

Process Mining

Generalized Alignment-Based Trace Clustering

Mathilde Boltenhagen¹

J.Carmona²

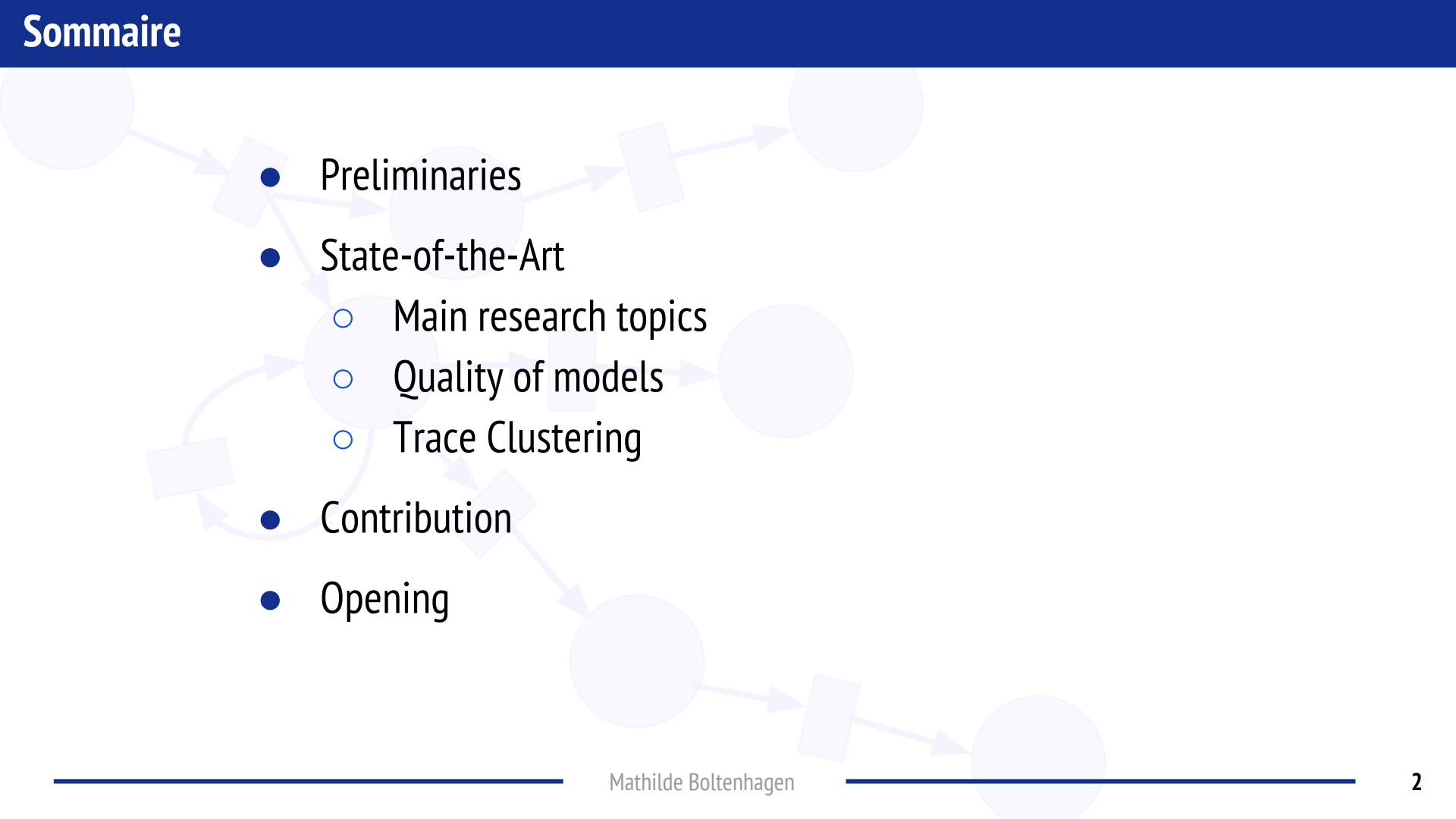
T.Chatain²

CIRM - 5 mars 2019

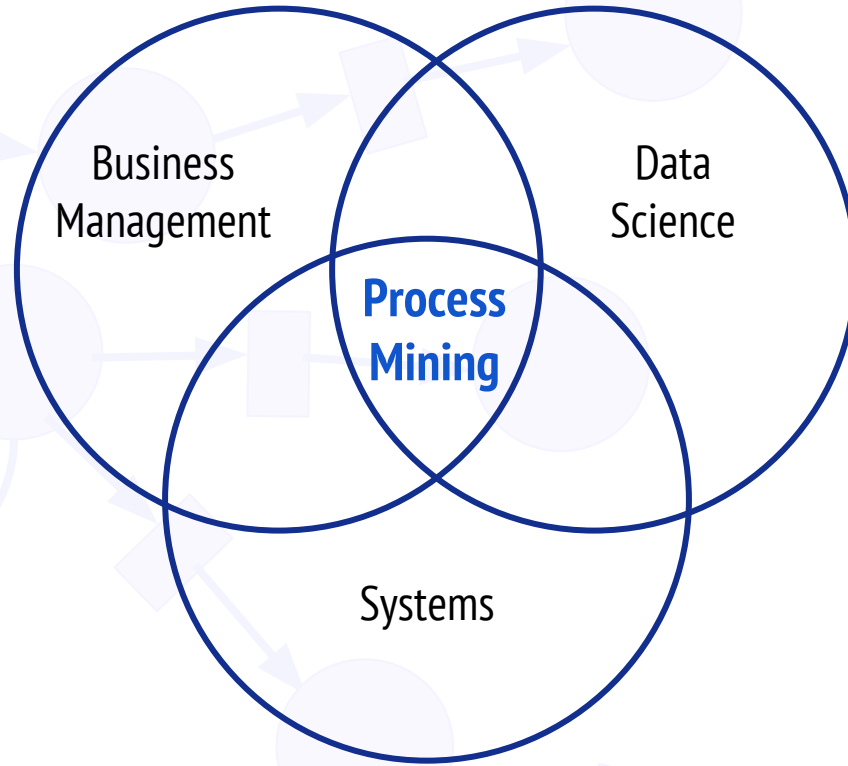
¹LSV, CNRS, ENS Paris-Saclay, Inria, Université Paris-Saclay

²Universitat Politècnica de Catalunya



- 
- Preliminaries
 - State-of-the-Art
 - Main research topics
 - Quality of models
 - Trace Clustering
 - Contribution
 - Opening

Position of the field



Example

a patient arrive at the emergencies e1

first aid is done e2

the patient is urgently operated e3

What are the behaviors?

