |
| 15 | CHEMBL103772 | CCCN1CC | 8 | = | nM |
| 10 | CHEMBL328107 | C1CN(CCN1) | 95.8 | = | nM |

KNIME
Open for Innovation

# Joining Columns of Data – Right Outer Join

Left Table

| molregno | chembl_id | SMILES |
|---|---|---|
| 22 | CHEMBL1794855 | CCCN(CCC) |
| 24 | CHEMBL278751 | CCN(C) |
| 15 | CHEMBL103772 | CCCN1CC |
| 10 | CHEMBL328107 | C1CN(CCN1) |

Right Table

| molregno | Ki_value | Ki_relation | Ki_unit |
|---|---|---|---|
| 17 | 76.0 | = | nM |
| 65 | 6.56 | = | nM |
| 35 | 100 | > | nM |
| 15 | 8 | = | nM |
| 10 | 95.8 | = | nM |

Right Outer Join

| molregno | chembl_id | SMILES | Ki_value | Ki_relation | Ki_unit |
|---|---|---|---|---|---|
| 17 | ? | ? | 76.0 | = | nM |
| 65 | ? | ? | 6.56 | = | nM |
| 35 | ? | ? | 100 | > | nM |
| 15 | CHEMBL103772 | CCCN1CC | 8 | = | nM |
| 10 | CHEMBL328107 | C1CN(CCN1) | 95.8 | = | nM |

KNIME
Open for Innovation

# Joining Columns of Data – Full Outer Join

**Left Table**

| molregno | chembl_id | SMILES |
|----------|-----------|--------|
| 22 | CHEMBL1794855 | CCCN(CCC) |
| 24 | CHEMBL278751 | CCN(C) |
| 15 | CHEMBL103772 | CCCN1CC |
| 10 | CHEMBL328107 | C1CN(CCN1) |

**Right Table**

| molregno | Ki_value | Ki_relation | Ki_unit |
|----------|----------|-------------|---------|
| 17 | 76.0 | = | nM |
| 65 | 6.56 | = | nM |
| 35 | 100 | > | nM |
| 15 | 8 | = | nM |
| 10 | 95.8 | = | nM |

**Full Outer Join**

Values missing in the left table

Values missing in the right table

| molregno | chembl_id | SMILES | Ki_value | Ki_relation | Ki_unit |
|----------|-----------|--------|----------|-------------|---------|
| 17 | ? | ? | 76.0 | = | nM |
| 65 | ? | ? | 6.56 | = | nM |
| 35 | ? | ? | 100 | > | nM |
| 15 | CHEMBL103772 | CCCN1CC | 8 | = | nM |
| 10 | CHEMBL328107 | C1CN(CCN1) | 95.8 | = | nM |
| 22 | CHEMBL1794855 | CCCN(CCC) | ? | ? | ? |
| 24 | CHEMBL278751 | CCN(C) | ? | ? | ? |

# Filter vs. Splitter



- User-defined criteria
- Only one output port
- Excluded rows or columns are not available for downstream processing

- User-defined criteria
- Two output ports
- Splits a dataset into two
- The 'excluded' data is available at the lower output port for downstream processing

# Missing value handling

- First tab: define table-wide action depending on data type
- Second tab: define action for individual columns

# Exclude duplicates

- identifies duplicate rows

- duplicate rows have identical values in certain columns

- user defines the column(s) for duplicate detection

# Change assigned data type

- modify existing columns using expressions

# Capitalize the letters in the sample ID

- Manipulates strings
- Many different functions available
- Can also be used for type converting

# Rename column

- Transformations of input columns
  - Renaming
  - Filtering
  - Re-ordering
  - Type changing
- Check data manipulation in the preview



**Table Manipulator**

# There isn't just one way of doing it…

**Column Rename**

- Rename columns

**Column Resorter**

- Change the column order

**Row Filter**   **Row Splitter**

- Exclude missing values

**String To Number**   **Double To Int**

- Convert data types

# Create classification REOS rules



*Define Condition* => *Outcome if condition is met*

*If not do this*

**IF… THEN… ELSE**

# GroupBy node

- Aggregate rows to summarize data
- Different aggregation methods available

| Product ID | REOS | AMW |
|------------|------------|--------|
| P001 | Pass | 462.55 |
| P002 | Pass | 397.92 |
| P003 | Complexity | 431.28 |
| P004 | Pass | 431.28 |
| P005 | MW | 371.46 |
| P006 | logP | 389.20 |

| Group | Count*(AMW) |
|------------|-------------|
| Pass | 3 |
| Complexity | 1 |
| MW | 1 |
| logP | 1 |

