# CSMDM21 - Data Analytics and Mining

# Introduction to the Data Science Platform KNIME

Module convenor
**Dr. Carmen Lam**
carmen.lam@reading.ac.uk
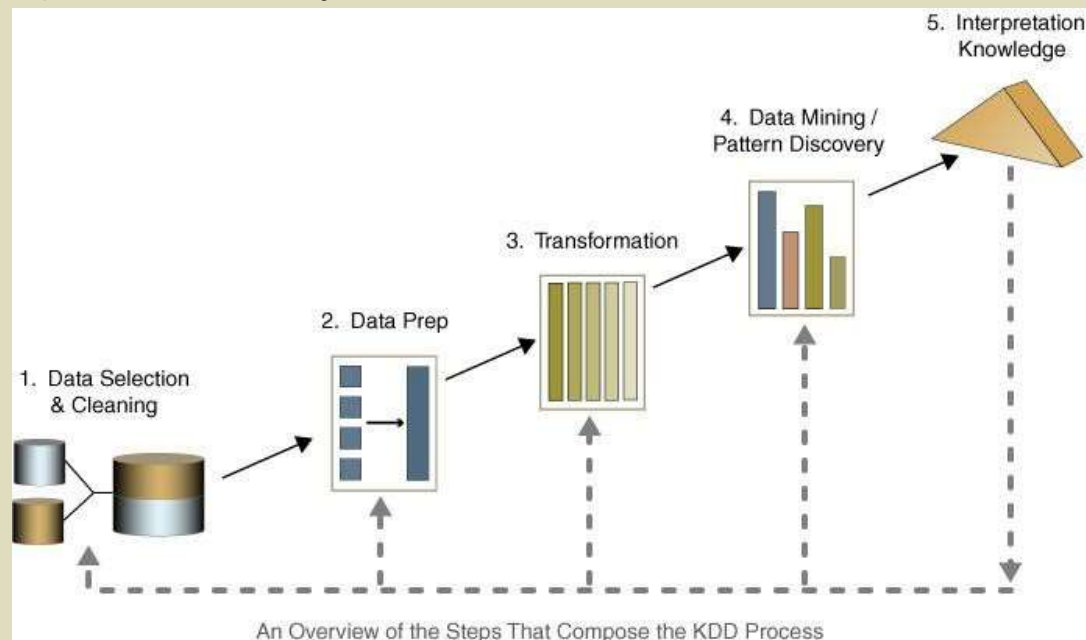Department of Computer Science

Lecture notes and videos created by
Prof. Giuseppe Di Fatta

# Development Environments for Data Science

❑ Data Science
❑ Knowledge Discovery Process (KDD)
❑ Data Mining and Machine Learning algorithms

- An open and integrated environment to facilitate the KDD process.
- A **workflow management system for data analytics and mining** that integrates methods for data management and information visualization with algorithms for prediction, pattern discovery, classification.



An Overview of the Steps That Compose the KDD Process
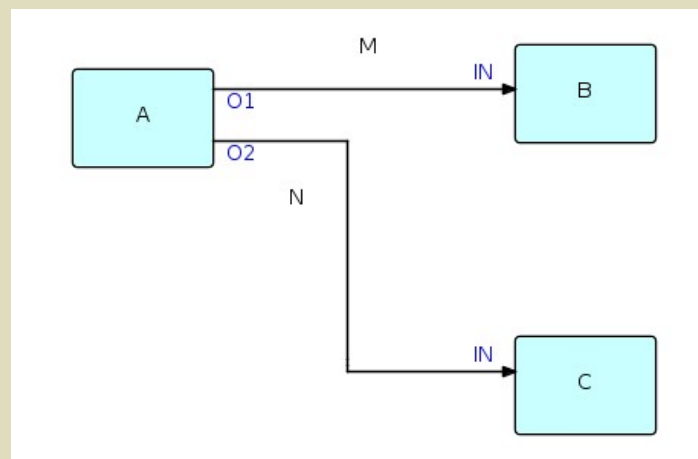
# Introducing KNIME



- Developed at the **ALTANA-Chair for Bioinformatics and Information Mining**, Department of Computer and Information Science, University of Konstanz (Germany)
- 1st release: ~140 MM effort in two years
- Under continuous evolution and extension
  - April 2006: 1st release, version 1.0.0
  - Current release: version 4.2.1 (July 2020)
    - "KNIME Analytics Platform" – An open source software for creating data science applications
      - extensions and integrations
      - KNIME SDK
    - KNIME server (proprietary, enterprise software for teams)
  - Available for Windows, Linux and Mac OS X
  - An open source license (GPL) allows KNIME to be downloaded, distributed, and used freely.