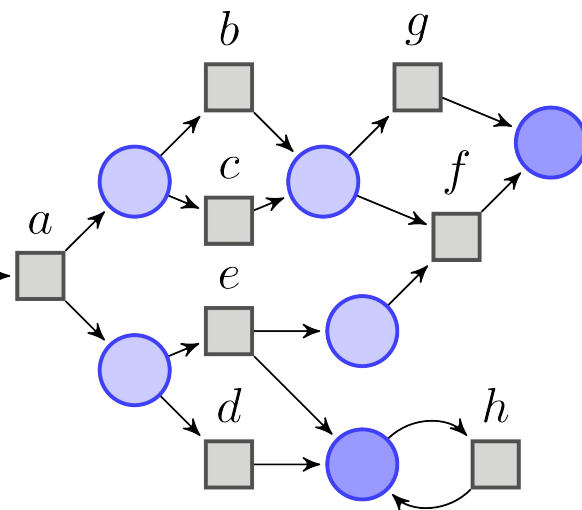
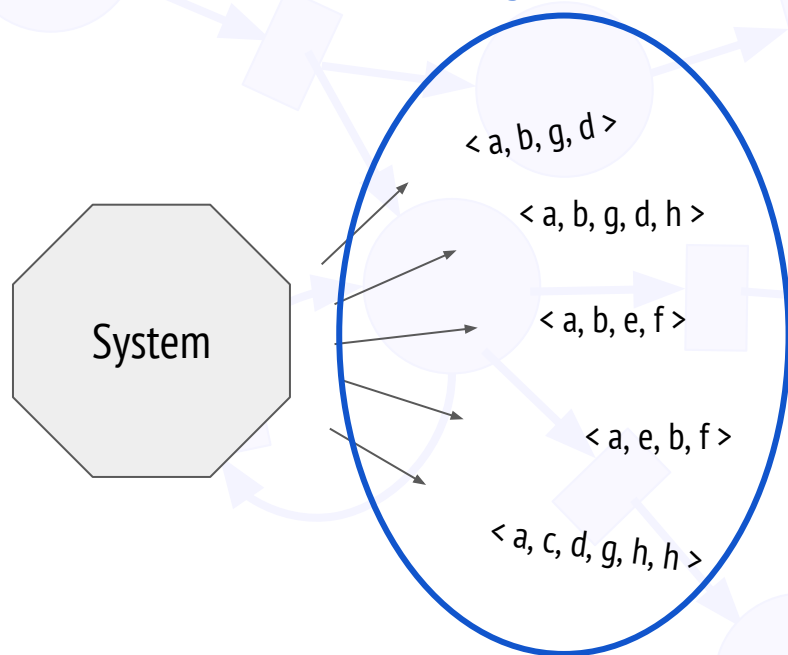
the patient must be moved to the department y

