# Pathway Enrichment with KNIME

🔒 InesAssum / **EnrichmentNodes** `Private`

👁 Watch ▾ | 1     ★ Star | 0     ⑂ Fork | 0

<> **Code**   ⊘ Issues 0   ⑂ Pull requests 0   ▥ Projects 0   ᐧ�◨ᐧ Insights   ⚙ Settings

R framework integrating different multiOMICs enrichment methods and translating them into KNIME nodes using genericworkflownodes

Edit

1.  **Pathway Enrichment using a modular R library**

2.  **Pathway Enrichment with KNIME for Users**

3.  **Pathway Enrichment with KNIME for Developers**

**InesAssum** added presentation and example workflow
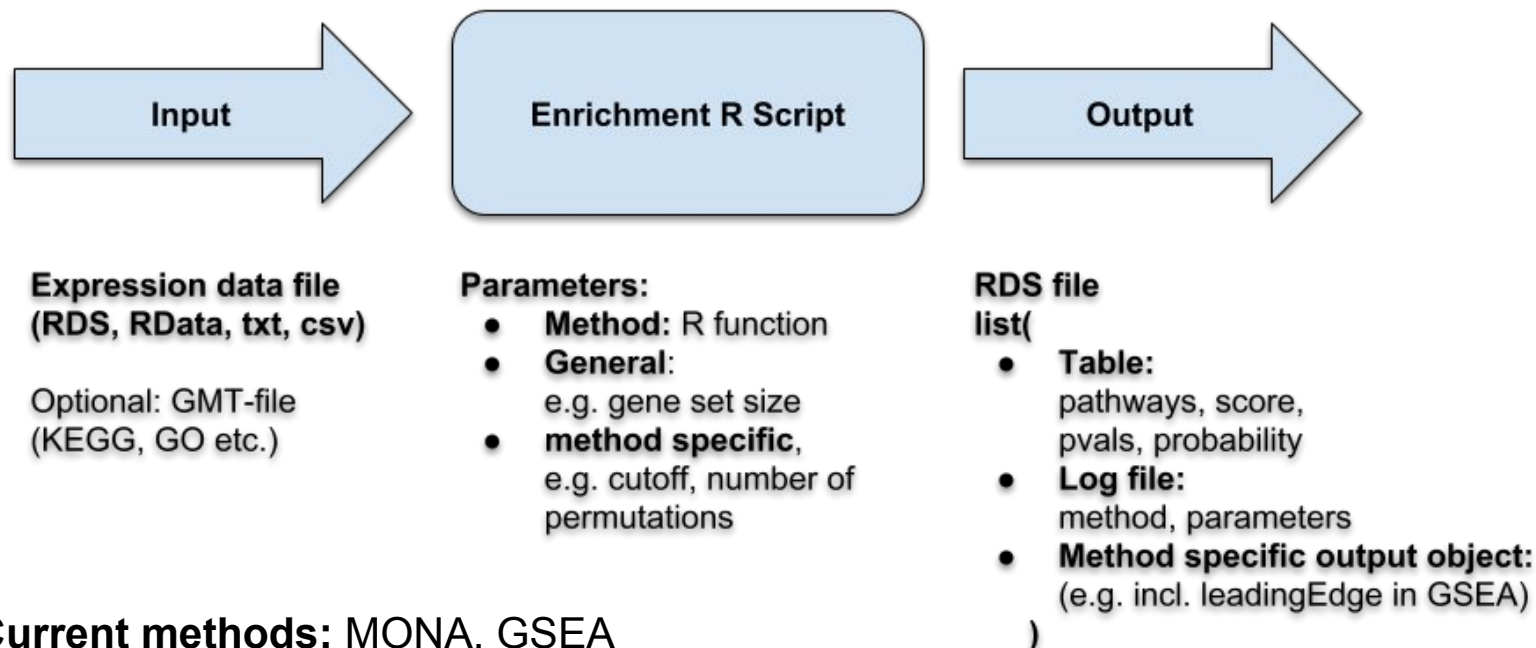
📁 examples/KNIME_workflows

📁 knime

📁 src

📁 tutorials

📄 .gitignore

📄 LICENSE

📄 README.md

2

# Modular R library integrating multiple Pathway Enrichment methods

# Modular R Framework



**Input**
**Expression data file (RDS, RData, txt, csv)**

Optional: GMT-file (KEGG, GO etc.)

