

# KNIME Pipeline User Manual

Eva Freckmann

## Contents

<b>1</b>	<b>Before starting</b>	<b>2</b>
1.1	KNIME installation . . . . .	2
1.2	R Integration . . . . .	3
1.3	Python Integration . . . . .	3
<b>2</b>	<b>Workflow structure</b>	<b>4</b>
<b>3</b>	<b>Data import and pre-processing</b>	<b>4</b>
3.1	Select Directory and match experiment key . . . . .	4
3.2	Manually Set Output Folder . . . . .	4
3.3	Select Samples and Frames to Include . . . . .	4
3.4	Create and Populate ShapeClassification column . . . . .	4
3.5	Duplicated Tracking Label Correction . . . . .	4
3.6	Edit Column Names, Create Unique-object_id, Filter Columns, and Normalise . . . . .	4
<b>4</b>	<b>Data summary files</b>	<b>4</b>
4.1	PCA of Replicates . . . . .	5
4.2	Calculate TimeChunks . . . . .	5
4.3	Measurements Over Time . . . . .	5
4.4	Heatmap/Line Plot of Spheroid Number Over Time, Heatmap of Initial Spheroid Number . . . . .	6
<b>5</b>	<b>User-defined Shape Classification</b>	<b>6</b>
5.1	Find Trajectories . . . . .	6
5.2	Compare Proportions of Trajectories . . . . .	6
5.3	Find Trajectories in One Image Sequence . . . . .	6
5.4	Trajectory Timechunk Plots (Frequency Motifs and Transitions) . . . . .	6
5.5	Get Representative spheroid per Trajectory . . . . .	6
5.6	Subsample or use Entire Dataset . . . . .	6
5.7	Perform tSNE . . . . .	6
5.8	tSNE Plotting . . . . .	6
5.9	Perform tSNE using different variable combinations . . . . .	6
5.10	Get representative outlines for each classification . . . . .	6
5.11	Colour and Overlay Outlines . . . . .	6
5.12	Calculate TimeChunks . . . . .	6
5.13	ShapeClassification Over Time . . . . .	6
5.14	ShapeClassification Mean Measurements . . . . .	6
<b>6</b>	<b>Unsupervised Shape Classification</b>	<b>6</b>
6.1	Subsample or use Entire Dataset . . . . .	7
6.2	Identify Subpopulations . . . . .	7
6.3	Find Trajectories . . . . .	7
6.4	Find Trajectories in One Movie . . . . .	7

6.5	Compare Proportions of Trajectories . . . . .	7
6.6	Trajectory Timechunk Plots (Frequency Motifs and Transitions) . . . . .	7
6.7	Get representative spheroid per Trajectory . . . . .	7
6.8	Compare Proportions of PhenoGraph Clusters . . . . .	7
6.9	Get representative outlines for each cluster . . . . .	7
6.10	Subsample or use Entire Dataset . . . . .	7
6.11	Perform tSNE using different variable combinations . . . . .	7
6.12	Perform tSNE . . . . .	7
6.13	tSNE Plotting . . . . .	7
6.14	Colour and Overlay Outlines . . . . .	7

## 1 Before starting

The computer you will be running the KNIME analysis from should have either access to, or a local, copy of the directory containing your phase images, and output from CellProfiler. For reference, this directory should have the following structure:

1. DatasetName
  - Experiment\_\_1
    - Data
    - Experiment Key
    - Phase
    - PhaseGrayCystOutlines
  - Experiment\_\_n
  - Output

The “Experiment Key” and “Phase” folders should be populated by the user prior to analysis with CellProfiler, which will in turn populate “Data” and “PhaseGrayCystOutlines” upon analysis completion.

The default output directory for KNIME analysis results will be the “Output” subdirectory of this folder.

Experiment Key and directory structure templates can be found in the [davebryantlab/MethodsPaper2020](#) Github repository.

**Tip for those new to KNIME:** Double-click on a meta-node/node to configure it.

### 1.1 KNIME installation

KNIME can be downloaded from [here](#).

#### 1.1.1 Notes for KNIME workflow

**Required KNIME version and extensions:**

Name	Version
KNIME Analytics Platform	4.0.2.v201909300912
KNIME Interactive R Statistics Integration	4.0.1.v201908131226
KNIME Core	4.0.2.v201909300912
KNIME Quick Forms	4.0.2.v201909242005
KNIME Math Expression (JEP)	4.0.2.v201909242005
KNIME Python Integration	4.0.0.v201906241606
KNIME Image Processing	1.8.0.201911140609
KNIME SVG Support	4.0.1.v201908131226
KNIME Virtual Nodes	4.0.0.v201905311239
KNIME Distance Matrix	4.0.2.v201909260824
KNIME Data Generation	4.0.0.v201905311239

Name	Version
KNIME File Handling Nodes	4.0.1.v201908131226
Vernalis KNIME Nodes	1.24.2.v201911141223
KNIME Excel Support	4.0.1.v201908131226
KNIME HCS Tools	4.0.0.v201906200802

## 1.2 R Integration

The KNIME workflow uses an R integration for some steps of the analysis. The workflow uses R version **version**, which can be downloaded here. Instructions for setting up the R integration can be found here. **Windows users:** Do not install Rserve version 1.7-3.1 as is suggested in point 1 of the “R packages installation” section of these instructions. Instead, go straight to point 2 of the section, to install Rserve v1.8-6. More information on installing Rserve can be found here.

