

# ICCL

---

Welcome! This manual contains instructions how to use the ICCL tool. The aim of this project is to provide a process discovery tool that enhances the comprehensibility of discovered object-centric Petri nets by clustering. The prerequisites to apply this tool are knowledge in Process Mining, object-centric Process Mining and unsupervised learning techniques such as clustering. The main pipeline of the application works the following way:

1. After the installation and successful start of the server, you can open the application. By default, you start on the data import page.
2. On the import page you can upload an object-centric event log from your local machine.
3. You then run through the clustering pipeline. This includes configuring the set of selected features, the graph embedding method, the clustering method as well as additional tools such as cluster evaluation.
4. You apply process discovery and access the process models based on the clustering before. You now have a more comprehensible representation of the original process model.

## Table of Content

1. [Installation](#)
2. [Application Layout](#)
3. [Data Management](#)
  - [Import OCEL](#)
  - [View OCEL](#)
4. [Configuration Pipeline](#)
  - [Feature Selection](#)
  - [Graph Embedding](#)
  - [Clustering](#)
  - [Cluster Evaluation](#)
  - [Process Executions](#)
5. [Process Discovery](#)
  - [Discover Process Model](#)
    - [Comprehensive Process Model](#)
    - [Cluster Process Models](#)
  - [Conformance Checking](#)
  - [Export Process Model](#)
6. [Algorithms](#)
  - [Process Execution Extraction](#)
  - [Graph Embedding Methods](#)
  - [Clustering Techniques](#)
  - [Evaluation Measures](#)
    - [Cluster Evaluation Measures](#)
    - [Model Evaluation Measures](#)
  - [Model Discovery](#)

# Installation

## Using Docker

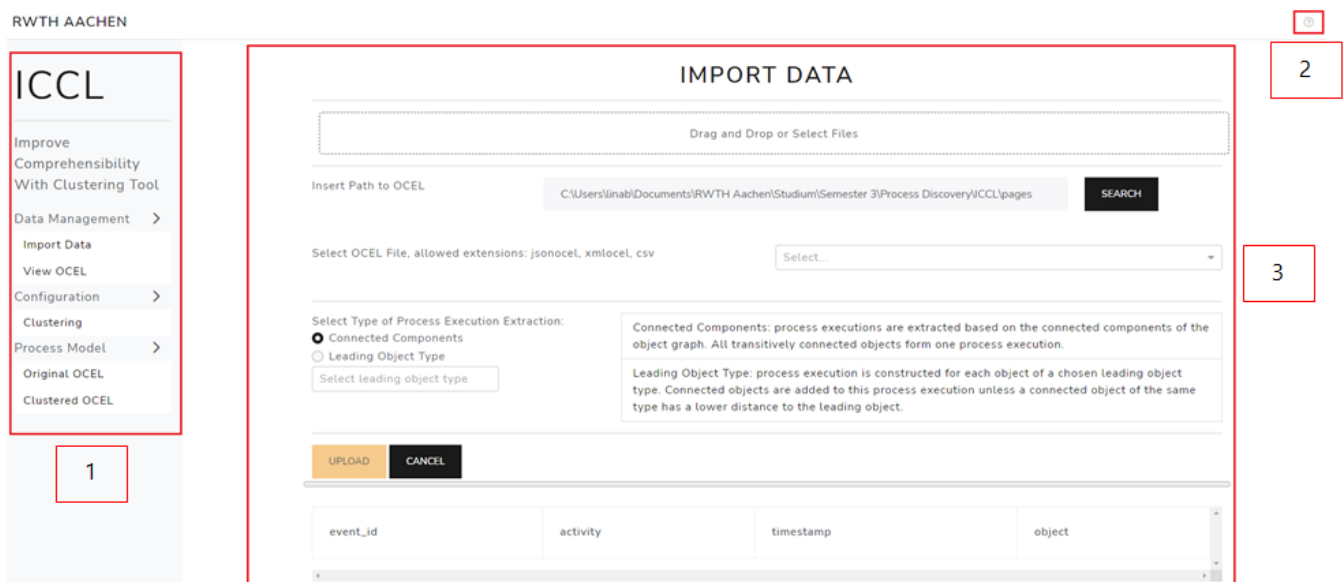
The code comes with a Dockerfile to build and run the application inside a docker container. To build the container run

```
docker build -t docker-iccl .
```

After the container is built, the webapp can be run using

```
docker run -p 8050:8050 docker-iccl
```

## Application Layout



The application layout is structured the following way:

1. The sidebar, to navigate through the application pipeline. The sidebar is visible throughout the whole workflow. You can navigate back-and-forth using these clickable page tabs.
2. The linked help page which provides the manual for this tool. Click on the question mark to access the manual.
3. In the middle the content of the page is located. By default, you start your ICCL journey with the data import page.