**Current methods:** MONA, GSEA

**Enrichment R Script**
Parameters:
- **Method:** R function
- **General**:
  e.g. gene set size
- **method specific**,
  e.g. cutoff, number of permutations

**Output**
**RDS file**
list(
- **Table:**
  pathways, score, pvals, probability
- **Log file:**
  method, parameters
- **Method specific output object:**
  (e.g. incl. leadingEdge in GSEA)
)

How to deal with different input data requirements,

i.e. raw data vs. summary statistics?

GSEA: Subramanian et al., PNAS, 2005,             MONA: Sass et al., NAR, 2013

# Modular R Framework

- Central runner script
- calls methods as R functions
- Idea: methods share common starting point (raw data)
- necessary preprocessing handled internally



```
#!/usr/bin/env Rscript
# -------------------------------------------------------------------------------
#' Runner script for multiOMICs pathway enrichment
#' @param method string // Select method for pathway enrichment
#' @param data string // file path to .RDS input file: named numerical vector, names=gene symbols#
#' @param result string // file path to save results as .RDS
#' @param gmt string // pathway definition to use: KEGG, GO... default: KEGG
#' @param minSize integer // min size of GSs to be considered, default: 15
#' @param maxSize integer // max size of GSs to be considered, default: 500
#' @param nperm integer // number of permutations, default: 10000
#' @param mygmt string // custom .gmt file // not required
# -------------------------------------------------------------------------------
```

# Modular R Framework

```
C:\Users\krisg\Documents>Rscript run_PEanalysis.R -m fgsea -input data.csv -o result.csv -g KEGG
[1] "Processing.........."
[1] "Completed successful! See log file for more information."
```

- Run from command line
- run locally in R or
- use Docker support
- Full code and documentation available on GitHub

**Suited for:**  **-** big/advanced projects (simulations, benchmarking)

**-** anyone, who is using R anyway

...but what if you don't have a programming background ?

# KNIME Enrichment for Users

*KNIME - The Konstanz Information Miner,* M. R. Berthold, B. Wiswedel, et al.

7

# KNIME Analytics Platform

KNIME is a free software for interactive data analysis



**Run**

# Introducing KNIME EnrichmentNodes

- KNIME Nodes are built on top of modular enrichment R scripts
- EnrichmentNodes use Docker containers that include
  - all software
  - all libraries and other dependencies
  - independent of local OS
- Docker image stored online at dockerhub.io and gets downloaded automatically

# KNIME Pathway Enrichment with EnrichmentNodes

Requirements:

