

Process Mining

Generalized Alignment-Based Trace Clustering for Process Behavior



Mathilde Boltenhagen¹, Thomas Chatain¹, Josep Carmona²

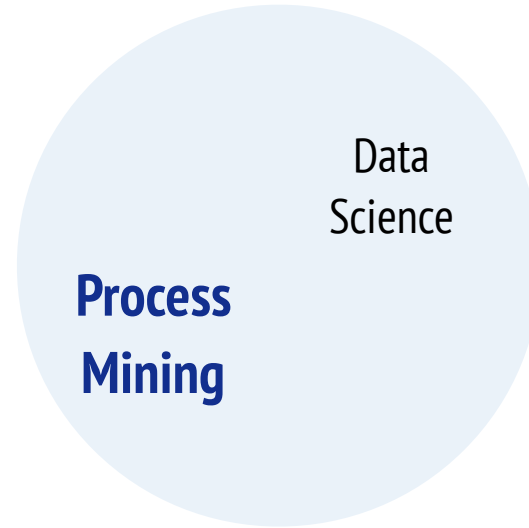


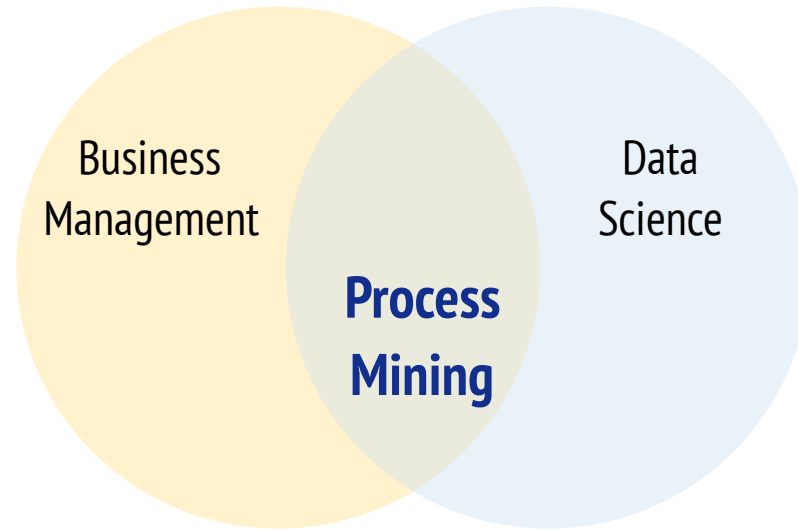
¹LSV, CNRS, ENS Paris-Saclay, Inria, Université Paris-Saclay