M. Berthold, N. Cebron, F. Dill, G. Di Fatta, T. Gabriel, F. Georg, T. Meinl, P. Ohl, C. Sieb, B. Wiswedel, "KNIME: the Konstanz Information Miner", Proc. of Workshop on Multi-Agent Systems and Simulation (MAS&S), 4th Annual Industrial Simulation Conference (ISC), Palermo, Italy, June 5-7, 2006, pp.58-61.
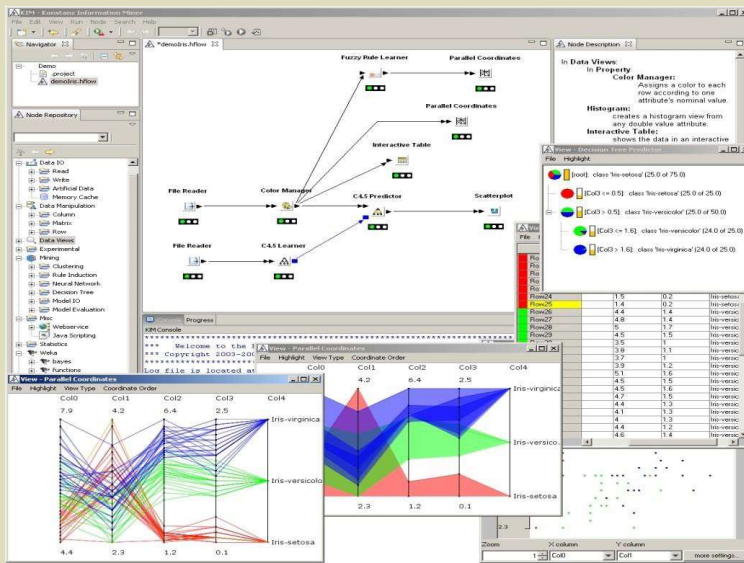
# Flow-based Programming

- Flow-based Programming (FBP) is a programming paradigm that defines applications as networks of "black box" processes, which exchange data across predefined connections by message passing, where the connections are specified externally to the processes.

- These black box processes can be reconnected endlessly to form different applications without having to be changed internally. FBP is thus naturally component-oriented.



- J. Paul Morrison, Flow-Based Programming, 2nd Edition: A New Approach to Application Development, CreateSpace, 2010

# KNIME

## Interactive Data Exploration



**Features:**

- Modular Data Pipeline Environment
- Large collection of Data Mining techniques
- Data and Model Visualizations
- Interactive Views on Data and Models
- Java Code Base as Open Source Project
- Seamless Integration: R Library, Weka, etc.
- Based on the Eclipse Plug-in technology