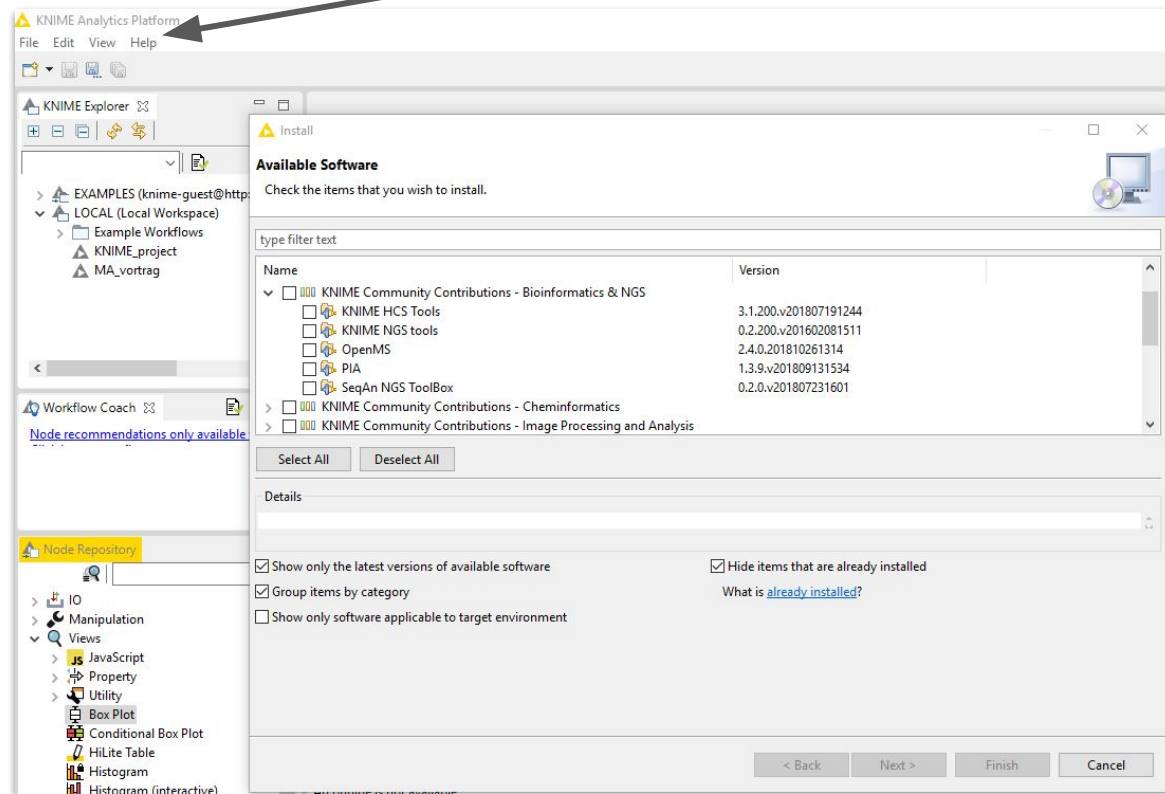Platform independent using Docker Images

CPU must support VT-x or AMD-v

- KNIME >= 3.1: http://www.knime.org
- Docker >= 1.9: https://www.docker.com/
- Generic KNIME Node (with Docker support): (https://github.com/genericworkflownodes/GenericKnimeNodes)

# KNIME Pathway Enrichment with "EnrichmentNodes"

Requirements: KNIME, Docker, CPU supporting VT-x or AMD-v

- Install extensions in KNIME: GKN, EnrichmentNodes (Help->Install software…)

- Docker loads automatically

- alternatively build your own version locally (temporary)

# Using KNIME

# KNIME Enrichment for Developers

# Integrating own Methods into KNIME Nodes

**First Step:** integrate your R script/tool into our modular R framework

```r
# -----------------------------------------------------------------
#' Script to run GSEA analysis
#'
#' @author Ines Assum
#' @param input string // .RDS input file: named numerical vector (gene symbols)
#' @param output string // file path to save results as .RDS
#' @param gmt string // Gene set definition (KEGG / GO / custom [mygmt])
#' @param minSize integer // min size of GSs to be considered, default: 15
#' @param maxSize integer // max size of GSs to be considered, default: 500
#' @param nperm integer // number of permutations, default: 10000
#' @param mygmt string // custom .gmt file // not necessary
# -----------------------------------------------------------------

run_fgsea <- function(input, output, gmt, minSize, maxSize, nperm, mygmt){
  library(fgsea)
  pathways <- gmtPathways(gmt)
  GSEA <- fgsea(pathways,
                rank,
                nperm,
                minSize=minSize,
                maxSize=maxSize)
}
```

fgsea: Sergushichev, bioRxiv, 2016

# Integrating own Methods into KNIME

**Step 2:  Provide a CTD-file**

```xml
<tool name="YourTool.R" >
    <description>Draw plot for a nucleotide sequence.</description>
        <cli>
        <clielement optionIdentifier="-sequence">
            <mapping referenceName="YourTool.sequence" />
        </clielement>
        <clielement optionIdentifier="-outfile">
            <mapping referenceName="YourTool.outputfile" />
        </clielement>
    </cli>
    <PARAMETERS >
        <NODE name="YourTool" description="Draw plot for a nucleotide sequence.">
            <ITEM name="sequence" value="" description="sequence filename" supported_formats="*.fasta"/>
            <ITEM name="outputfile" value="" description="Output file." supported_formats="*.ps"/>
        </NODE>
    </PARAMETERS>
</tool>
```

# Integrating own Methods into KNIME

**Step 3:   Specify plugin.properties file**

```
1   # the package of the plugin
2   pluginPackage=de.enrichment
3
4   # the name of the plugin
5   pluginName=EnrichmentNodes
6
7   # the version of the plugin
8   pluginVersion=1.0.0.0
9
10  # the path (starting from KNIMEs Community Nodes node)
11  nodeRepositoyRoot=community
12
13  executor=com.genericworkflownodes.knime.execution.impl.LocalDockerToolExecutor
14  commandGenerator=com.genericworkflownodes.knime.execution.impl.DockerCommandGenerator
15
16  #docker specific configurations
17  dockerMachine=default
18
19  #tool specific configurations
20  tool.monal.dockerImage=enrich
21  tool.gsea.dockerImage=enrich
22  tool.enrich.dockerImage=enrich
```