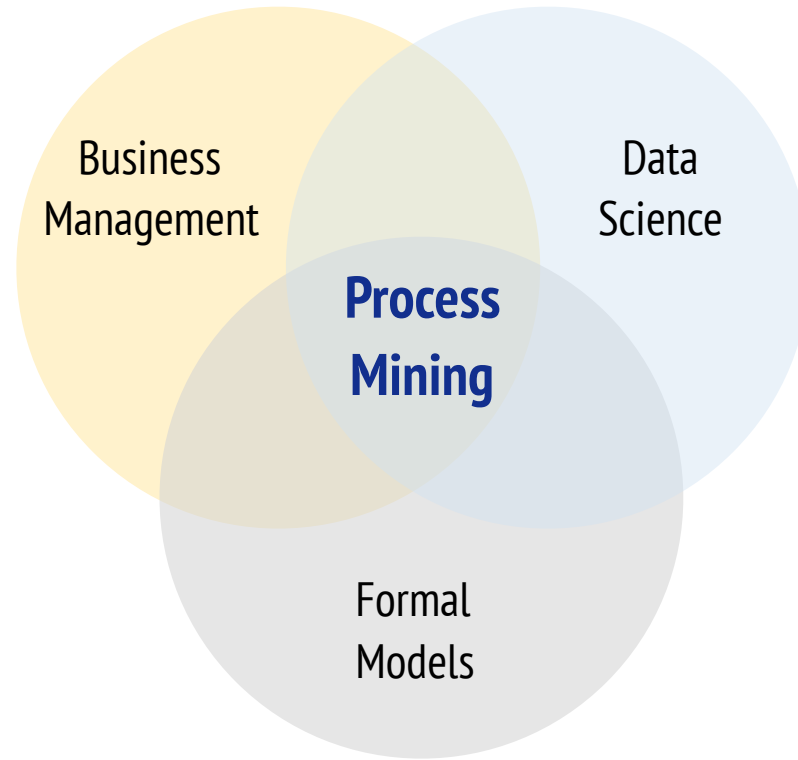
²Universitat Politècnica de Catalunya



**Process
Mining**











Start the form





Start the form

S



Bad rating

B

Good rating

G



Start the form

S



Bad rating

B



Good rating

G

Get apologies

A



Start the form

S



Bad rating

B



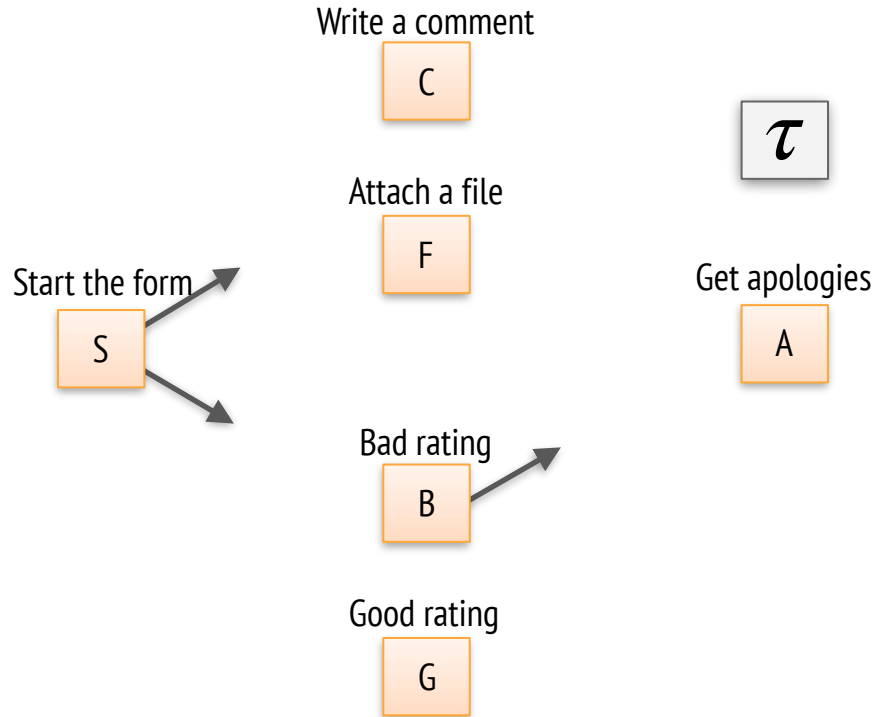
Good rating

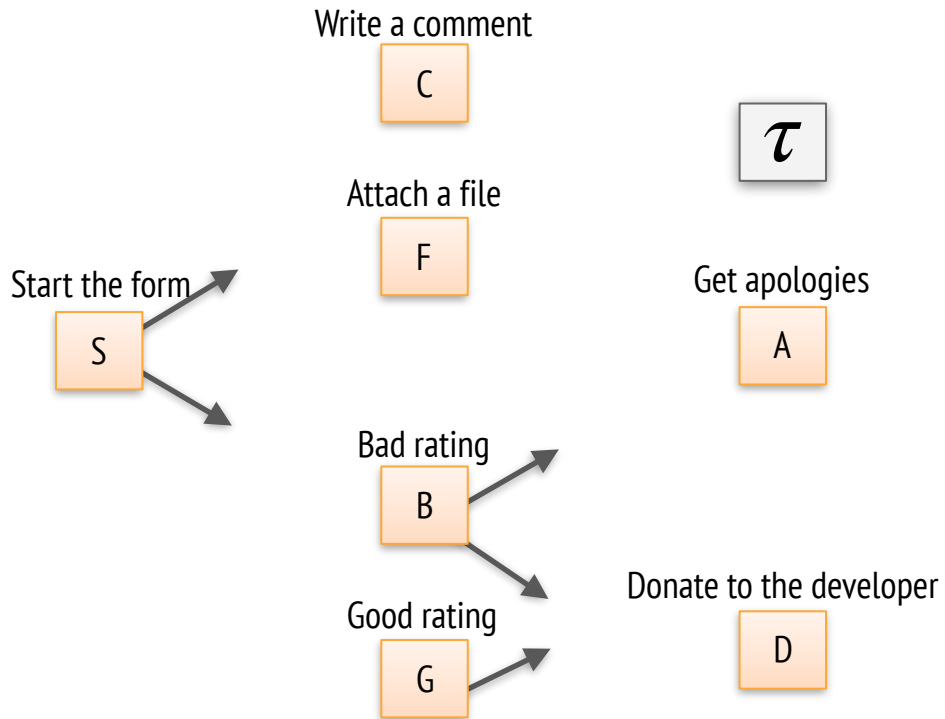
G

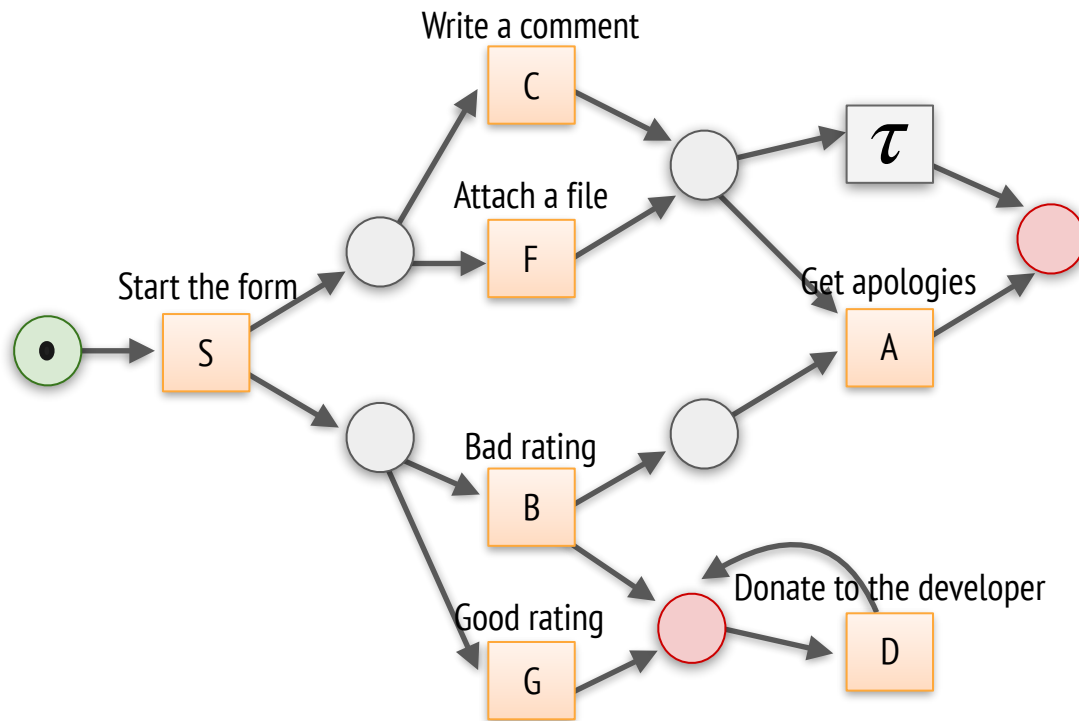
τ

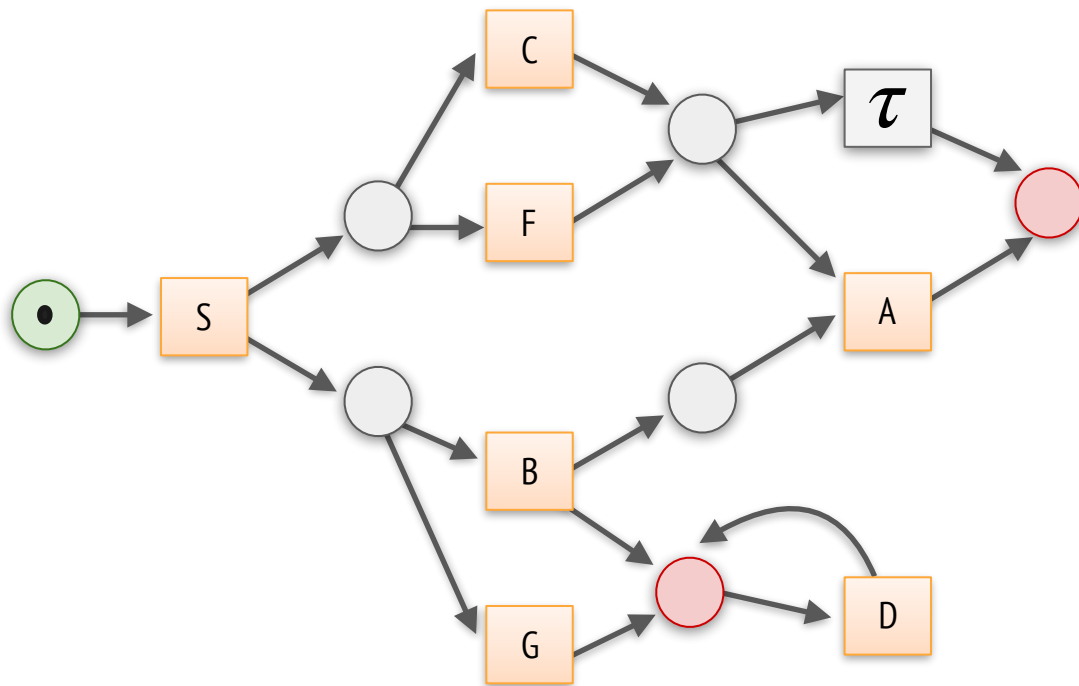
Get apologies

A









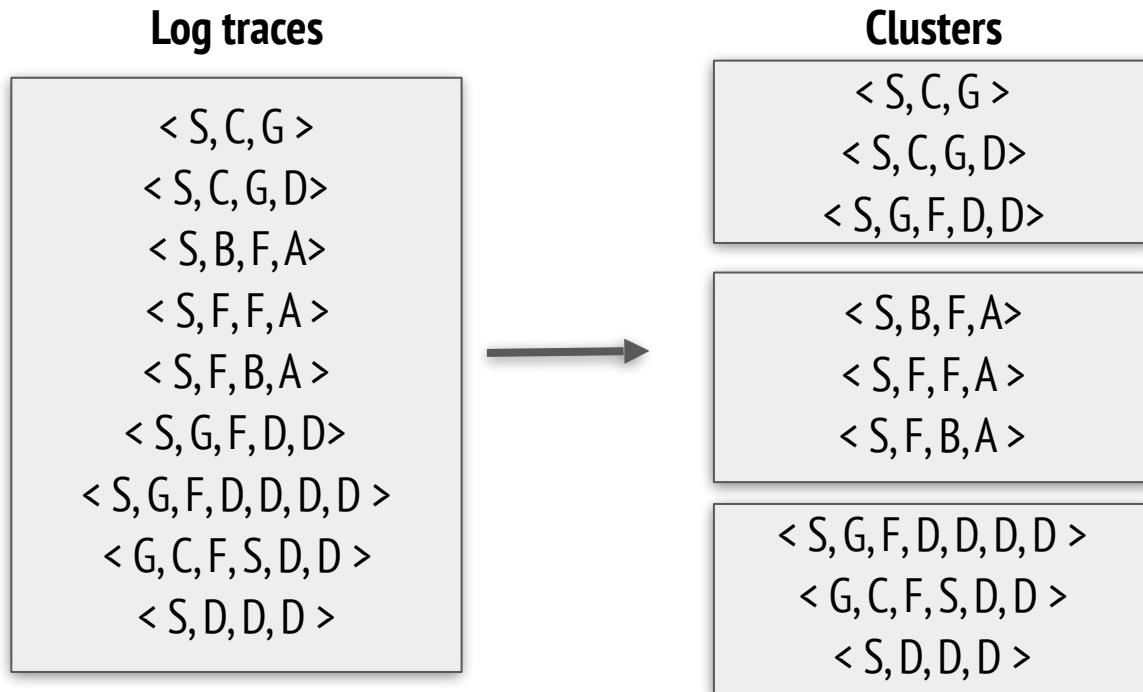
Log traces

< S, C, G >
< S, C, G, D >
< S, B, F, A >
< S, F, F, A >
< S, F, B, A >
< S, G, F, D, D >
< S, G, F, D, D, D, D >
< G, C, F, S, D, D >
< S, D, D, D >

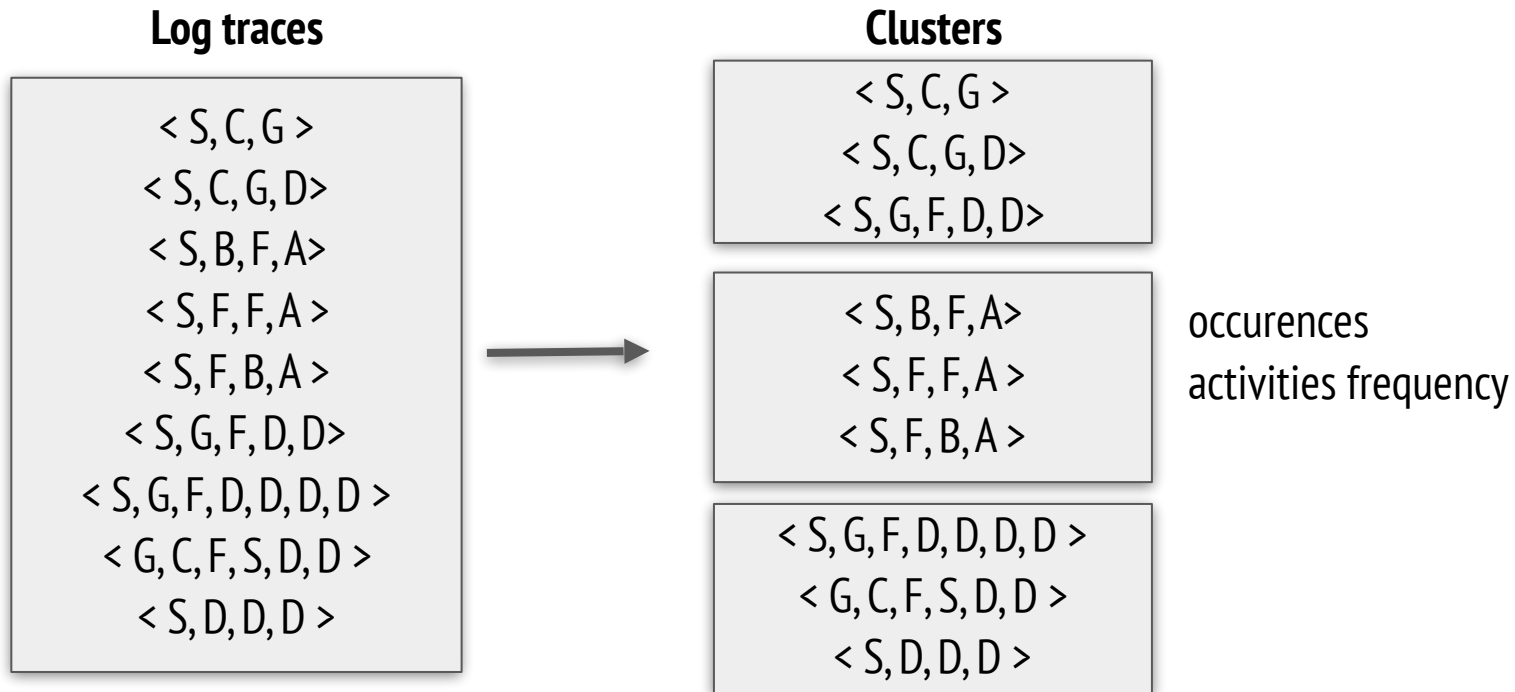
Log traces

< S, C, G >
< S, C, G, D >