the doctor fill a form about the patient e4

$\langle e1, e2, e3, e4 \rangle$

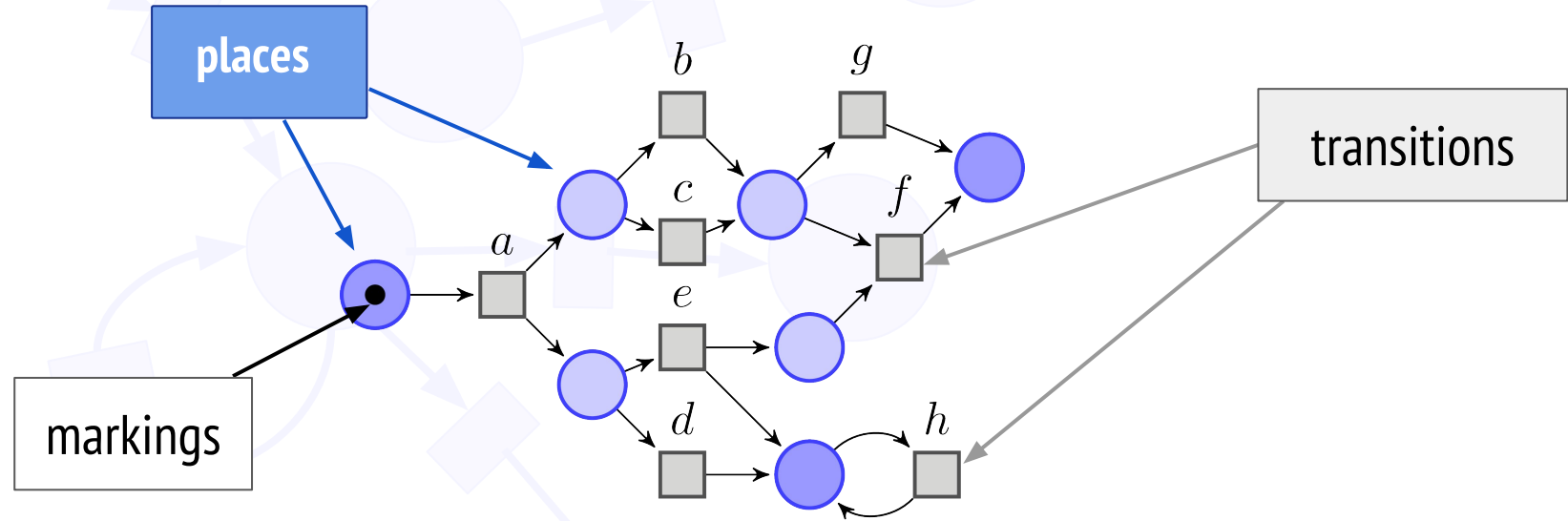
Hospital

log traces



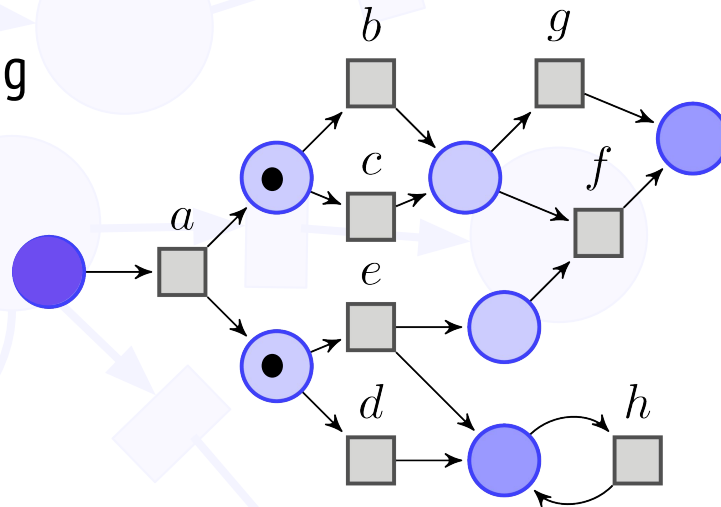
Process Model

Petri Nets



Token game

- > Firing transition
- > Reachable marking





> Process model discovery

> Conformance checking

> Model repair

How good is the discovered model ?

> Fitness

*Has the recorded behavior been modelled ?**

> Precision

*Has the modelled behavior been recorded ?**

> Generalization

Will the model fit future behaviors ?

> Simplicity

*How complex is the model ?**

[Buijs 2012] ; [Carmona 2018]

Log traces

< a, b, g, d >
< a, b, g, d, h >
< a, b, e, f >
< a, e, b, f >
< a, c, d, g, h, h >
< a, c, d, g, h, h, h, h, h >
< a, f, f, f >