**All users:** In KNIME Analytics Platform go to File → Preferences. From the list on the left, select R under KNIME. Set the “Rserve receiving buffer size limit” to 0.

The KNIME workflow should automatically download and install any missing R packages that are required for the analysis. **Check that cytofit does this** The following packages are utilised:

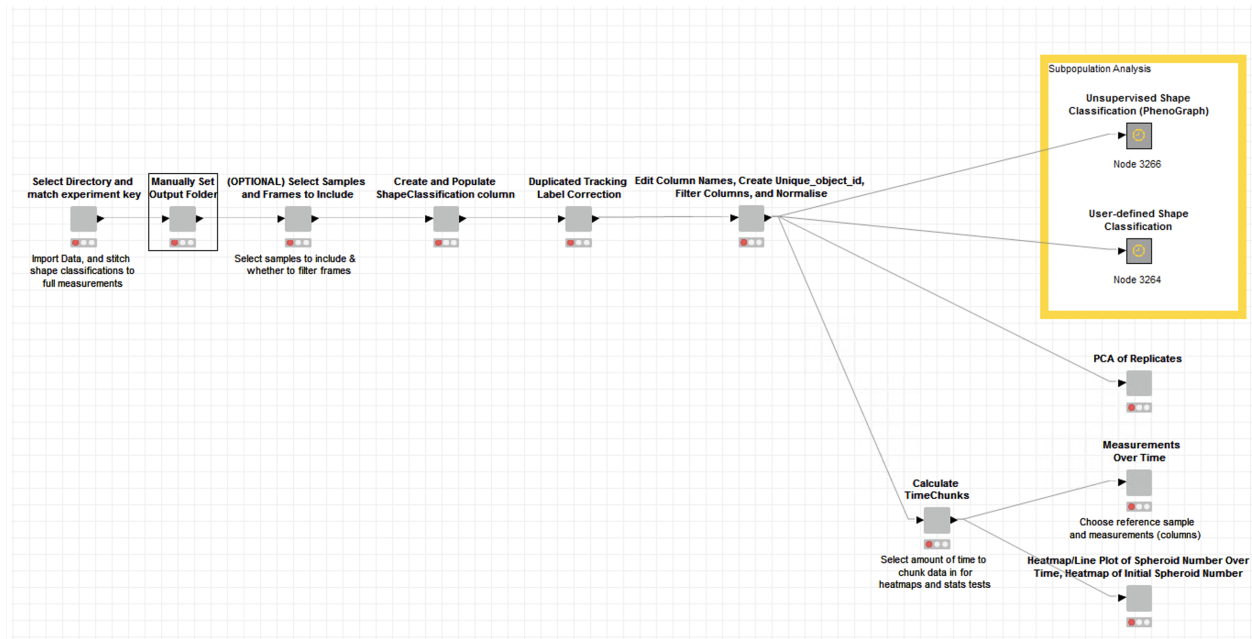
Package	Version
ggplot2	
reshape2	
ggnewscale	

## 1.3 Python Integration

The KNIME workflow requires Python in order to run the GeoSketch algorithm for subsampling data. Instructions for setting up the Python integration can be found here. First, follow the “Quickstart” and “Anaconda Setup” instructions, and download the Python environments provided on the **davebryant-lab/MethodsPaper2020** Github repository. Load these environments into Anaconda - this can be easily done using the Anaconda Navigator application, instructions for this can be found here under “Importing an environment”. Then follow the “Setting up the KNIME Python Integration” instructions, and select the provided environments within your KNIME Python preferences.

Package	Version
GeoSketch	

## 2 Workflow structure



The KNIME workflow is comprised of four parts, each addressed in one section of this manual:

Data Import and Pre-processing

Data Summary Files

User-Defined Shape Classification

Unsupervised Shape Classification

## 3 Data import and pre-processing

Placeholder

### 3.1 Select Directory and match experiment key

### 3.2 Manually Set Output Folder

### 3.3 Select Samples and Frames to Include

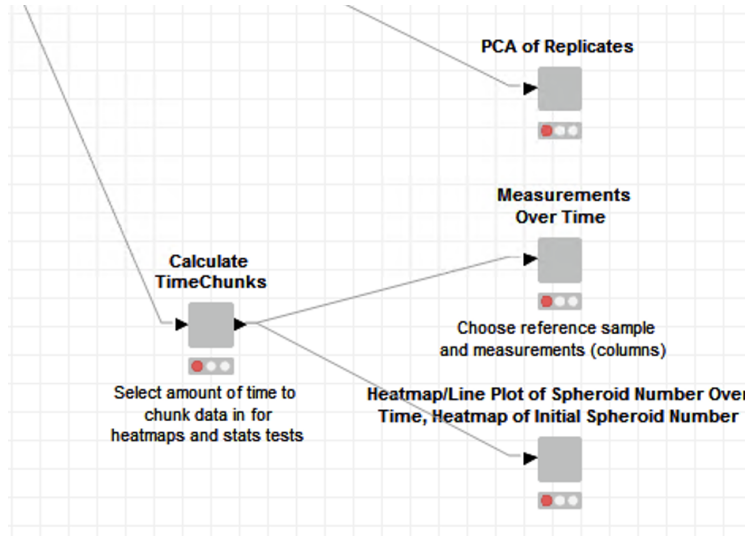
### 3.4 Create and Populate ShapeClassification column