< S, B, F, A >
< S, F, F, A >
< S, F, B, A >
< S, G, F, D, D >
< S, G, F, D, D, D, D >
< G, C, F, S, D, D >
< S, D, D, D >

Data clustering is the task of grouping objects by similarity.



[Greco et al. 2006] ; [Ferreira et al. 2007] ; [Bose et al. 2009]



[Greco et al. 2006] ; [Ferreira et al. 2007] ; [Bose et al. 2009]

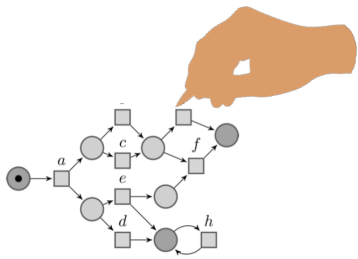
New idea : to cluster data based on an existing process model

- > highlight parts of models that are executed
- > show deviating traces
- > model repair

Existing models

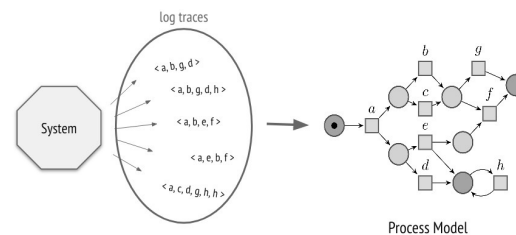
> Available process model

- system modelled by hand
- model used at design phase



> Discovered process model

- data are used to create the process model : process discovery (Alpha Miner, ILP miner, Inductive miner..)