# Data Management

## Import OCEL

**IMPORT DATA**

Drag and Drop or Select Files

Insert Path to OCEL

C:\Users\linab\Documents\RWTH Aachen\Studium\Semester 3\Process Discovery\Samples

SEARCH

Path successfully searched.

Select OCEL File, allowed extensions: jsonocel, xmlocel, csv

1b.3

gen\_log.csv

x

PLEASE SPECIFY THE NECESSARY PARAMETERS FOR OCEL CSV IMPORT

Select event id column:

Select timestamp column:

Select start timestamp column (optional):

Select event activity column:

Enter object names:

eid

timestamp

Select...

activity

x
item
x
order

2b PARSE CSV PARAMETERS

3 Select Type of Process Execution Extraction:

☒ Connected Components

☐ Leading Object Type

Select leading object type

Connected Components: process executions are extracted based on the connected components of the object graph. All transitively connected objects form one process execution.

Leading Object Type: process execution is constructed for each object of a chosen leading object type. Connected objects are added to this process execution unless a connected object of the same type has a lower distance to the leading object.

4 UPLOAD

CANCEL

5

eid	activity	timestamp	item	order
0	place_order	2022-01-01T11:00:00.000Z	['i1','i2']	['o1']
1	send_invoice	2022-01-01T11:01:00.000Z	['i1']	

You have two options to import your OCEL:

- 1a: click on the drag and drop field and select a file from the file browser
- 1b.1: insert the path in which the OCEL is located
- 1b.2: click the SEARCH button
- 1b.3: select the OCEL from the list. Note: only files with .csv, .jsonocel or .xmlocel extensions are listed

If you have selected an OCEL with .csv format, the first rows of the .csv are displayed in table 6 and you have to fill out the form 2:

- 2a: specify the OCEL parameters by selecting the column names in the dropdown lists. Select the object types in the last dropdown list.
- 2b: parse the parameters by clicking the PARSE CSV PARAMETERS button

After selecting a file, proceed as follows:

- 3: specify the type of process execution extraction. If leading object type is selected, the possible object types are loaded into the list. Select one of the object types.

- 4: click the UPLOAD button. Note: a progress bar will be displayed, showing you the progress of uploading the OCEL

Optional:

- 5: click the CANCEL button to stop the uploading process.

After successfully uploading the OCEL, the first five rows are displayed in the table 6.

Note: only OCEL files with .csv, .jsonocel or .xmlocel extensions are supported.

View OCEL

ICCL

Improve Comprehensibility With Clustering Tool

Data Management >

Import Data

1

View OCEL

Configuration >

Clustering

Process Model >

Original OCEL

Clustered OCEL

VIEW OCEL

SUMMARY OF OCEL

2

OBJECT TYPES:

["order", "item"]

2.a

ACTIVITY COUNT:

["send\_invoice": 30, "receive\_review": 20, "place\_order": 20, "pick\_item": 20, "delivery\_received": 10, "pay\_order": 10, "send\_delivery": 10, "archive\_order": 10, "cancel\_order": 10, "item\_unavailable": 10]

OBJECT TYPE COUNT:

["order": 20, "item": 30]

NUMBER OF EVENTS:

150

NUMBER OF ACTIVITIES:

10

NUMBER OF OBJECTS:

50

2.b

1

OCEL (LIMITED TO 1000 RECORDS)

event_id	event_activity	event_timestamp	event_start_timestamp	order	item
0	place_order	2022-01-01T12:00:00	2022-01-01T12:00:00	o1	i2i1
1	send_invoice	2022-01-01T12:01:00	2022-01-01T12:01:00		i1
2	send_invoice	2022-01-01T12:02:00	2022-01-01T12:02:00		i2

To view the statistics and gain insight into your OCEL you can visit the 'View OCEL' page. The page can be accessed via the sidebar clicking on 'View OCEL' (1). Inside 2 you can see the summary statistics of your OCEL, providing information such as the object types 2.a or the activity count for each activity 2.b. The table 3 shows the entries of the OCEL, whereas the number of rows is limited to the first 1000 rows.