### 3.5 Duplicated Tracking Label Correction

### 3.6 Edit Column Names, Create Unique-object\_id, Filter Columns, and Normalise

## 4 Data summary files

*think of a better name for this section*



## 4.1 PCA of Replicates

This metanode performs Principal Component Analysis (PCA) on the dataset replicates, and outputs results for two PCAs: one in which each point is a well, and another in which each point is an experiment. A plot of PCA Loadings is also output for each analysis. Four files are saved in the main Output directory: “PCAofReplicates\_perWell.pdf”, “PCAofReplicates\_perWell\_Loadings.pdf”, “PCAofReplicates\_perExperiment.pdf”, “PCAofReplicates\_perExperiment\_Loadings.pdf”.

No configuration required.

## 4.2 Calculate TimeChunks

Visualisation of the change in a variable over time can be complex to display when many timepoints are present. To simplify this, data can be presented in grouped segments of time, rather than at individual timepoints. This metanode calculates time interval chunks for the input data table.

Configure the metanode to indicate how many timepoints (frames) to include in each timechunk. If chunking is not required, set this value to 1. **check this works**

## 4.3 Measurements Over Time

This metanode generates a heatmap showing change in measurements of area, shape, and movement over time (for a set of user-specified measurements). For each sample, the average value of the measurement is calculated at each of the previously defined time chunk intervals. For the purposes of presentation, resulting values are Z-score normalised per measurement. A t-tests are performed to compare samples at each time interval, to the specified reference sample. A Bonferroni adjustment is applied to adjust for multiple testing. Two files are generated in the main Output directory: “MeasurementsOverTime\_CONTROLSAMPLEControl.pdf”, “MeasurementsOverTime\_CONTROLSAMPLEControl.csv”

Configure the metanode to set a reference sample for statistical comparison. Measurements to be plotted in the heatmap should be included in the green “Include” box - all others should remain in the red “Exclude” box. Finally, indicate how rows (samples) in the heatmap should be ordered: alphabetically or in user-specified ordering. If a user-specified order is to be used, use the text box to define this order (from top row of the heatmap, to bottom). Sample names should be written as they appear in the Experiment Key, separated by only commas. The easiest way to ensure this is done correctly is by copy and pasting sample names directly from the Experiment Key. **Tip:** Be aware of any leading or trailing spaces in sample names in the Experiment Key, and if present, ensure not to delete them when configuring this text box.

## 4.4 Heatmap/Line Plot of Spheroid Number Over Time, Heatmap of Initial Spheroid Number

*This metanode outputs the initial counts of spheroids in each treatment group/sample per experiment (CSV file), and the relative change over time (line plot). Three files are generated in the main Output directory: “InitialSpheroidNumbersPerExperiment.csv”, “SpheroidNumberOverTime.csv”, “SpheroidNumberOverTime\_LinePlot.pdf”*

No configuration required.

## 5 User-defined Shape Classification

Placeholder

### 5.1 Find Trajectories

### 5.2 Compare Proportions of Trajectories

### 5.3 Find Trajectories in One Image Sequence

### 5.4 Trajectory Timechunk Plots (Frequency Motifs and Transitions)

### 5.5 Get Representative spheroid per Trajectory

### 5.6 Subsample or use Entire Dataset

### 5.7 Perform tSNE

### 5.8 tSNE Plotting

### 5.9 Perform tSNE using different variable combinations

### 5.10 Get representative outlines for each classification

### 5.11 Colour and Overlay Outlines

### 5.12 Calculate TimeChunks

### 5.13 ShapeClassification Over Time

### 5.14 ShapeClassification Mean Measurements

## 6 Unsupervised Shape Classification

Placeholder

- 6.1 Subsample or use Entire Dataset
- 6.2 Identify Subpopulations
- 6.3 Find Trajectories
- 6.4 Find Trajectories in One Movie
- 6.5 Compare Proportions of Trajectories
- 6.6 Trajectory Timechunk Plots (Frequency Motifs and Transitions)
- 6.7 Get representative spheroid per Trajectory
- 6.8 Compare Proportions of PhenoGraph Clusters
- 6.9 Get representative outlines for each cluster
- 6.10 Subsample or use Entire Dataset
- 6.11 Perform tSNE using different variable combinations
- 6.12 Perform tSNE
- 6.13 tSNE Plotting
- 6.14 Colour and Overlay Outlines