Log traces

< S, C, G >

< S, C, G, D >

< S, B, F, A >

< S, F, F, A >

< S, F, B, A >

< S, G, F, D, D >

< S, G, F, D, D, D, D >

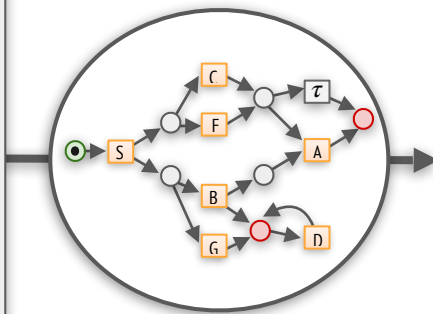
< G, C, F, S, D, D >

< S, D, D, D >

[Chatain et al. 2017]

Log traces

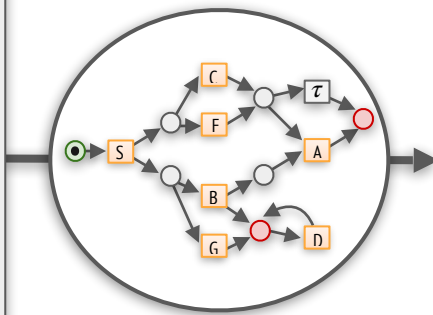
< S, C, G >
< S, C, G, D >
< S, B, F, A >
< S, F, F, A >
< S, F, B, A >
< S, G, F, D, D >
< S, G, F, D, D, D, D >
< G, C, F, S, D, D >
< S, D, D, D >



[Chatain et al. 2017]

Log traces

< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >



Clusters

< S, C, G >
 < S, C, G, D >

< S, B, F, A >
 < S, F, F, A >

< S, F, B, A >

< S, G, F, D, D >

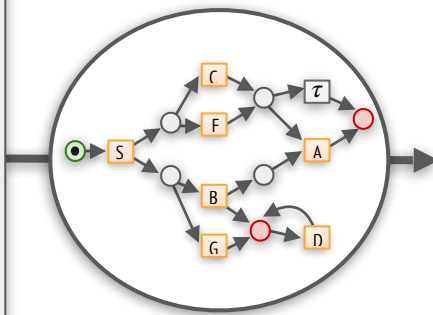
< S, G, F, D, D, D, D >

< G, C, F, S, D, D >
 < S, D, D, D >

[Chatain et al. 2017]

Log traces

< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >

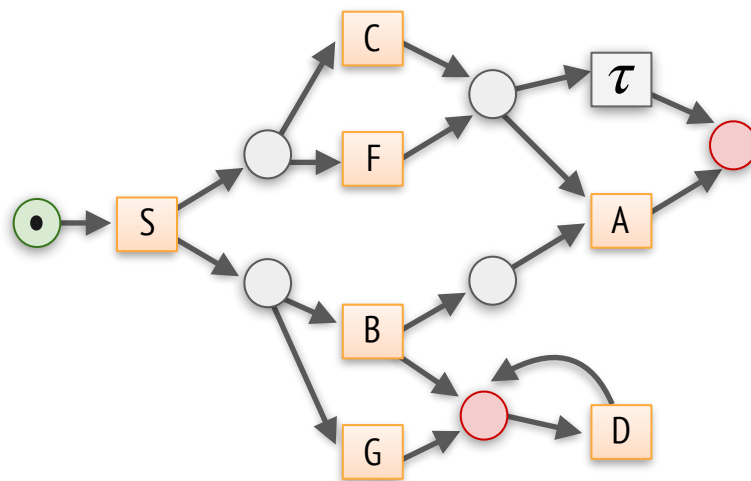


Clusters



[Chatain et al. 2017]

Full runs as centroids



Example of full run : $\langle S, C, \tau, G \rangle$

Full runs as centroids

\mathcal{N} a process model, \mathcal{L} a log :

Find $\mathcal{C} = (\{u_1 \dots u_n\}, \chi)$:

$\{u_1 \dots u_n\} \in \text{Runs}(\mathcal{N})$

$\chi : \mathcal{L} \rightarrow \{nc, u_1 \dots u_n\}$

Such that for every $\sigma \in \mathcal{L}$, $\text{dist}(\sigma, \chi(\sigma))$ is small

Full runs as centroids

\mathcal{N} a process model, \mathcal{L} a log :

Find $\mathcal{C} = (\{u_1 \dots u_n\}, \chi)$:

$\{u_1 \dots u_n\} \in \text{Runs}(\mathcal{N})$

$\chi : \mathcal{L} \rightarrow \{nc, u_1 \dots u_n\}$

Such that for every $\sigma \in \mathcal{L}$, $\text{dist}(\sigma, \chi(\sigma))$ is small

Example :

$\sigma_1 = \langle S, C, G \rangle$

$\sigma_2 = \langle S, C, G, D \rangle$

Full runs as centroids

\mathcal{N} a process model, \mathcal{L} a log :

Find $\mathcal{C} = (\{u_1 \dots u_n\}, \chi)$:

$\{u_1 \dots u_n\} \in \text{Runs}(\mathcal{N})$

$\chi : \mathcal{L} \rightarrow \{nc, u_1 \dots u_n\}$

Such that for every $\sigma \in \mathcal{L}$, $\text{dist}(\sigma, \chi(\sigma))$ is small

Example :

$\sigma_1 = \langle S, C, G \rangle$

$\sigma_2 = \langle S, C, G, D \rangle$

$u = \langle S, C, \tau, G \rangle$

$\text{dist}(\sigma_1, u) = 0$

$\text{dist}(\sigma_2, u) = 1$

Full runs as centroids

\mathcal{N} a process model, \mathcal{L} a log :

Find $\mathcal{C} = (\{u_1 \dots u_n\}, \chi)$:

$\{u_1 \dots u_n\} \in \text{Runs}(\mathcal{N})$

$\chi : \mathcal{L} \rightarrow \{nc, u_1 \dots u_n\}$

Such that for every $\sigma \in \mathcal{L}$, $\text{dist}(\sigma, \chi(\sigma))$ is small

Example :

$\sigma_1 = \langle S, C, G \rangle$

$\sigma_2 = \langle S, C, G, D \rangle$

$u = \langle S, C, \tau, G \rangle$

$\text{dist}(\sigma_1, u) = 0$

$\text{dist}(\sigma_2, u) = 1$

dist is a distance between words (Hamming distance, Edit distance..)

- > Pseudo SAT formulas
- > solver : minisat+

> Pseudo SAT formulas

> solver : minisat+

Example of parts of the formula :

$$\bigwedge_{i=1}^n \bigvee_{a \in \Sigma} (\tau_{i,a} \wedge \bigwedge_{a' \in \Sigma \setminus t} \neg \tau_{i,a'})$$

transition a fires at instant i

> Pseudo SAT formulas

> solver : minisat+

Example of parts of the formula :

$$\bigwedge_{i=1}^n \bigvee_{a \in \Sigma} (\tau_{i,a} \wedge \bigwedge_{a' \in \Sigma \setminus t} \neg \tau_{i,a'})$$