Easy extendibility
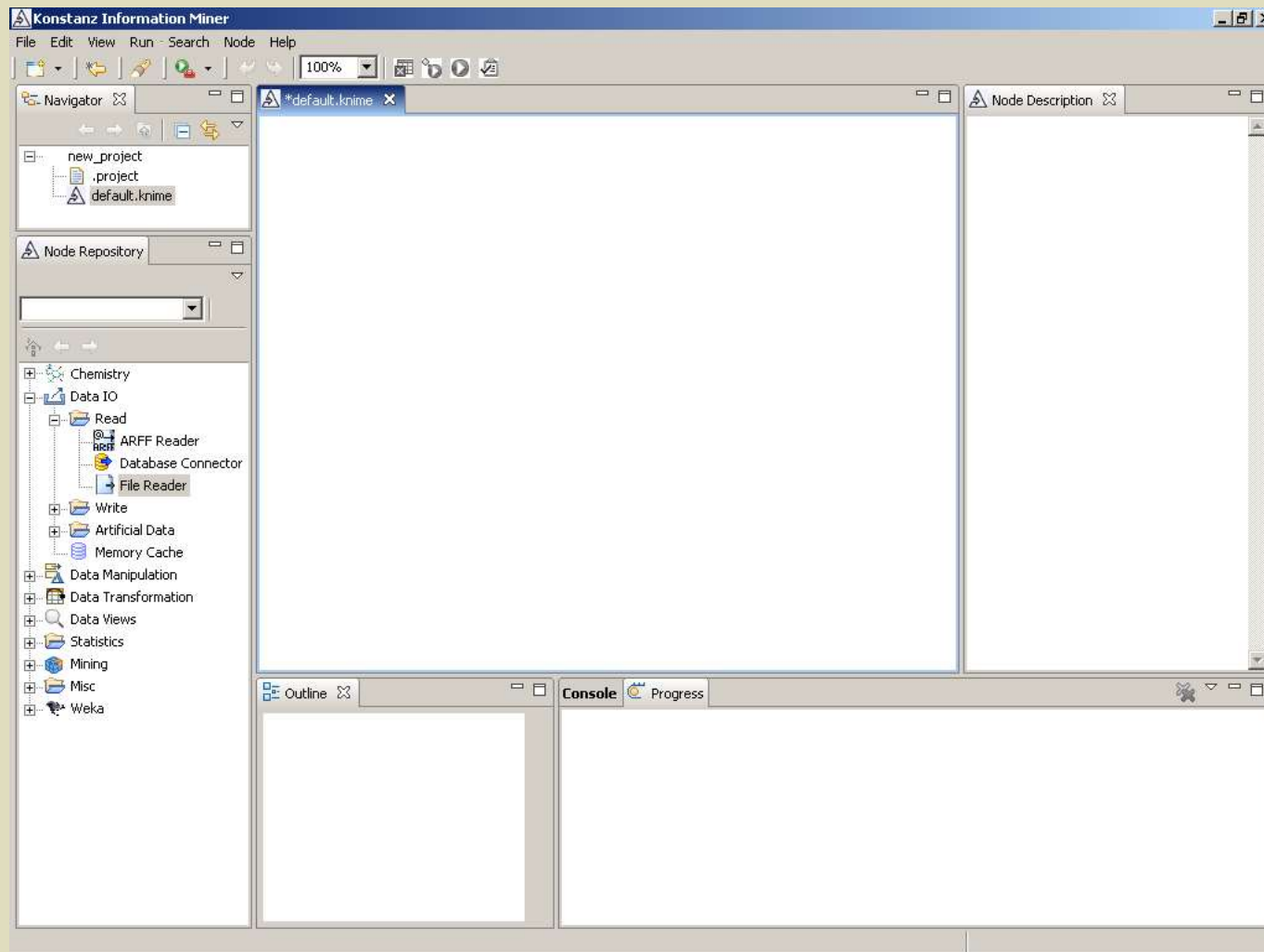    New nodes via open API and integrated wizard