4 / 15

## Configuration Pipeline

The configuration pipeline to cluster your process executions is as follows:

- 1: Choose the option Clustering from the sidebar
- 2: Select your features within the tab Features
- 3: Choose the graph embedding technique within the tab Embedding
- 4: Choose the clustering technique and parameters within the tab Clustering

Optional:

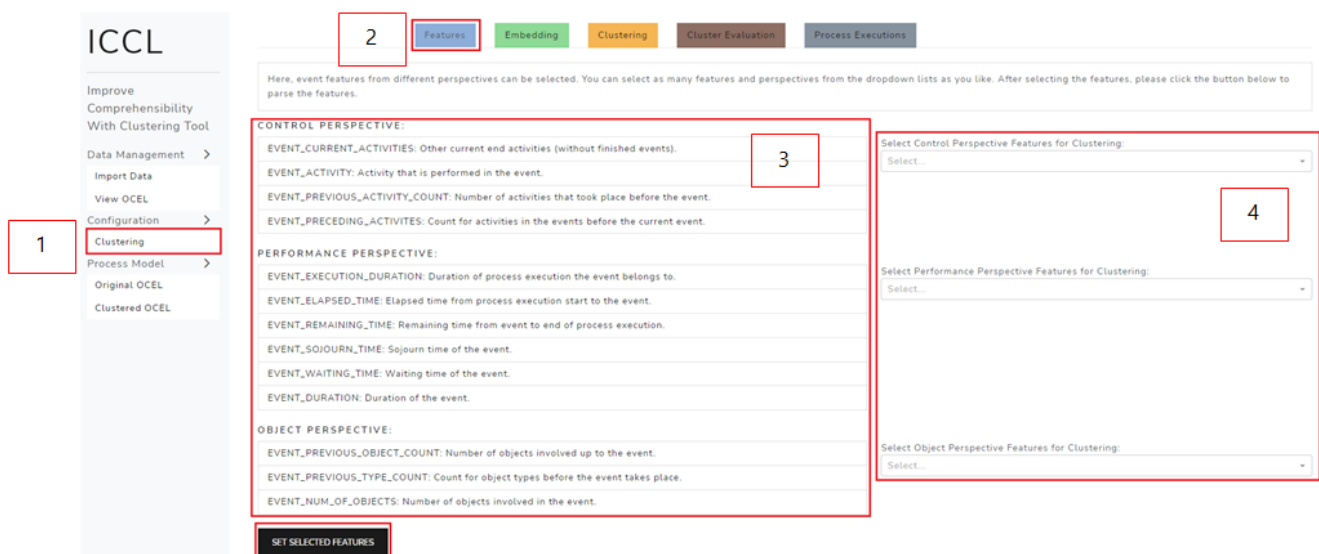
- 5: evaluate the clustering techniques in the Cluster Evaluation tab
- 6: visualize the process execution graphs in the Process Executions tab

### Feature Selection

The starting point of the pipeline is the feature selection. Each process execution is represented as a graph, and each node (=event) can carry different features.

We provide several possible features from different perspectives from which you can choose. Note that you don't have to choose any features at all and can choose as many features as you like.

However, you must click the SET SELECTED FEATURES button so that the chosen selection is parsed.