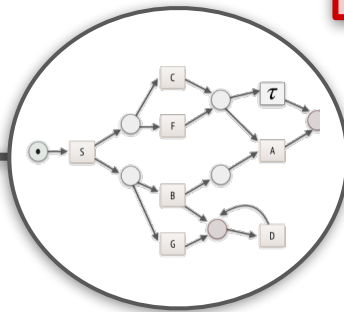
transition a fires at instant i

$$\bigwedge_{n=1}^k \bigwedge_{\sigma \in L} (\chi_{n,\sigma} \wedge \bigwedge_{n \neq m} \neg \chi_{m,\sigma})$$

trace σ is in cluster n

Log traces

< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >



Clusters

The use of full runs as centroids does not allow concurrency

< S, B, F, A >

< S, F, F, A >

< S, F, B, A >

< S, G, F, D, D >

< S, G, F, D, D, D, D >

< G, C, F, S, D, D >

< S, D, D, D >

Centroids : runs

< tau, G >

< S, B, F, A >

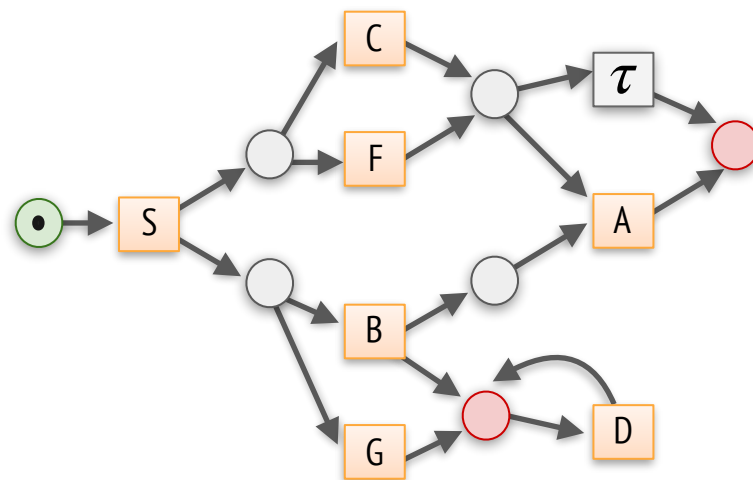
< S, F, B, A >

< S, G, F, D, D >

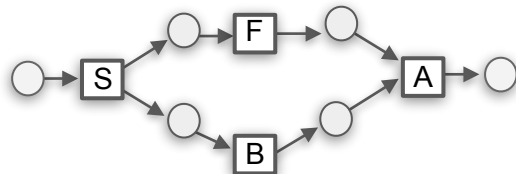
< S, G, F, D, D, D, D >

Non-clustered

Processes as centroids



Example of process :



Linearizations of the process :

$\langle S, B, F, A \rangle$

$\langle S, F, B, A \rangle$

[Engelfriet 1991]

Processes as centroids

\mathcal{N} a process model, \mathcal{L} a log :

Find $\mathcal{C} = (\{P_1 \dots P_n\}, \chi)$:

$\{P_1 \dots P_n\} \in \text{Processes}(\mathcal{N})$

$\chi : \mathcal{L} \rightarrow \{nc, P_1 \dots P_n\}$

Such that for every $\sigma \in \mathcal{L}$, $\text{dist}(\sigma, \chi(\sigma))$ is small

\mathcal{N} a process model, \mathcal{L} a log :

Find $\mathcal{C} = (\{P_1 \dots P_n\}, \chi)$:

$\{P_1 \dots P_n\} \in \text{Processes}(\mathcal{N})$

$\chi : \mathcal{L} \rightarrow \{nc, P_1 \dots P_n\}$

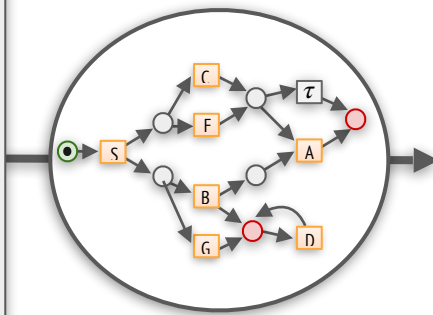
Such that for every $\sigma \in \mathcal{L}$, $\text{dist}(\sigma, \chi(\sigma))$ is small

dist is the minimal distance between a **linearization of P** and the trace (computed as distance between words : Hamming distance, Edit distance..)

Processes as centroids

Log traces

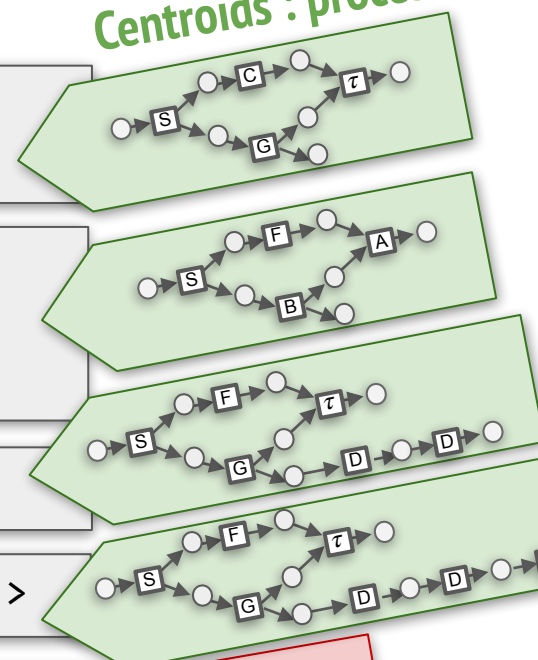
< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >



Clusters

< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >

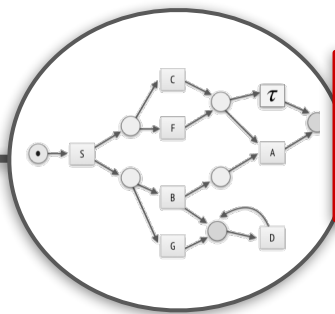
Centroids : processes



Non-clustered

Log traces

- < S, C, G >
- < S, C, G, D >
- < S, B, F, A >
- < S, F, F, A >
- < S, F, B, A >
- < S, G, F, D, D >
- < S, G, F, D, D, D, D >
- < G, C, F, S, D, D >
- < S, D, D, D >

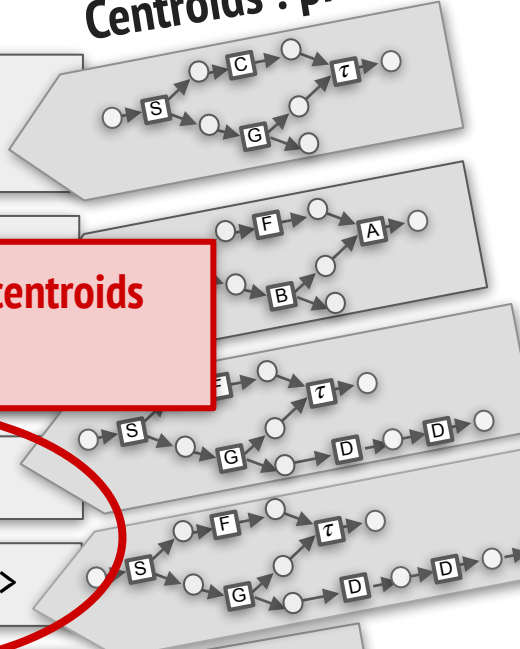


The use of processes as centroids does not allow loops

Clusters

- < S, C, G >
- < S, C, G, D >

Centroids : processes

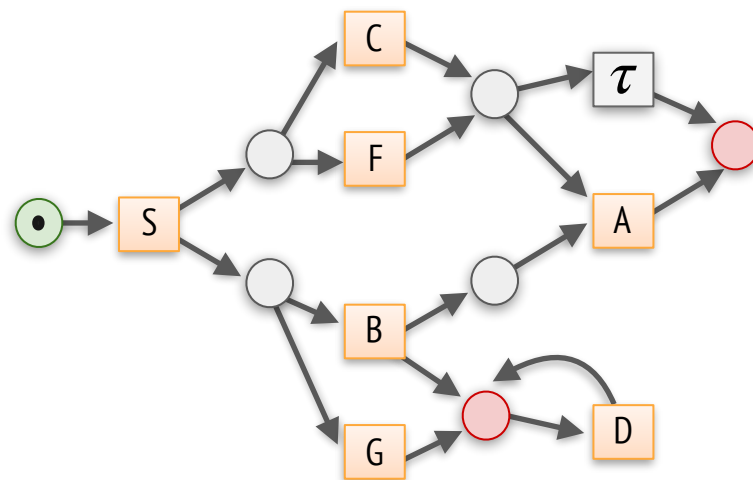


- < S, G, F, D, D >
- < S, G, F, D, D, D, D >

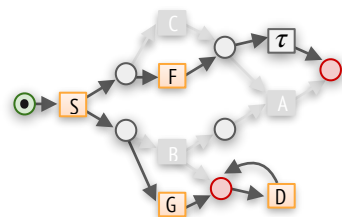
- < G, C, F, S, D, D >
- < S, D, D, D >

Non-clustered

Subnets as centroids



Example of subnet :