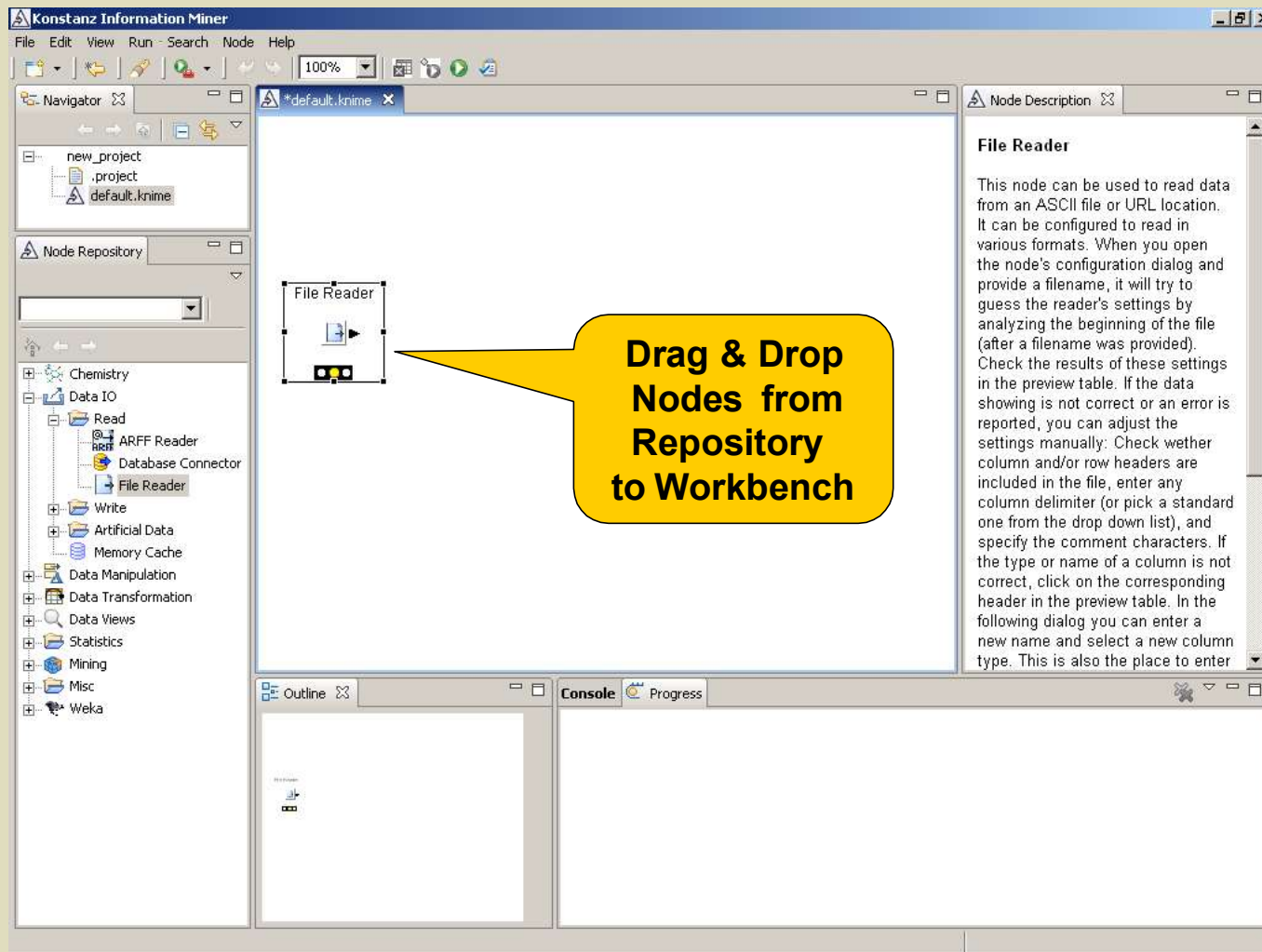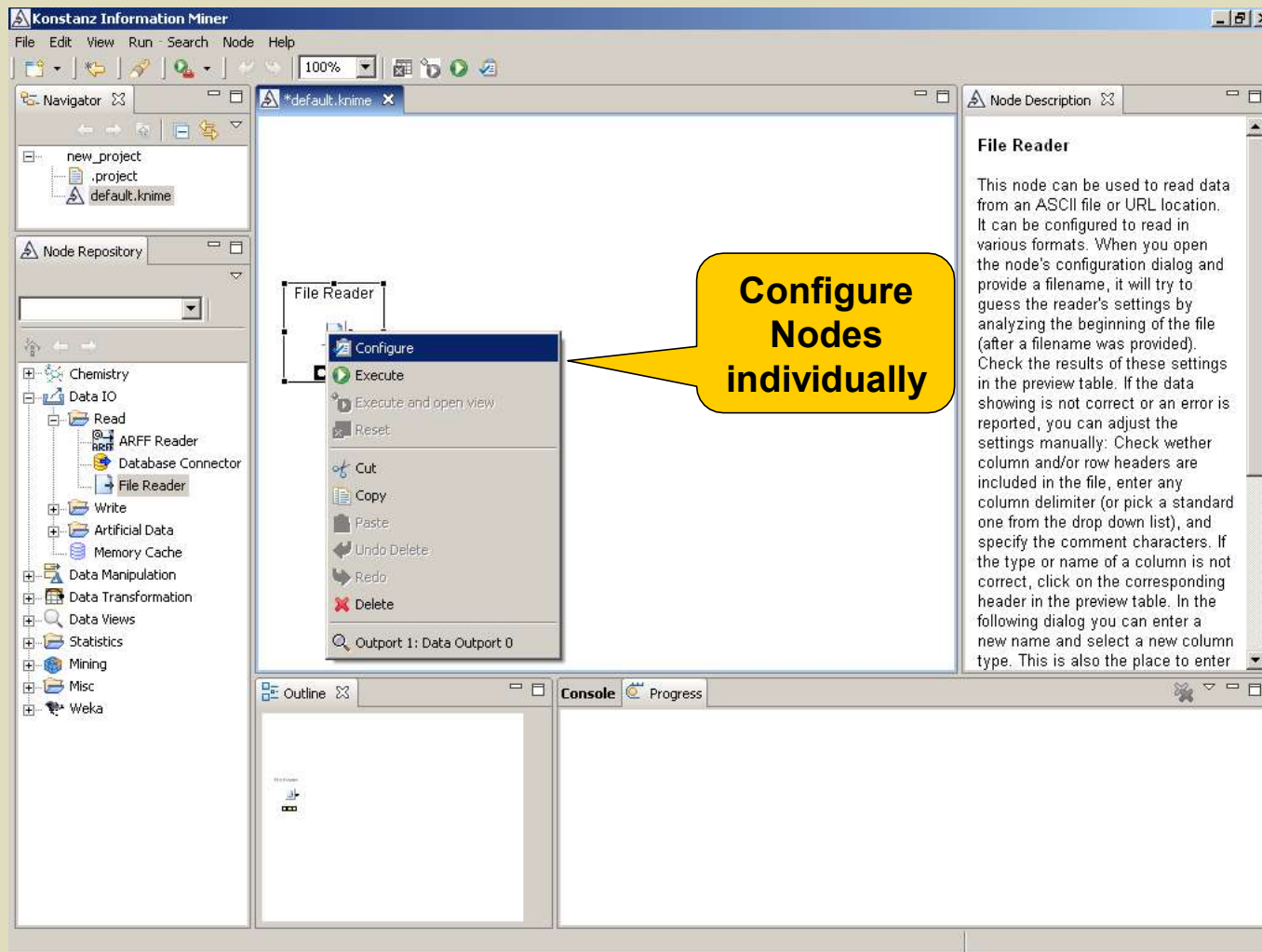
Next:

