# KNIME Pros Learnathon
*Building Reliable and Reusable Components*

## Group 2: Life Sciences

KNIME Team

**@KNIME**
**#Learnathon**

**tinyurl.com/KNIME-Pros-Stuff**

# See you soon in a Breakout Room!



**Main Zoom Session**

### Group 1
**Financial Analysis**

**Maarit**
KNIME Team Member

**Lada**
KNIME Team Member

### Group 2
**Life Sciences**

**temesgen-dadi**
KNIME Team Member

**Francosinus**
KNIME Team Member

### Group 3
**Automation**

**paolotamag**
KNIME Team Member

**Mpattadkal**
KNIME Team Member

**tinyurl.com/KNIME-Pros-Stuff**

KNIME
Open for Innovation

# Three Parallel Tracks via Zoom Breakout Rooms!

**Main Zoom Session**

**Welcome to:**
**Group 2.**
**Life Sciences**

**Group 2**
**Life Sciences**

temesgen-dadi
**KNIME Team Member**

Francosinus
**KNIME Team Member**

Group_2-Life_Sciences
Group_2-Building_a_Component
Group_2-Building_a_Component-Solution

- Download exercises from:
  **tinyurl.com/KNIME-Pros-Stuff**

**OR**

- Search **hub.knime.com** :
  "KNIME Pros Learnathon Group 2"

**tinyurl.com/KNIME-Pros-Stuff**

KNIME
Open for Innovation

# Life Sciences Verified Components



Available on
[hub.knime.com](hub.knime.com)

# FASTA File Format

- common file format in bioinformatics and biochemistry

- text-based format for representing either nucleotide sequences or amino acid (protein) sequences

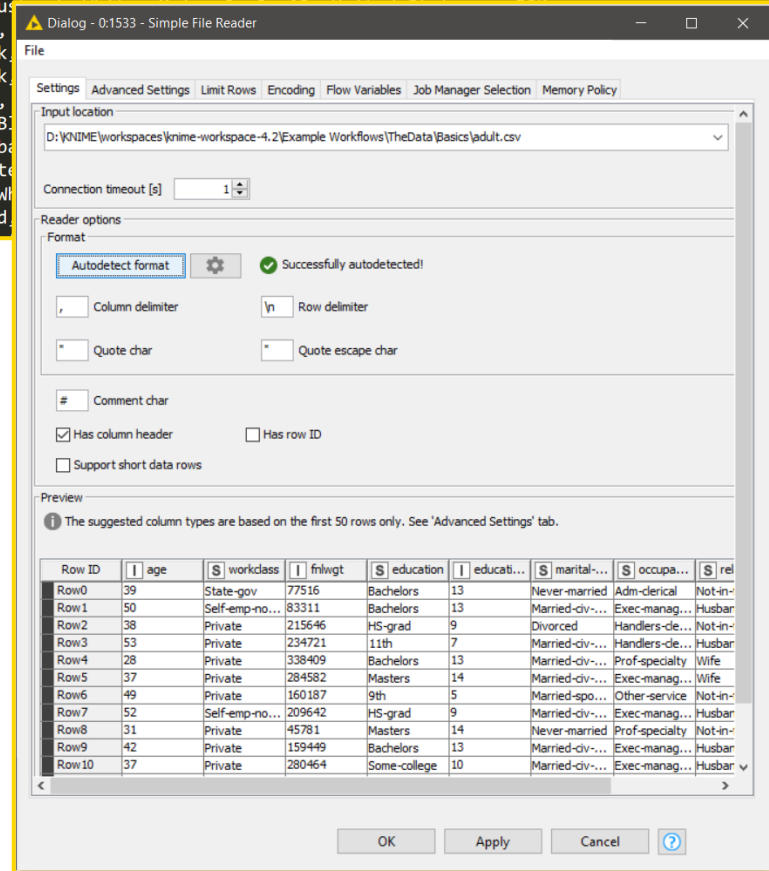- allows for sequence names and comments to precede the sequences

# KNIME Nodes (File Readers)



- Well suited to read tabular data

# FASTA File as Tabular Data



- Re-think Row/Column delimiters
  - ">" Row delimitr
  - "\n" Column delimeter

# Today's Challenge: FASTA Reader Component

- Download Exercises from **link** and import .knar file to your KNIME Analytics Platform LOCAL Workspace!



**FASTA Reader**

Select a File

Group_2-Life_Sciences
   Group_2-Building_a_Component
   Group_2-Building_a_Component-Solution

**OR**

- Download from KNIME Hub!

hub.knime.com/knime/spaces/Education/latest/Learnathons
or Search for "**KNIME Pros Learnathon - Group 2**"

Open for Innovation
KNIME