Runs of the subnet :

$\langle S, G, F, \tau, D \rangle$

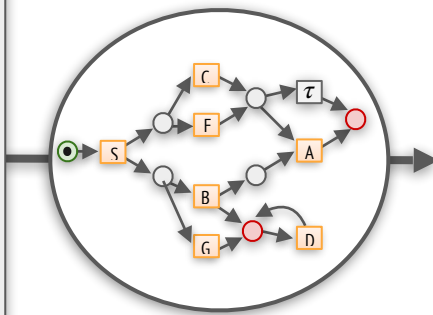
$\langle S, F, \tau, G, D, D \rangle$

$\langle S, G, D, D, D, D, F, \tau \rangle$

Subnets as centroids

Log traces

< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >



Clusters

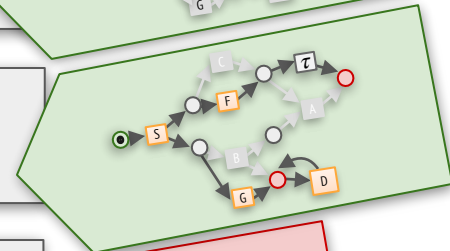
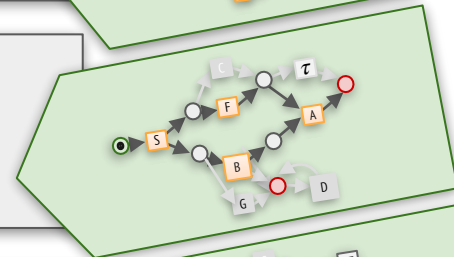
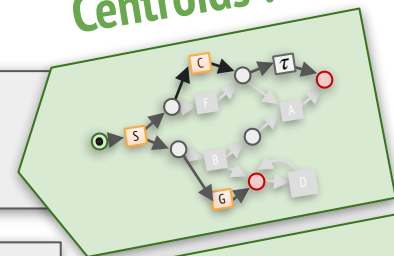
< S, C, G >
 < S, C, G, D >

< S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >

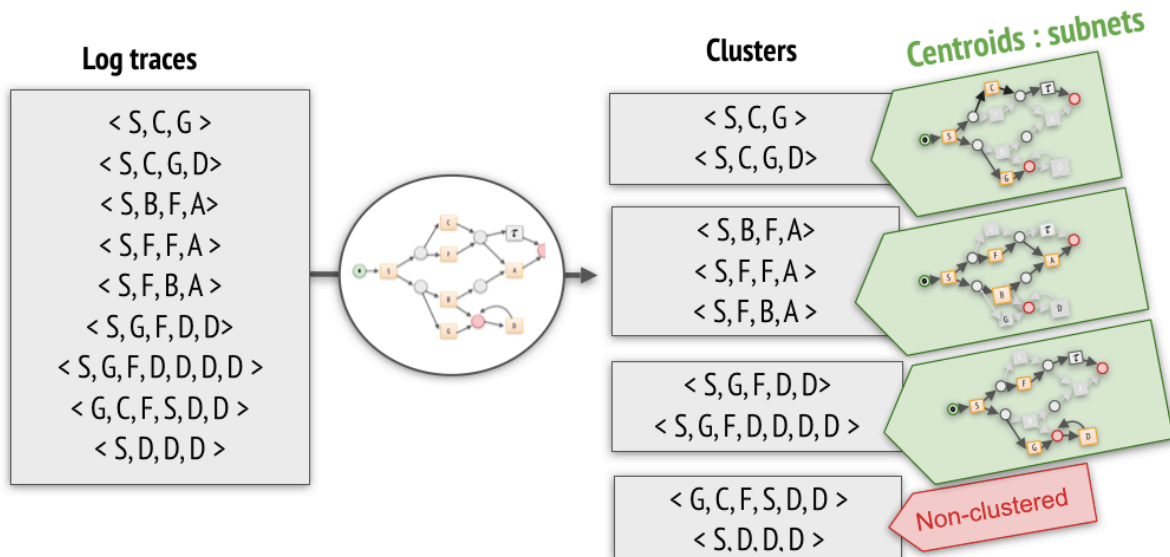
< S, G, F, D, D >
 < S, G, F, D, D, D, D >