# Integrating own Methods into KNIME Nodes

**4. Provide predefined directory**

```
Plugin_dir
    |
    |--- plugin.properties   (description file)
    |
    |--- descriptors
    |       |-- yourTool.ctd
    |       |-- mime.types
    |
    |--- DESCRIPTION         (short description of project)
    |
    |--- LICENSE             (Licensing information)
    |
    |--- COPYRIGHT           (Copyright information)
```

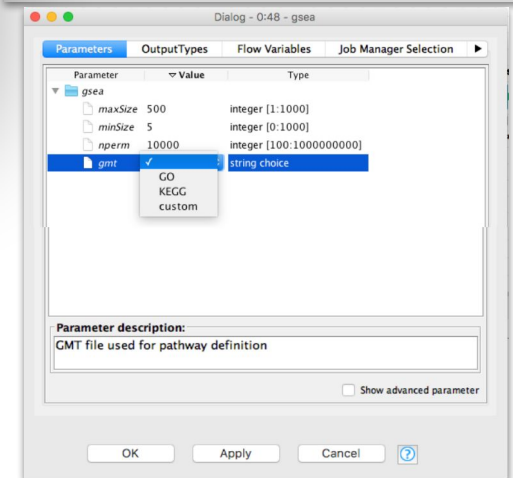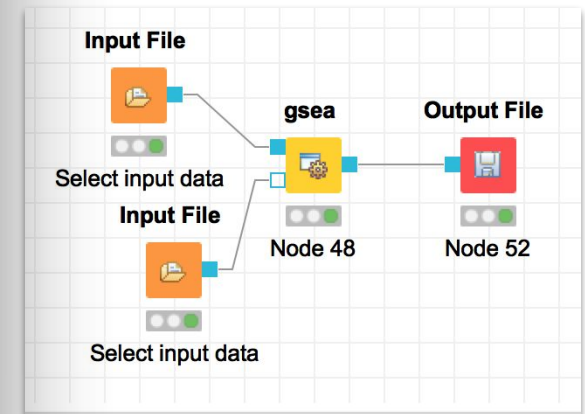**5. Get GenericKnimeNodes² and run ant**

**6. Import and run on KNIME SDK**

**Voila!**

[2]: GenericKNIMENodes (GKN)