➢Introduction to KNIME - Demo

# Demo - start

# The User Interface of KNIME 1.0

The original user interface was so simple and intuitive that is still very similar in the latest version.

Execute one or more nodes

Many more views and also other types available…

# Demo - end

# CSMDM16 - Data Analytics and Mining

# Introduction to KNIME - continued

**Dr. Carmen Lam**

Department of Computer Science

carmen.lam@reading.ac.uk

# The Four States of a KNIME Node

- **Not Configured** (<u>red light</u>): node has not been configured yet or the configuration is not valid, or the node has been executed unsuccessfully.

- **Configured** (<u>yellow light</u>): node configuration is valid

- **Executed** (<u>green light</u>): node has been executed successfully and, if applicable, output data and views are ready to be used.

- **Error**: the node is not available: the platform is missing the extension that provides the node.

- ❑ Progress bar: most nodes typically provide information about the progress of their execution (%).

# Data Table

- Contains meta information (spec)
  - data types
  - domains
  - # of rows/cols
- Large tables are buffered on disc
- Blob cell support for large data cells e.g. images



**Bitvector Generator**

# Data Types

- **Common data types**
  - Double Value `D Double ... 0.2`
  - Int Value `I Integer... 1`
  - String Value `S String Col Iris-setosa`
  - Collections `(...) Collecti... [0.2,1]`
    - Sets
    - Lists
  - Bit vectors `BitVectors 10`

- **Additional data types**
  - Terms and Documents
  - Image
  - Network
  - Chemical types
    - Molecules i.e. CDK, Smiles, SDF, …
  - Distance Matrix
  - Custom data types

# KNIME Features

- Node (algorithms) repository
- Node types
  - I/O
  - Data manipulation
  - Learners
  - Predictors
  - Views
- Highlighting
- Metanodes
- Loops and flow variables
- Error handling: "try-catch" nodes
- Extensions (KNIME plugins)

# Preferences

- In the KNIME Analytics Platform
  - In menu File, select Preferences

- E.g., set the level of
  log messages shown
  in the console ➔
    - DEBUG
    - INFO
    - WARN
    - ERROR