Clusters

< a, b, g, d >
< a, b, g, d, h >

< a, b, e, f >
< a, e, b, f >
< a, f, f, f >

< a, c, d, g, h, h >
< a, c, d, g, h, h, h, h, h >

- > Data analysis
- > Model discovery

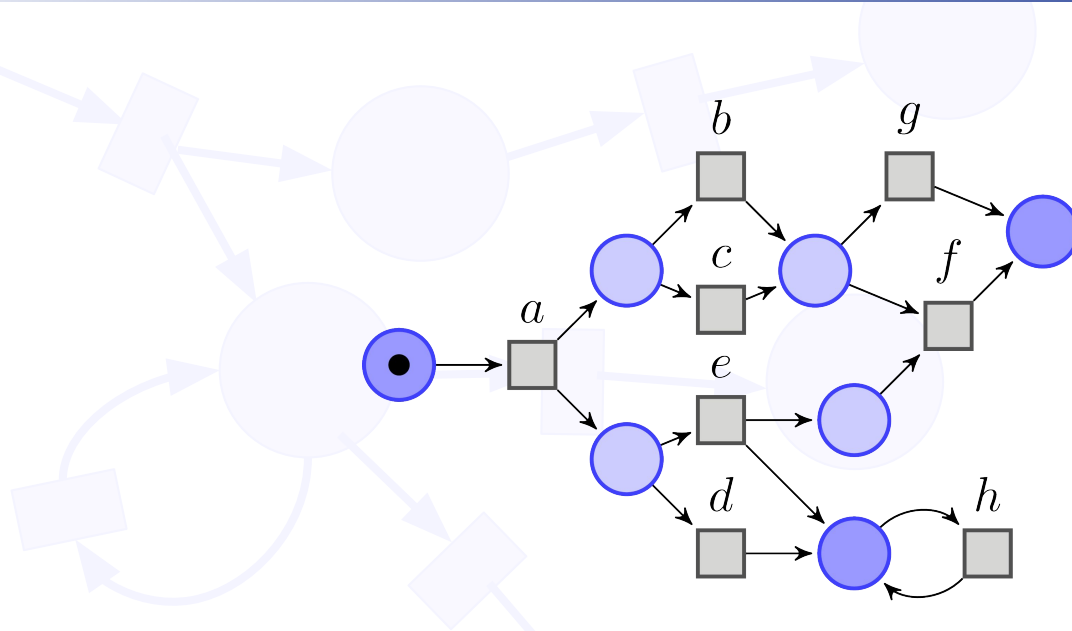
[Greco 2006] ; [Ferreira 2007]



Idea : to cluster data with the existing model

- > highlight parts of models that are executed
- > show deviated traces
- > model repairs

Full Runs



Example of full run : $\langle a, b, d, g, h \rangle$

Full Runs as Centroids

Centroids (full runs)	Traces	Distances
< a, b, g, d >	< a, b, g, d >	0
	< a, b, g, d, h >	1
< a, b, e, f >	< a, b, e, f >	0
< a, e, b, f >	< a, e, b, f >	0
< a, c, d, g, h, h >	< a, c, d, g, h, h >	0
< a, c, d, g, h, h, h, h >	< a, c, d, g, h, h, h, h >	0
Non-clustered	< a, f, f, f >	

[Chatain 2017]

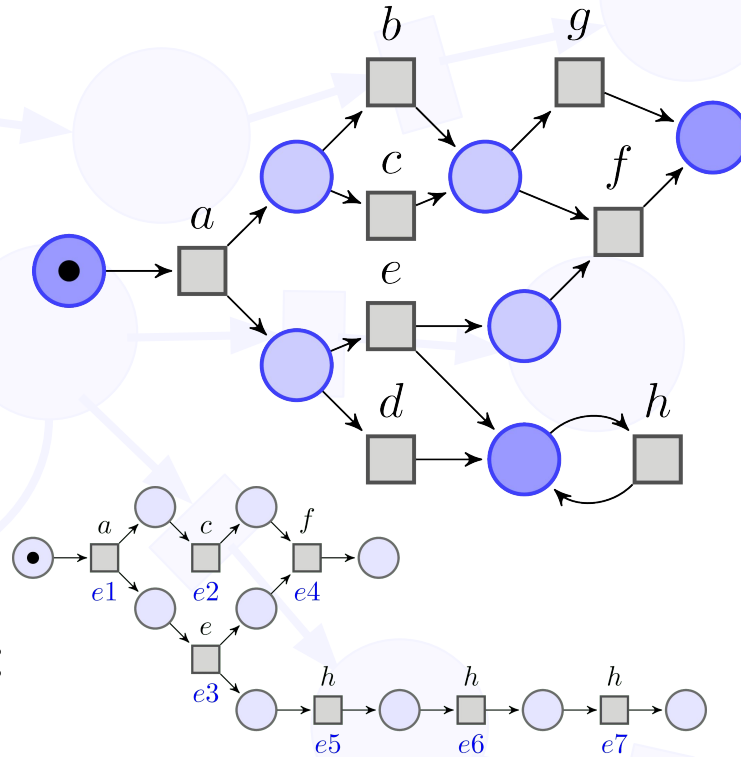
Limits

Centroids (full runs)	Traces	Distances
< a, b, g, d >	< a, b, g, d >	0
	< a, b, g, d, h >	1
< a, b, e, f >	< a, b, e, f >	0
< a, e, b, f >	< a, e, b, f >	0
< a, c, d, g, h, h >	< a, c, d, g, h, h >	0
< a, c, d, g, h, h, h, h >	< a, c, d, g, h, h, h, h >	0
Non-clustered	< a, f, f, f >	