# Example

# Example

# Example

```
# -----------------------------------------------------------------
#' Script to run GSEA analysis
#'
#' @author Ines Assum
#' @param input string // .RDS input file: named numerical vector (gene symbols
#' @param output string // file path to save results as .RDS
#' @param gmt string // Gene set definition (KEGG / GO / custom [mygmt])
#' @param minSize integer // min size of GSs to be considered, default: 15
#' @param maxSize integer // max size of GSs to be considered, default: 500
#' @param nperm integer // number of permutations, default: 10000
#' @param mygmt string // custom .gmt file // not necessary
# -----------------------------------------------------------------

run_fgsea <- function(input, output, gmt, minSize, maxSize, nperm, mygmt){
  library(fgsea)
  pathways <- gmtPathways(gmt)
  GSEA <- fgsea(pathways,
                rank,
                nperm,
                minSize=minSize,
                maxSize=maxSize)
}
```



```
1  <tool name="gsea">
2      <description>Runs GSEA using fgsea R package.</description>
3      <manual>Detailed description goes here.</manual>
4      <executableName>run_fgsea.R</executableName>
5      <PARAMETERS >
6          <NODE name="gsea" description="Node the runs GSEA">
7              <ITEM description="GMT file used for pathway definition" name="gmt" restrictions="GO,KEGG,c
8              <ITEM description="min Size" name="minSize" restrictions="0:1000" type="int" value="15"/>
9              <ITEM description="max Size" name="maxSize" restrictions="1:1000" type="int" value="500"/>
10             <ITEM description="Permutations" name="nperm" restrictions="100:1000000000" type="int" valu
11         </NODE>
12     </PARAMETERS>
13     <cli>
14         <clielement optionIdentifier="-gmt">
15             <mapping referenceName="gsea.gmt"/>
16         </clielement>
17         <clielement optionIdentifier="-min">
18             <mapping referenceName="gsea.minSize"/>
19         </clielement>
20         <clielement optionIdentifier="-max">
21             <mapping referenceName="gsea.maxSize"/>
22         </clielement>
23         <clielement optionIdentifier="-p">
24             <mapping referenceName="gsea.nperm"/>
25         </clielement>
26     </cli>
27  </tool>
```
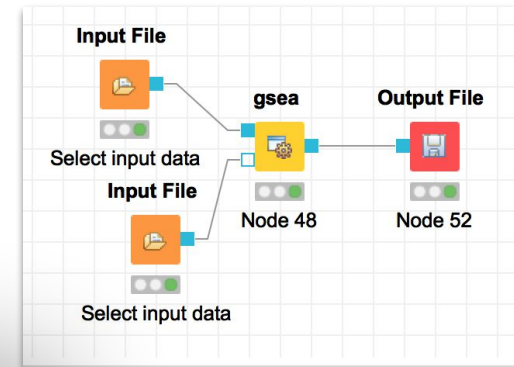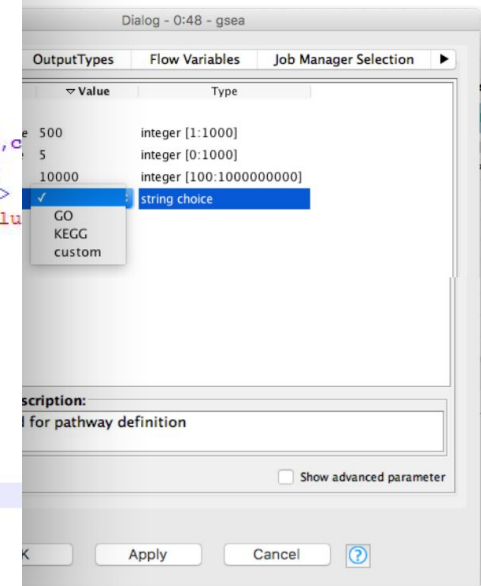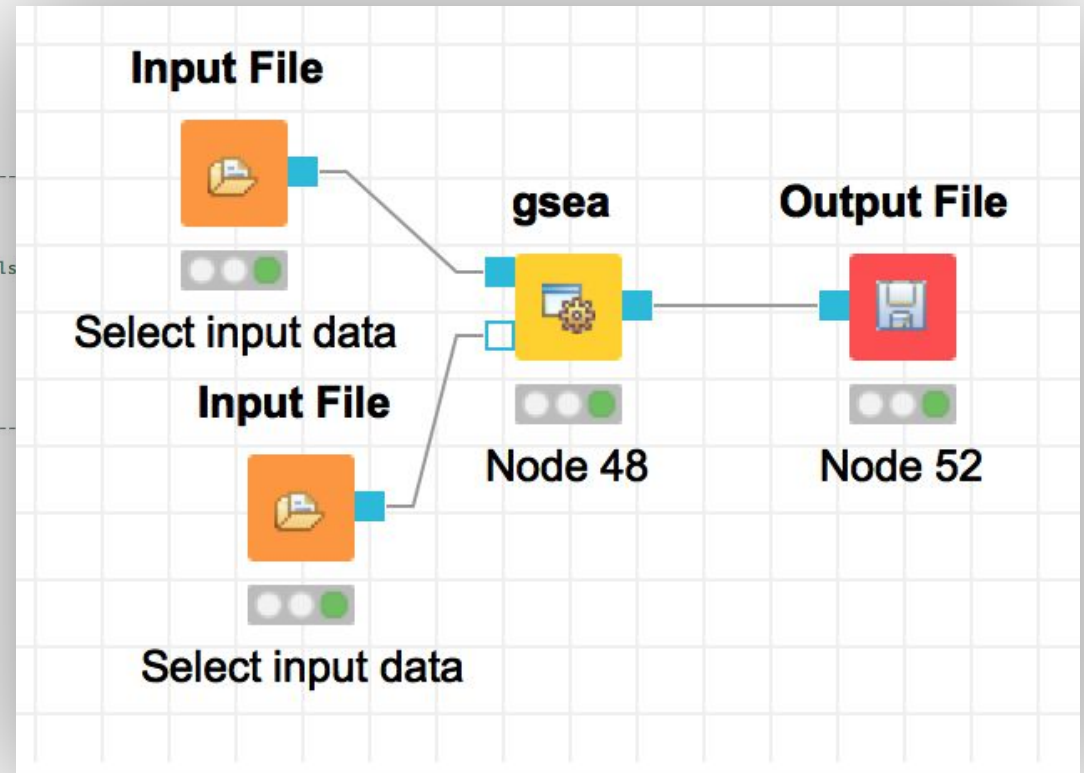