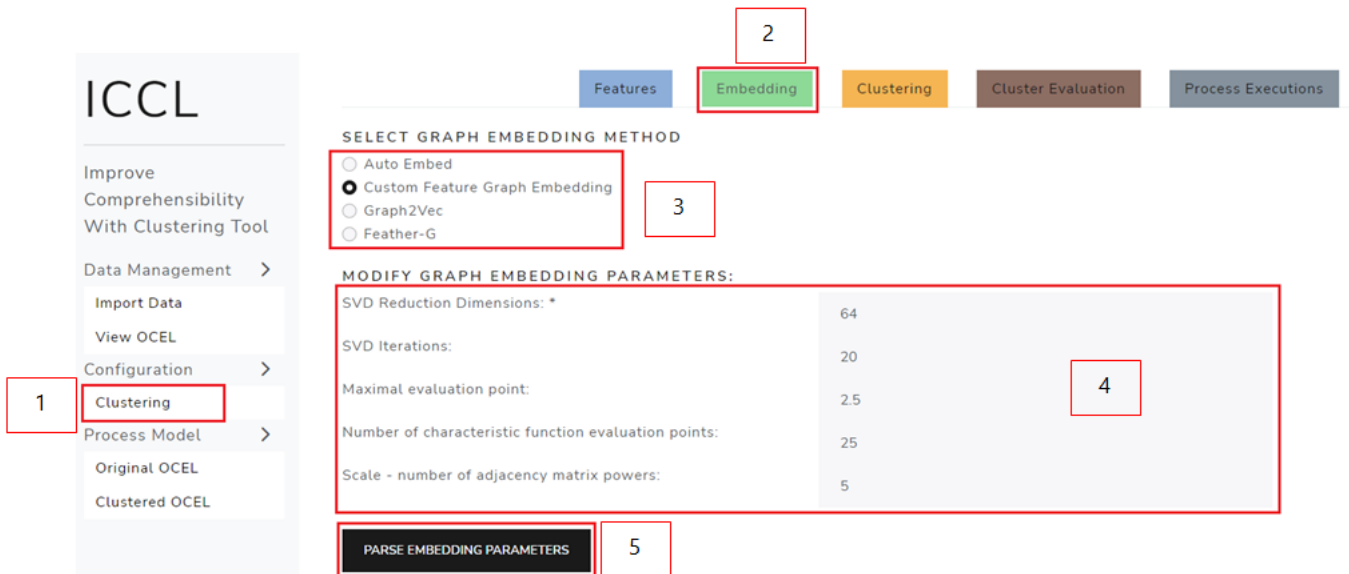


Conduct the feature selection:

1. click on the Clustering page
2. click on the Features tab
3. list of the different features with short explanations, grouped by perspective
4. select the desired features in the dropdown lists - note: you don't have to select features for each perspective
5. click the SET SELECTED FEATURES button

## Graph Embedding

The process execution graphs get enriched with features and are then embedded, so that a clustering algorithm can handle them. In the graph embedding tab, you can specify which kind of embedding technique you want to use. Please see [Graph Embedding Methods](#) for more information about the embedding algorithms.



To specify the graph embedding technique, follow those steps:

1. click on the Clustering page
2. click on the Embedding tab
3. select one of the embedding techniques
4. depending on your selection in 3, a set of adjustable parameters is displayed - you can insert custom parameter values or leave it at the default settings
5. click the PARSE EMBEDDING PARAMETERS button

## Clustering

After the graph embedding is created, you can cluster the process executions. To do so, you have to specify a clustering technique. Furthermore, you can adjust the algorithm's parameters or use the automatic clustering mode. Please see [Clustering Techniques](#) for more information about the clustering algorithms.



Here's how you can execute the clustering:

1. click on the Clustering page
2. click on the Clustering tab
3. choose a clustering technique
4. depending on your selection 3., a set of adjustable parameters is displayed - if you have chosen K-Means or Hierarchical, you can change the number of clusters by manipulating the slider
5. click the PARSE CLUSTERING PARAMETERS button
6. click the START CLUSTERING button

After clicking the START CLUSTERING button, a progress bar is displayed, informing you about how far the process of clustering has progressed. If you want to cancel the process, click the CANCEL button 7.