concurrency

loops

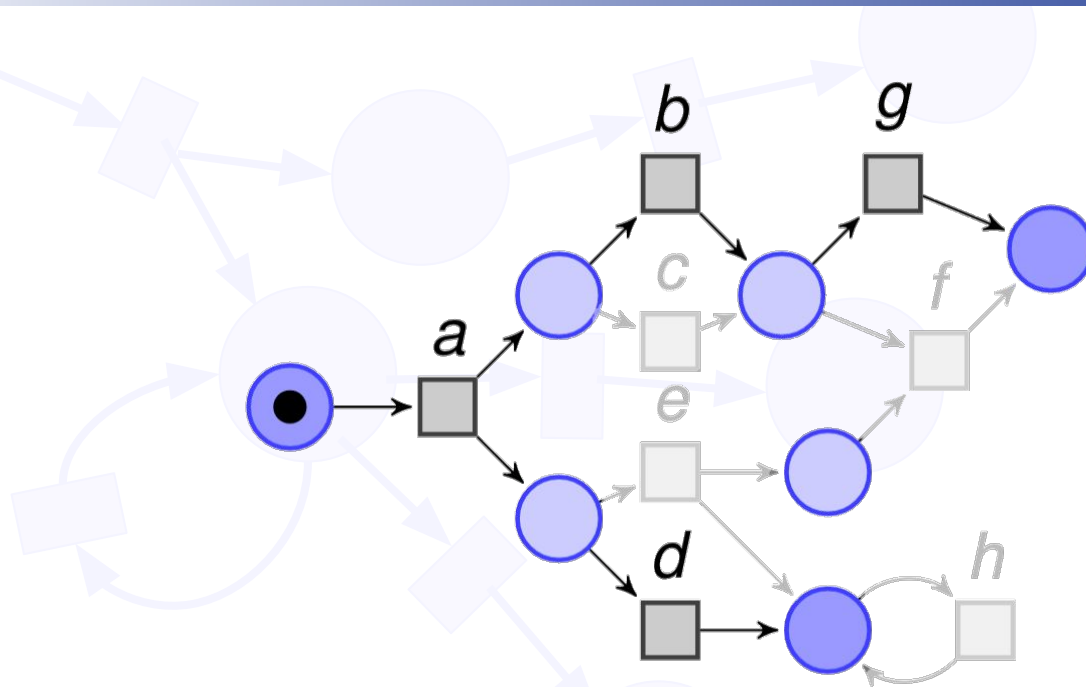
[Chatain 2017]



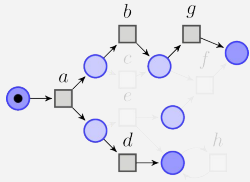
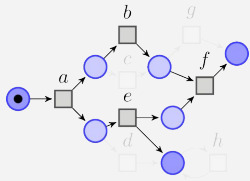
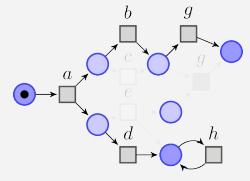
Example of process :