# Example

```
# ----------------------------------------------------------------
#' Script to run GSEA analysis
#'
#' @author Ines Assum
#' @param input string // .RDS input file: named numerical vector (gene symbols)
#' @param output string // file path to save results as .RDS
#' @param gmt string // Gene set definition (KEGG / GO / custom [mygmt])
#' @param minSize integer // min size of GSs to be considered, default: 15
#' @param maxSize integer // max size of GSs to be considered, default: 500
#' @param nperm integer // number of permutations, default: 10000
#' @param mygmt string // custom .gmt file // not necessary
# ----------------------------------------------------------------

run_fgsea <- function(input, output, gmt, minSize, maxSize, nperm, mygmt){
  library(fgsea)
  pathways <- gmtPathways(gmt)
  GSEA <- fgsea(pathways,
                rank,
                nperm,
                minSize=minSize,
                maxSize=maxSize)
}
```

```xml
1  <tool name="gsea">
2      <description>Runs GSEA using fgsea R package.</description>
3      <manual>Detailed description goes here.</manual>
4      <executableName>run_fgsea.R</executableName>
5      <PARAMETERS >
6          <NODE name="gsea" description="Node the runs GSEA">
7              <ITEM description="GMT file used for pathway definition" name="gmt" restrictions="GO,KEGG,c
8              <ITEM description="min Size" name="minSize" restrictions="0:1000" type="int" value="15"/>
9              <ITEM description="max Size" name="maxSize" restrictions="1:1000" type="int" value="500"/>
10             <ITEM description="Permutations" name="nperm" restrictions="100:1000000000" type="int" valu
11         </NODE>
12     </PARAMETERS>
13     <cli>
14         <clielement optionIdentifier="-gmt">
15             <mapping referenceName="gsea.gmt"/>
16         </clielement>
17         <clielement optionIdentifier="-min">
18             <mapping referenceName="gsea.minSize"/>
19         </clielement>
20         <clielement optionIdentifier="-max">
21             <mapping referenceName="gsea.maxSize"/>
22         </clielement>
23         <clielement optionIdentifier="-p">
24             <mapping referenceName="gsea.nperm"/>
25         </clielement>
26     </cli>
27  </tool>
```
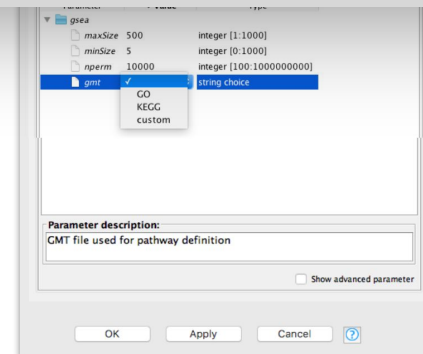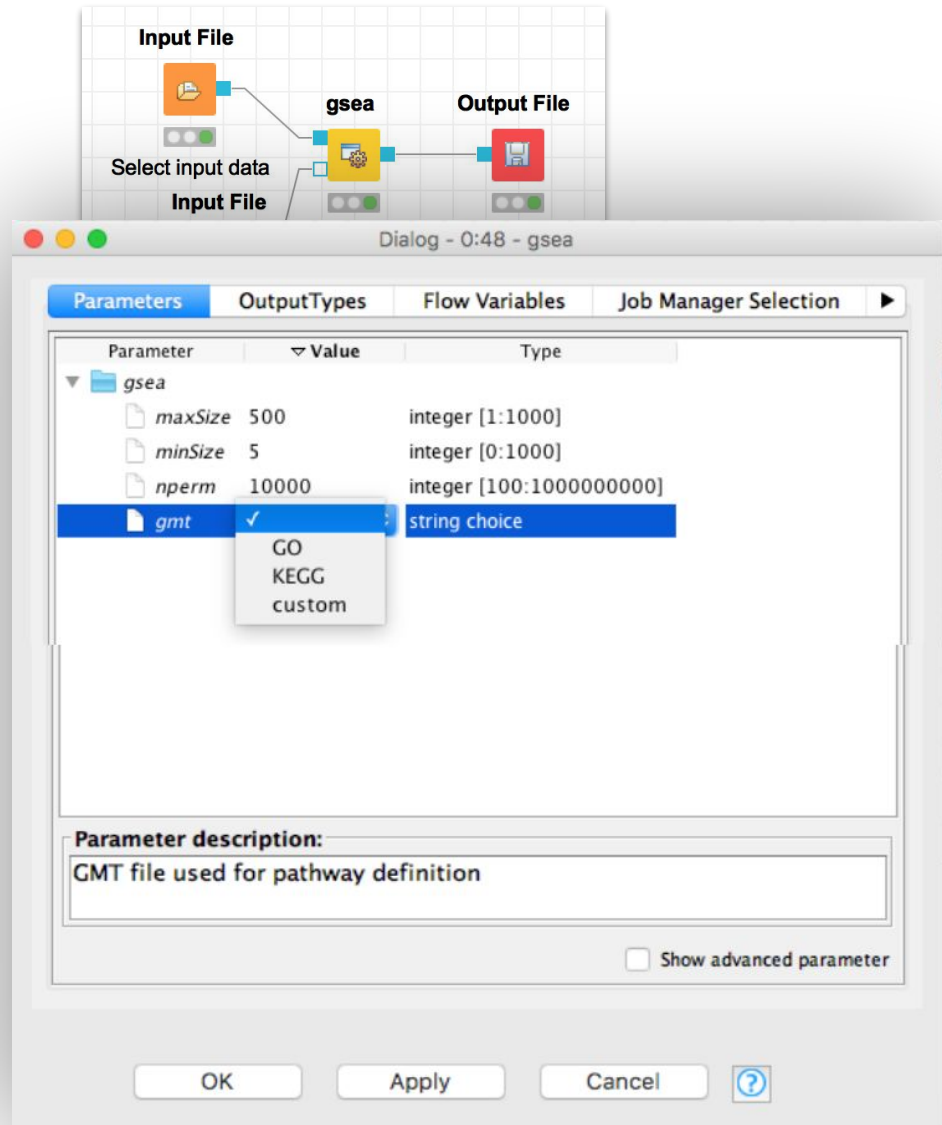
**Input File**