< G, C, F, S, D, D >
 < S, D, D, D >

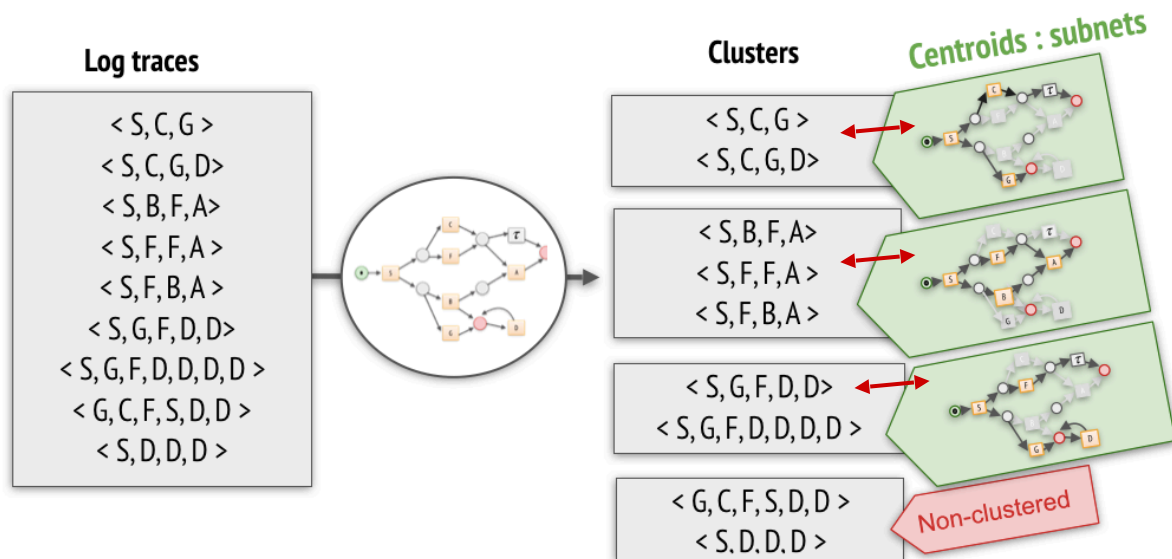
Centroids : subnets



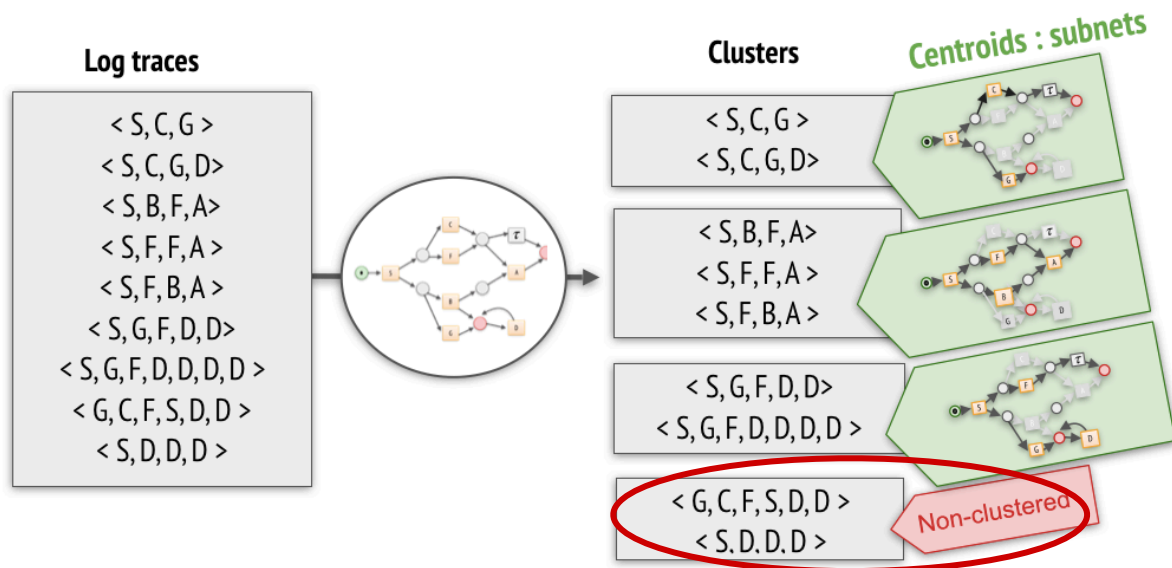
Non-clustered



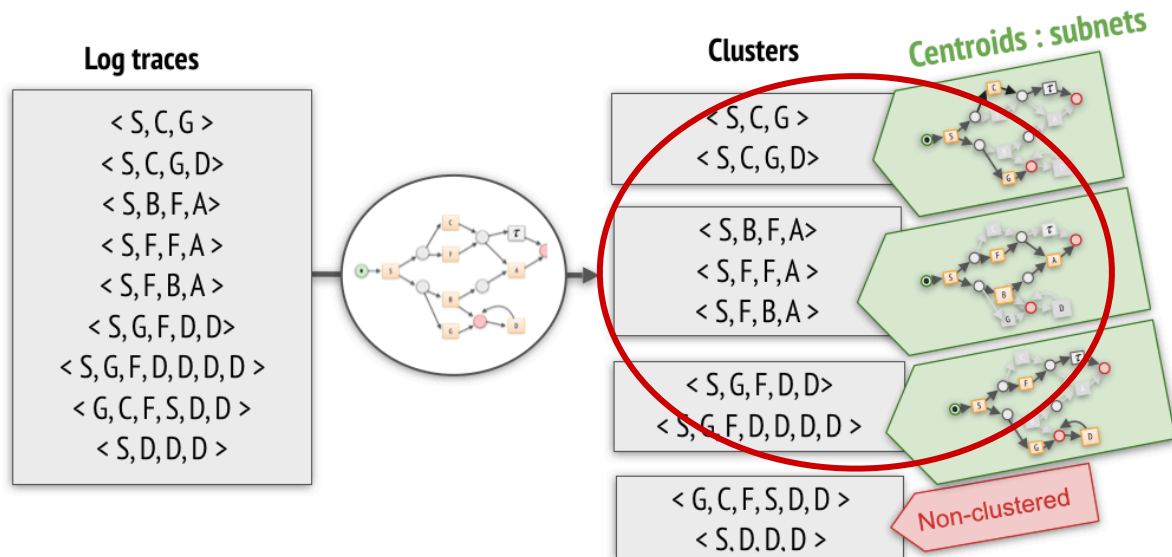
> Distance between traces and their centroid



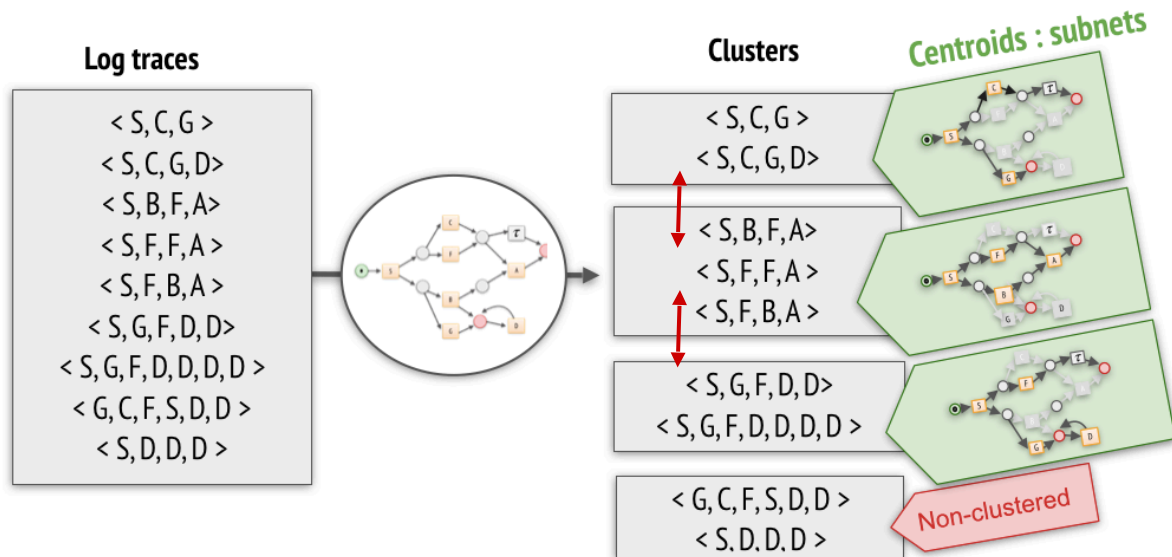
- > Distance between traces and their centroid
- > Number of non-clustered traces