# File knime.ini and JVM memory

- The file **knime.ini** (in the KNIME installation folder) sets options used by the Java Virtual Machine when KNIME Analytics Platform is launched.
  - You can allocate more memory for the JVM to run KNIME by changing the **Xmx** JVM argument (default: 512 MB)
- Track memory used
  - File -> Preferences -> General -> Tick "Show heap status" checkbox

```
-startup
plugins/org.eclipse.equinox.launcher_1.4.0.v20161219-1356.jar
--launcher.library
plugins/org.eclipse.equinox.launcher.win32.win32.x86_1.1.551.v20171108-1834
--launcher.defaultAction
openFile
-vm
plugins/org.knime.binary.jre.win32.x86_1.8.0.152-01/jre/bin
-vmargs
-server
-Dsun.java2d.d3d=false
-Dosgi.classloader.lock=classname
-XX:+UnlockDiagnosticVMOptions
-XX:+UnsyncloadClass
-Dsun.net.client.defaultReadTimeout=0
-XX:CompileCommand=exclude,javax/swing/text/GlyphView,getBreakSpot
-Xmx512m
-Dorg.eclipse.swt.browser.IEVersion=10001
-Dsun.awt.noerasebackground=true
-Dequinox.statechange.timeout=30000
```

# Importing/Exporting KNIME Workflows

❑ Many examples are available at KNIME Hub:

- https://hub.knime.com/
- You can search and import nodes and workflows

❑ Users can share (export and import) workflows via archive files (*archive.knwf*)

- File -> Export KNIME workflow
- File -> Import KNIME workflow

# Perspective and Preferences (in KNIME SDK)

- In the SDK version (Eclipse), you need to select the KNIME perspective.

- In Menu Windows, select Preferences

# KNIME Updates and Extensions

- Extensions are additional features (w.r.t. the basic installation) and are installed via the **KNIME update site**.
  - In KNIME click "Help -> Install KNIME extension...".
  - select the KNIME update site (*http://update.knime.com/analytics-platform/3.7)*
  - select the features you want to install in the dialog.
  - You need to restart KNIME after installing new extensions in order to get them activated.

  ➢ If you import an external workflow and a node is "missing" (warning), you then need to install the specific extension that provides that node.

# KNIME Extensions (Plugins/Dropins)

Some available extensions include:

- Chemistry types and features
- Distance Matrix
- Ensemble Learning
- Item Set Mining
- R Statistics Integration
- Python integration
- Weka Data Mining Integration
- HTML/PDF Writer
- Report Designer
- Webservice Client
- XLS Support
- XML Processing
- Cloud connectors (Amazon, Azure, etc.)
- etc.

# KNIME Extensions

- Experimental "Lab" Extensions: http://tech.knime.org/knime-labs
  - JavaScript views
  - Modular Data Generators
  - Network Mining
  - Perl Scripting
  - Text Processing
  - etc.

- Community Contributions: http://tech.knime.org/community
  - Chemoinformatics
  - High Content Screening
  - Image Processing
  - Next Generation Sequencing
  - R/Groovy/Matlab/Python Scripting
  - STARK
  - etc.

- KNIME is designed to be extended!
  - You can create your own KNIME nodes (extensions) by using the KNIME SDK version.

# Conclusions on KNIME

- ## Modularity and extendibility
  - General and extendible data structure (DataTable and DataCell)
  - Nodes encapsulate computational processing tasks (algorithms)

- ## A workflow management system
  - directed edges connects nodes to create data pipelines
  - a workflow is, in general, a directed acyclic graph
  - multi-threading
  - Meta-nodes (nested workflows)

- ## New releases
  - Enhanced GUI and performance
  - Include more and more modules and features

# KNIME Useful Resources

KNIME user ➔ desktop version: "KNIME Analytics Platform"

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

KNIME resources ➔ learning material, examples, etc.:

- https://www.knime.com/learning
- https://www.knime.com/learning-hub
- https://www.knime.com/knimepress
- https://www.youtube.com/user/KNIMETV

For advanced programmers (based on Java Reflection metaprogramming):

KNIME developer ➔ SDK version (Eclipse):

- https://www.knime.com/developers
- https://github.com/knime/knime-sdk-setup
- http://tech.knime.org/developer-guide
- http://tech.knime.org/developer/example
- API: for example see the DataTable interface in
  http://tech.knime.org/docs/api/org/knime/core/data/package-summary.html

# Next:

➢ P01: practical on KNIME Basics

# Next week:

➢ Proximity Measures