# GroupBy node

- Aggregate rows to summarize data
- Different aggregation methods available

| Product ID | REOS | AMW |
|------------|------------|--------|
| P001 | Pass | 462.55 |
| P002 | Pass | 397.92 |
| P003 | Complexity | 431.28 |
| P004 | Pass | 431.28 |
| P005 | MW | 371.46 |
| P006 | logP | 389.20 |

| Group | Mean(AMW) |
|------------|-----------|
| Pass | 430.58 |
| Complexity | 431.28 |
| MW | 371.46 |
| logP | 389.20 |

KNIME
Open for Innovation

# Groupby node configuration

- Define the column to group on in the first tab

**GroupBy**

# Groupby configuration

- Define the aggregation method in the second tab

- Information about the methods is provided in the Description tab

# Visualization in components

- Mark the view nodes you want to combine

- Right-click > create component

- Use the Javascript-based nodes

# Components – Combined Views

- Multiple JavaScript View nodes can be combined in components

- Selections are transmitted to all other views

- Also for use on the KNIME WebPortal (commercial product)

# JS-based nodes: Selection and Filter Events



- Selection and Filter tab in many Javascript-based view nodes

# JS-based nodes: Selection and Filter Events



- Listens to the selection made in other view

# JS-based nodes: Selection and Filter Events



- Makes the other views listen to the selection made in this view

# Color Manager

- Colors by nominal or continuous values
- Syncs colors between views using the color model port and Color Appender node



Discrete colors for nominal values

Color range for numerical values

**Color Manager**

Hit classification, true or false
green fluorescence vs red fluorescence

# Configure Content and Views Layout

- Click layout button when inside Component to assign views to rows and columns

- Add views and rows via drag&drop

- Add columns using **+** buttons

That's it for workflow I…

# Workflow II – web service data retrieval

- Instead of retrieving data manually, I can do it in an automated way
- Especially useful for large amounts of data
- Many data sources offer the 'endpoint' to do so ($\rightarrow$ REST API, web service)

**MANUALLY**



user

User interacts with the browser

browser

Browser gets the data from the database in the back

laptop

database

# Workflow II – web service data retrieval

- Instead of retrieving data manually, I can do it in an automated way

- Especially useful for large amounts of data

- Many data sources offer the 'endpoint' to do so ($\rightarrow$ REST API, web service)

**AUTOMATICALLY**

user

browser

laptop

Let the laptop talk to the database directly to retrieve the data

database

Open for Innovation
KNIME

# Workflow II – web service data retrieval

- Examples of data sources with an API:
    - https://docs.mygene.info/en/latest/index.html
    - https://chembl.gitbook.io/chembl-interface-documentation/web-services/chembl-data-web-services
- Construct the according URL to retrieve data for according query

## Quick start

MyGene.info provides two simple web services: one for gene queries and the other for gene annotation retrieval. Both return results in JSON format.

### Gene query service

**URL**

```
http://mygene.info/v3/query
```

Steady part

**Examples**

```
http://mygene.info/v3/query?q=cdk2
http://mygene.info/v3/query?q=cdk2&species=human
http://mygene.info/v3/query?q=cdk?
http://mygene.info/v3/query?q=1L*
http://mygene.info/v3/query?q=entrezgene:1017
http://mygene.info/v3/query?q=ensemblgene:ENSG00000123374
http://mygene.info/v3/query?q=cdk2&fields=symbol,refseq
```

# Workflow II – web service data retrieval

- Join the steady part of the URL with your query using the String Manipulation
- The GET Request node 'talks' to the DB and retrieves the data
- The result is often not very human-friendly (JSON or XML data format)
- Parse the results to get a
      neat table

That's it for workflow II…

# KNIME courses

- **Self-paced courses**
  - Courses are organized by level L1 (basic) - L4 (specialized)
  - lessons with ~5 minutes videos, hands-on exercises, and knowledge-check questions

- **Instructor-lead online courses**
  - One week, 1h15 per day
  - Possibility to ask questions
  - Discount available for academics



https://www.knime.com/knime-courses

# KNIME book

- KNIME Beginners Luck



https://www.knime.com/knimepress/beginners-luck

# More resources

- ## KNIME Community Hub
  - Find example workflows you can adapt
  - Find out e.g. how others use a certain nodes



- ## KNIME Forum
  - Get and find answers to your questions



… and **Youtube**
https://www.youtube.com/knimetv

# Cheat Sheets

- Building a KNIME Workflow for Beginners

- Building Components for your Team or the KNIME Community

- Control and Orchestration with KNIME AP

- Data Wrangling with KNIME AP

- Connectors with KNIME AP

- Machine Learning with KNIME AP



Cheat Sheet: Data Wrangling with KNIME Analytics Platform

https://www.knime.com/cheat-sheets

**Thank You!**

**Questions? Please reach out**

alice.krebs@knime.com

zachary.durso@knime.com