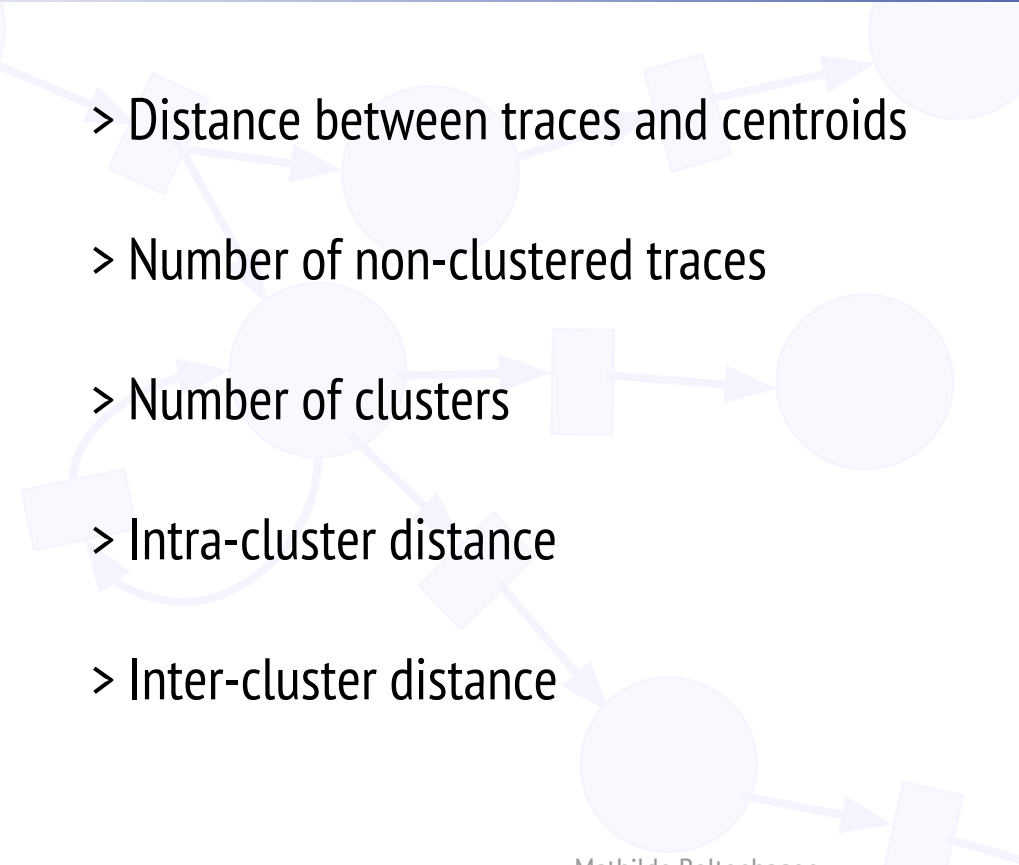
Processes as Centroids

Centroids (process)	Traces	Distances
	$\langle a, b, g, d \rangle$	0
	$\langle a, b, g, d, h \rangle$	1
	$\langle a, b, e, f \rangle$	0
	$\langle a, e, b, f \rangle$	0
	$\langle a, c, d, g, h, h \rangle$	0
	$\langle a, c, d, g, h, h, h, h \rangle$	0
Non-clustered	$\langle a, f, f, f \rangle$	



Subnets as Centroids

Centroids (subnets)	Traces	Distances
	$\langle a, b, g, d \rangle$	0
	$\langle a, b, g, d, h \rangle$	1
	$\langle a, b, e, f \rangle$	0
	$\langle a, e, b, f \rangle$	0
	$\langle a, c, d, g, h, h \rangle$	0
	$\langle a, c, d, g, h, h, h, h \rangle$	0
Non-clustered	$\langle a, f, f, f \rangle$	

- 
- > Distance between traces and centroids
 - > Number of non-clustered traces
 - > Number of clusters
 - > Intra-cluster distance
 - > Inter-cluster distance