Select input data

**Input File**

gsea        **Output File**

**Dialog - 0:48 - gsea**

| Parameters | OutputTypes | Flow Variables | Job Manager Selection ▶ |

| Parameter | ▽ Value | Type |
|---|---|---|
| ▼ 📁 gsea | | |
| 📄 maxSize | 500 | integer [1:1000] |
| 📄 minSize | 5 | integer [0:1000] |
| 📄 nperm | 10000 | integer [100:1000000000] |
| 📄 gmt | ✓ | string choice |
| | GO | |
| | KEGG | |
| | custom | |

**Parameter description:**

GMT file used for pathway definition

☐ Show advanced parameter

OK        Apply        Cancel        ⑦

```
# --------------------------------------------------------------------
#' Script to run GSEA analysis
#'
#' @author Ines Assum
#' @param input string // .RDS input file: named numerical vector (gene symbols)
```



Dialog - 0:48 - gsea

| Parameters | OutputTypes | Flow Variables | Job Manager Selection | ▶ |

| Parameter | ▽ Value | Type |
|---|---|---|
| ▼ 📁 *gsea* | | |
| 📄 *maxSize* | 500 | integer [1:1000] |
| 📄 *minSize* | 5 | integer [0:1000] |
| 📄 *nperm* | 10000 | integer [100:1000000000] |
| 📄 *gmt* | ✓ | string choice |
| | GO | |
| | KEGG | |
| | custom | |

**Parameter description:**

GMT file used for pathway definition

☐ Show advanced parameter

OK     Apply     Cancel     ⑦

# Summary

- Modular R library integrated into easy-to-use KNIME nodes
- Expand functionality with minimal effort
- We will soon make the "EnrichmentNodes" project on GitHub public:

  -> https://github.com/InesAssum/EnrichmentNodes
- Template for exemplary node provided
- Collect and share methods
- Join our slack workspace for discussion and support

  -> Slack.com/KNIME-setup



- Check out Benni's ImmunoNodes:  github.com/FRED-2/ImmunoNodes/

# Acknowledgements

**Epigenereg group**

Matthias Heinig

Ines Assum

Thomas Walzthöni

**ICB (MONA)**

Nikola Müller

Florian Büttner

Andreas Kopf

**GenericKnimeNodes**     &     **ICB**

Julianus Pfeuffer

Benjamin Schubert

-> ImmunoNodes