After the clustering was successfully executed, the following shows up:

Cluster ID	Number of Process Executions
0	7
1	10
2	3

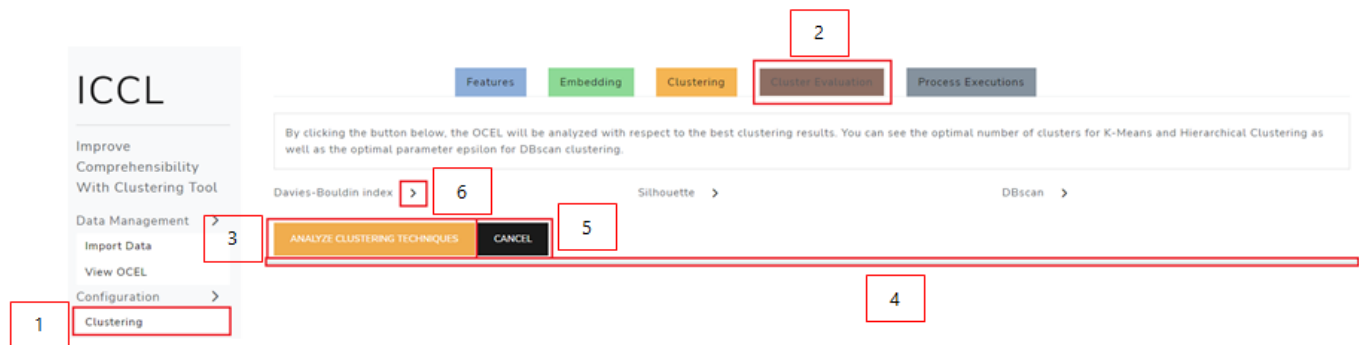
DISCOVER PROCESS MODELS

In the table 1, you can see the resulting clusters and the number of process executions that belong to the cluster ID. If you click the button 2, you get forwarded to the process discovery page.

## Cluster Evaluation

If you want to cluster the process executions using K-Means or Hierarchical Clustering, you might not know which number of clusters is suitable. On the other hand, if you apply DBscan, where you don't have to specify the number of clusters, you must specify the epsilon parameter. To give you more insight into the effect of changing those parameters, you can try the cluster evaluation analysis that is included in ICCL. You can get more information about the scores that are used in [Cluster Evaluation](#).

Note: before launching the cluster evaluation, you have to parse the features and graph embedding parameters.



Execute the cluster evaluation:

1. click on the Clustering page
2. click on the Cluster Evaluation tab
3. click the ANALYZE CLUSTERING TECHNIQUES button
4. look at the progress bar to see how far the process progressed
5. optional: click the CANCEL button to stop the cluster evaluation

To get more information about the scores that are calculated, you can click on the arrows next to the name (6.) After the cluster evaluation was successfully executed, you can see different plots:

- Hierarchical Clustering: silhouette score, Davies-Bouldin index score for average and ward linkage
- K-Means: silhouette score and Davies-Bouldin index score
- DBscan: silhouette score for different epsilon

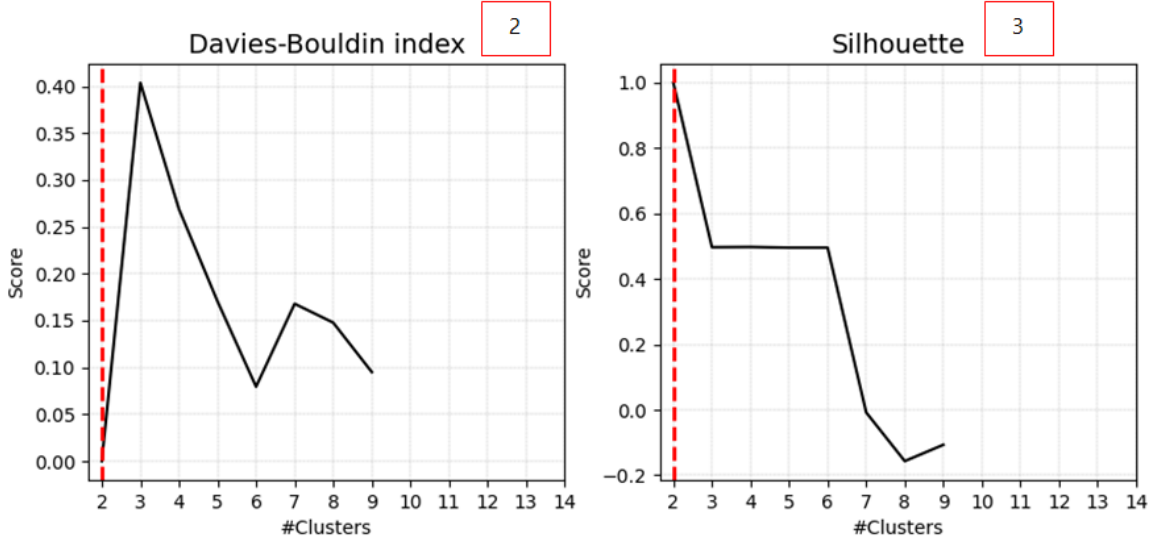


Look at the following plots by way of example:

1

HIERARCHICAL CLUSTERING

Ward Linkage

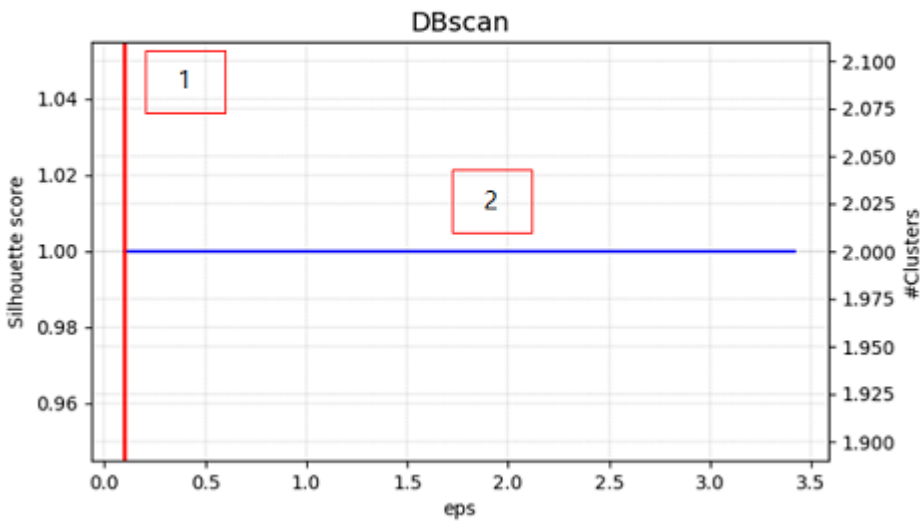


- 1. shows the type of clustering algorithm that was used
- 2. result for Davies-Bouldin index for different #clusters
- 3. result for silhouette score for different #clusters

The red dashed line indicates the best number of cluster, i.e. the number of clusters that eventuated in the optimal score.

The resulting plot for DBscan evaluation looks a bit different:

DBSCAN CLUSTERING



- 1. indicates the optimal value for epsilon
- 2. shows the resulting number of clusters

## Process Executions

One of the key aspects in ICCL is the process execution. If you want to gain more insights into the process executions of your OCEL, you can do so on the Process Executions page.

Note: you can visualize the process executions and inspect its features at any point after uploading the OCEL.

To visualize the graph of one process execution and see its features, perform the following steps:

The screenshot shows the ICCL interface. On the left, a sidebar contains 'Clustering' (1). The top navigation bar has 'Process Executions' (2). Below the navigation bar, 'Number of Process Executions: 20' (7) is displayed. A 'START PROCESS EXECUTION FEATURE EXTRACTION' button (3) is visible. To its right, a 'Select Process Execution' dropdown shows '19' (4). Below these, a table (5) lists features: Number of Events (4), Number of Ending Events (1), Throughput Duration (180), Number of Objects (2), Unique Activities (4), Number of Starting Events (1), and Duration of Last Event (60). At the bottom, a process execution graph (6) is shown with nodes 'place\_order', 'cancel\_order', 'send\_invoice', and 'send\_goods'.

1. click on the Clustering page
2. click on the Process Executions tab
3. click on the START PROCESS EXECUTION FEATURE EXTRACTION button
4. select a process execution in the dropdown list

The results are the following:

5. feature table with the features that were extracted from the process execution
6. graph of the process execution:

- to move the position of a node, click on the node and pull it to the desired location
- to move the whole graph, click on an edge and pull it to the desired location
- to zoom in and out, locate the cursor on the graph and use the scroll wheel of your mouse

At any point after uploading the OCEL, you can see the number of process executions in 7.