> Tool *DARK SIDER**

> Optimal clusterings

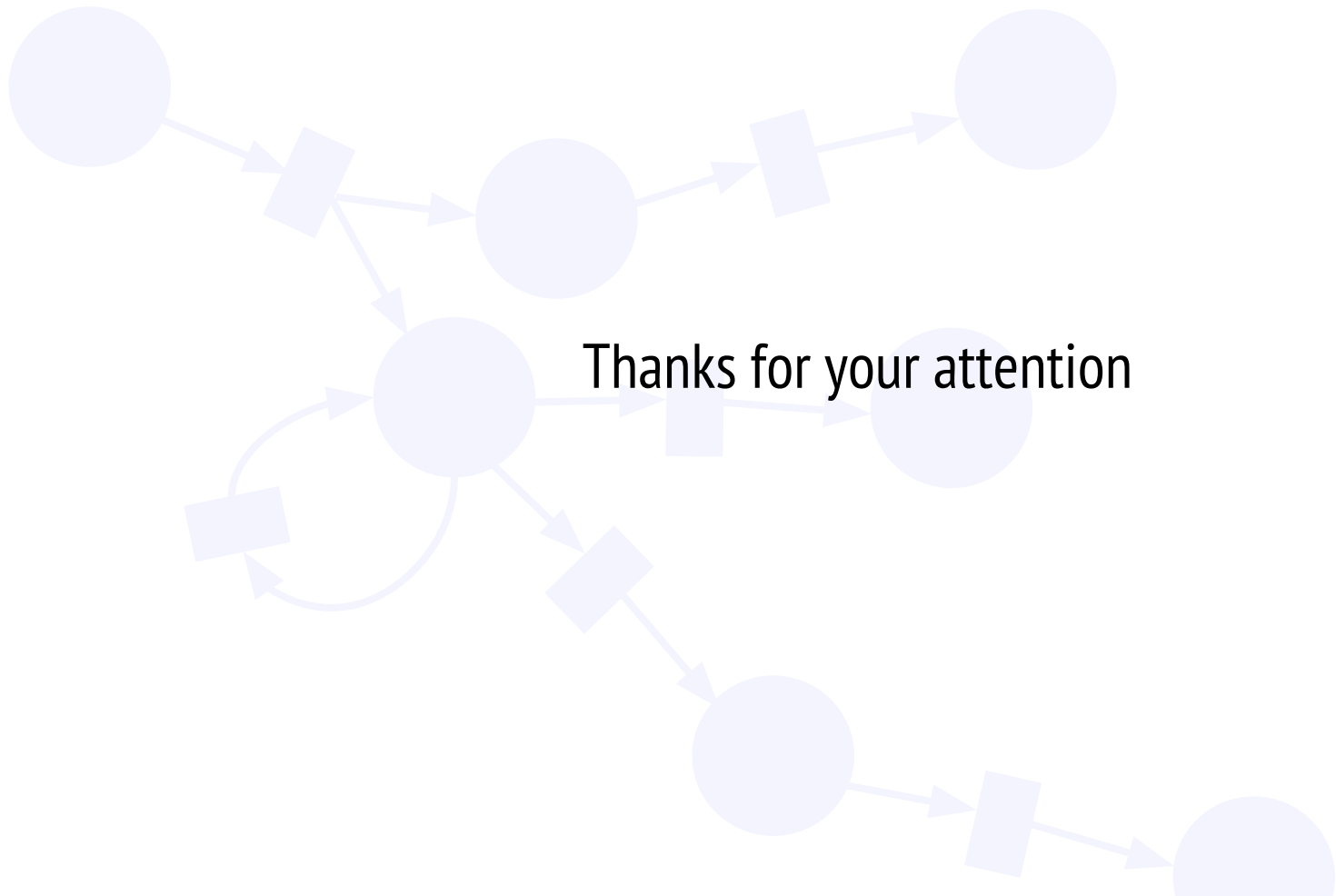
*<https://github.com/BoltMaud/darksider>

At least a better fitness

Not too far from the original model

How can one repair an existing process model ?

....



Thanks for your attention