# Process Discovery

## Discover Process Model

You can discover the Process Model for the original OCEL as well as for the sub-logs that are created by clustering the OCEL. Before discovering the latter Process Models, you have to go through the ICCL pipeline, so that clusters are already created. You can discover the comprehensive Process Model at any point after uploading the OCEL.

### Comprehensive Process Model

ICCL

Improve Comprehensibility With Clustering Tool

Data Management >

Import Data

View OCEL

Configuration >

Clustering

Process Model >

Original OCEL

Clustered OCEL

PROCESS DISCOVERY ORIGINAL OCEL

Select Conformance Measure: Fitness x

CALCULATE

Result:

2 START PROCESS DISCOVERY

Conducting the process discovery:

1. click on the 'Original OCEL' page
2. click the 'START PROCESS DISCOVERY' button

After the process discovery is performed, the resulting model is displayed below. You can manipulate the view of the process model:

- to move the view of the model, click on the model and pull it to the desired location
- to zoom in and out, locate the cursor on the model and use the scroll wheel of your mouse

### Cluster Process Models

ICCL

Improve Comprehensibility With Clustering Tool

Data Management >

Import Data

View OCEL

Configuration >

Clustering

Process Model >

Original OCEL

Clustered OCEL

PROCESS DISCOVERY CLUSTERED OCEL

Select Conformance Measure: Fitness x

CALCULATE

Result:

2 START PROCESS DISCOVERY

Process Models successfully discovered.

Select Cluster: 0 x

3

PROCESS EXECUTION FEATURES AVERAGE:

Number of Events	Number of Ending Events	Throughput Duration	Number of Objects	Unique Activities	Number of Starting Events	Duration of Last Event
11	3	600	3	8	1	60

4

### Conducting the process discovery:

1. click on the 'Clustered OCEL' page
2. click the 'START PROCESS DISCOVERY' button
3. select the process model of a specific cluster by selecting the cluster's number in the dropdown list
4. the table displays the average feature values of the process executions that belong to the model

You can manipulate the view of the process model:

- to move the view of the model, click on the model and pull it to the desired location
- to zoom in and out, locate the cursor on the model and use the scroll wheel of your mouse

### Conformance Checking

To compare the process model with the OCEL, you can apply two conformance checking methods: fitness and precision. Before checking the conformance of the clustered Process Models, you have to go through the ICCL pipeline, so that clusters are already created. You can check the conformance of the comprehensive Process Model at any point after uploading the OCEL.

1. select a measure: fitness or precision
2. click the CALCULATE button
3. result of the measure (in case of conformance checking on the clustered OCEL, a table is displayed, showing the result for each cluster)

### Export Process Model

After discovering a process model for the original OCEL or for the clustered sub-logs (and a cluster is selected), this segment shows up:

To save the displayed process model, follow these steps:

1. insert the path to the location at which the image should be saved
2. enter a name for the image
3. select the file extension: png or svg
4. click the SAVE button

# Algorithms

Here you can find a short description of the algorithms that are used in ICCL.

## Process Execution Extraction

A process execution is a set of events of connected objects and resembles the case notion in traditional Process Mining. You can choose between two different types of extraction methods:

- **Connected Components:** the technique uses the object graphs and extracts a process execution based on all transitively connected objects. This way, complex event logs might lead to large process executions.
- **Leading Object Type:** after choosing one object type as the leading type, a process execution is constructed by looking at each object of the object type. Objects that are connected to the leading object are added to the process execution unless another object of that type has a lower distance. As a result, the process executions are more limited in size.

## Graph Embedding Methods

- **Auto Embed:** uses Custom Feature Graph Embeddings but optimizes the number of dimensions automatically. The method tries to reduce the number of dimensions as far as possible but ensures that there are still enough dimensions to fully express the complexity of the feature graphs. This is done by comparing the normalized embedding loss (<https://doi.org/10.1038/s41467-021-23795-5>) for different numbers of dimensions. The fewer dimensions, the more efficient the clustering can be performed. The final choice of dimensions is displayed to the user. If no features were selected, graph2vec is automatically used, since our Custom Feature Graph Embedding method requires at least one feature.
- **Custom Feature Graph Embedding:** especially designed for ICCL. The algorithm first creates a node embedding via FeatherNode which uses characteristic functions of node features with random walk weights to describe node neighborhoods. In the second step, the node embeddings are averaged over each dimension, resulting in a vectorized embedding of the graph. Focusses **features**, the structure of the process execution graphs is only implicitly considered. Check out the [karateclub documentation](#) for more information about the FeatherNode parameters that can be configured.

Note: The algorithm requires that at least one feature was selected.

- **Graph2Vec:** first identifies subgraphs surrounding each node in the feature graphs. By means of the Weisfeiler-Lehman's algorithm, the subgraphs are considered as the vocabulary for a doc2vec SkipGram model. Since the graph's structure is captured within the algorithm, feature graphs that are similar in structure will be close in the embedding space. Focusses the **graph structure** and additionally allows one feature per node. Check out the [karateclub documentation](#) for more information about the parameters that can be configured.
- **Feather-G:** uses characteristic functions defined on graph vertices to describe the distribution of vertex attributes. Feather-G extracts node-level features that are pooled and then used to create a description of the feature graph. Focusses only the **graph structure**. Check out the [karateclub documentation](#) for more information about the parameters that can be configured.

## Clustering Techniques

ICCL makes use of the `sklearn.cluster` module to apply different clustering algorithms to the embedding. Please refer to the [sklearn.cluster documentation](#) to gather more information about the clustering parameters that can be configured in ICCL.

- **Auto Cluster:** This technique tries to find the best clustering technique automatically and optimizes their hyperparameters. Currently, we compare MeanShift, DBscan and KMeans with different parameter settings and our goal is to find a setting which maximizes the silhouette score. After the clustering has been performed, the final selected method as well as its parameters are displayed to the user.
- **K-Means:** clusters data by trying to separate samples in  $k$  groups of equal variance. It aims to choose centroids that minimize the within-cluster sum-of-squares. The number of clusters ( $k$ ) has to be specified beforehand.
- **Hierarchical Clustering:** builds nested clusters by merging or splitting them successively. This hierarchy of clusters is represented as a tree, whereas the root is the unique cluster gathering all samples and the leaves are clusters containing only one sample. In ICCL, you have to specify the number of clusters beforehand.
- **Mean-Shift:** aims to discover blobs in a smooth density of samples. It is a centroid based algorithm, which works by updating candidates for centroids to be the mean of the points within a given region.
- **Affinity-Propagation:** Affinity-Propagation creates clusters by sending messages between pairs of samples until convergence. The algorithm finds members of the data points that are representatives of the clusters.
- **DBscan:** Density-Based Spatial Clustering of Applications with Noise is an clustering approach that finds core samples of high density and expands clusters from them. The parameter epsilon is specifying the radius of a neighborhood with respect to some point, in which the number of neighboring points is counted.

## Evaluation Measures

### Cluster Evaluation Measures

- **Silhouette Score:** quantifies the space between different clusters. For each number of clusters, a clustering algorithm is performed. Afterwards, it is measured how similar the observation are to the assigned cluster and how dissimilar they are to the observation of the nearest cluster. The measure has the range  $[-1, +1]$ , whereas a score near  $+1$  indicates that the clusters are well separated and negative scores indicate that the samples might be wrongly separated. Generally, the silhouette score is calculated for each datapoint and then averaged over the whole dataspace. You can find the calculation steps here: [click](#)
- **Davies-Bouldin Index:** measure of the ratio between within-cluster distances, and between cluster distances. The score is bounded between  $[0, 1]$ . The lower the value, the tighter the clusters and the separation between clusters. The steps of calculation can be found here: [click](#)

## Model Evaluation Measures

- **Fitness:** measures to what extend the observed traces can be replayed by the model.
- **Precision:** measures the extend of behaviour that is not captured in the event log, but allowed in the model.

## Model Discovery

The discovery of an object-centric Petri net works as follows: for each object type, a flattened event log is created and subsequently an accepting Petri net is discovered. For the discovery part, ICCL supports the Inductive Miner. In the second step, the Petri nets are merged into one Petri net. The object types are assigned to places and the variable arcs are identified.

**Inductive Miner:** the algorithm recursively performs the following steps

- build a directly follows graph from the log
- find a cut and return the cut-operator and cut-partition
- split the log into sublogs