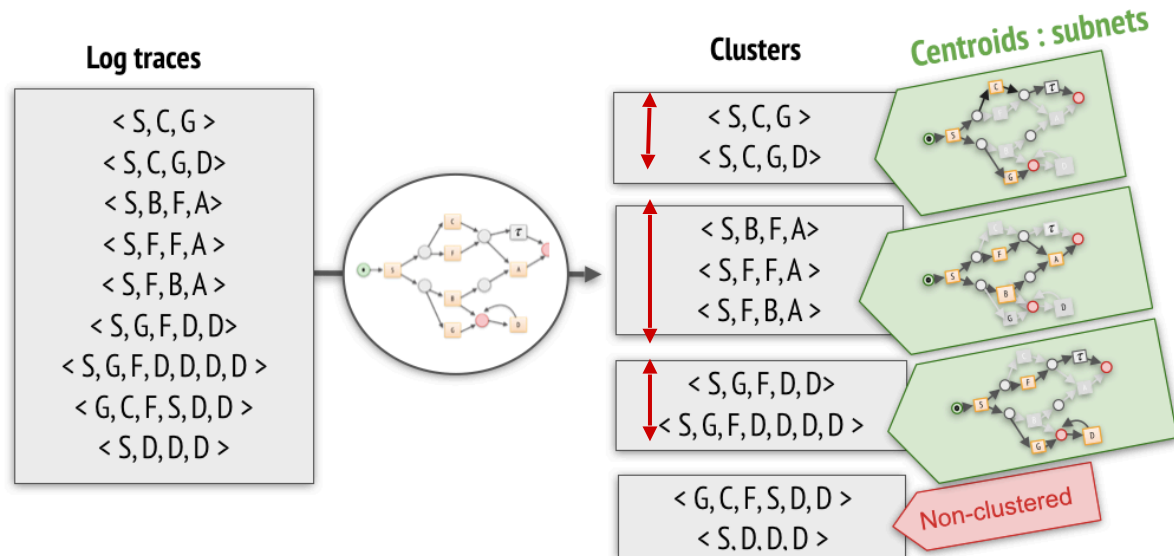
- > Distance between traces and their centroid
- > Number of non-clustered traces
- > Number of clusters



- > Distance between traces and their centroid
- > Number of non-clustered traces
- > Number of clusters
- > Inter-cluster distance



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



Inter-cluster distance

> maximise differences between clusters

Full runs centroids

$$\Phi(\mathcal{C}) \stackrel{\text{def}}{=} \min_{i \neq j} \text{dist}(u_i, u_j)$$

distance between centroids
computed as distance
between words (Hamming
distance, Edit distance..)

Inter-cluster distance

> maximise differences between clusters

Full runs centroids

$$\Phi(\mathcal{C}) \stackrel{\text{def}}{=} \min_{i \neq j} \text{dist}(u_i, u_j)$$

distance between centroids
computed as distance
between words (Hamming
distance, Edit distance..)

Processes centroids

$$\Phi(\mathcal{C}) = \min_{i \neq j} \text{dist}(\mathcal{P}_i, \mathcal{P}_j)$$

$$\text{dist}(\mathcal{P}, \mathcal{P}') \stackrel{\text{def}}{=} \min_{\substack{u \in \text{Runs}(\mathcal{P}) \\ u' \in \text{Runs}(\mathcal{P}')}} \text{dist}(u, u')$$

distance between centroids is
the minimal distance for any
linearizations of the processes

Subnets centroids

$$\Phi(\mathcal{C}) = \min_{i \neq j} \text{dist}(\mathcal{N}_i, \mathcal{N}_j)$$

$$\text{dist}(\mathcal{N}, \mathcal{N}') \stackrel{\text{def}}{=} \min_{\substack{u \in \text{Runs}(\mathcal{N}) \\ u' \in \text{Runs}(\mathcal{N}')}} \text{dist}(u, u')$$

distance between centroids is
the minimal distance for any
full runs of the subnets

Intra-cluster distance

- > only for subnet centroids
- > minimize differences in a cluster

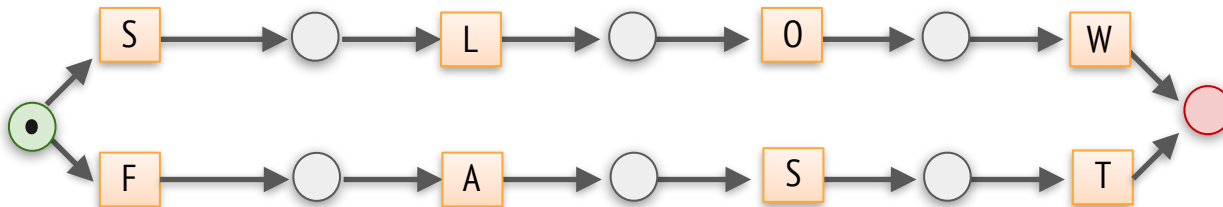
$$\Theta(\mathcal{C}) \stackrel{\text{def}}{=} \max_k \left(\sup_{\mathcal{P}, \mathcal{P}' \in Proc(\mathcal{N}_k)} dist(\mathcal{P}, \mathcal{P}') \right)$$

Intra-cluster distance

- > only for subnet centroids
- > minimize differences in a cluster

$$\theta(\mathcal{C}) \stackrel{\text{def}}{=} \max_k \left(\sup_{\mathcal{P}, \mathcal{P}' \in \text{Proc}(\mathcal{N}_k)} \text{dist}(\mathcal{P}, \mathcal{P}') \right)$$

A subnet with very different behaviors :

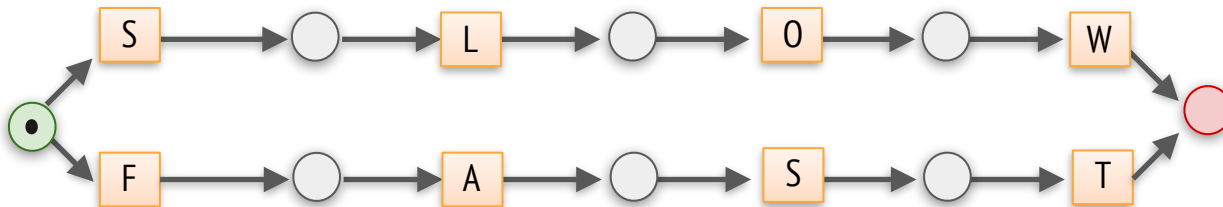


Intra-cluster distance

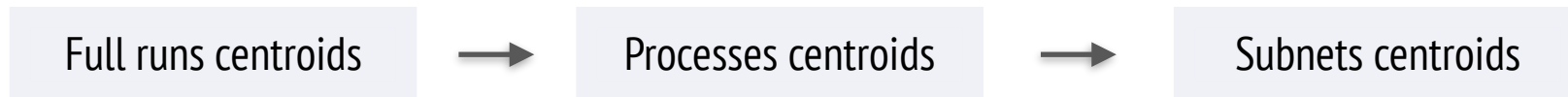
- > only for subnet centroids
- > minimize differences in a cluster

$$\Theta(\mathcal{C}) \stackrel{\text{def}}{=} \max_k \left(\sup_{\mathcal{P}, \mathcal{P}' \in \text{Proc}(\mathcal{N}_k)} \frac{\text{dist}(\mathcal{P}, \mathcal{P}')}{(1 + \epsilon)^{\max(|\mathcal{P}|, |\mathcal{P}'|)}} \right)$$

A subnet with very different behaviors :

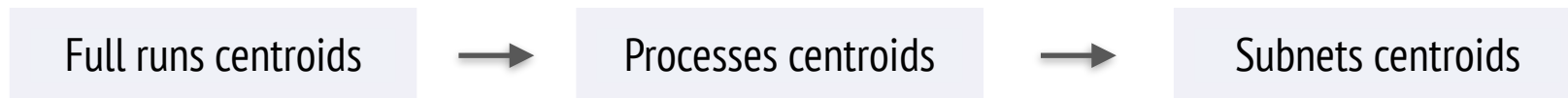


> Relating clustering by generalizing centroids



- > Decrease the distance between the traces and the centroids
- > Decrease the inter-cluster distance

> Relating clustering by generalizing centroids



$\langle S, B, F, A \rangle$

$\langle S, F, B, A \rangle$

> Relating clustering by generalizing centroids

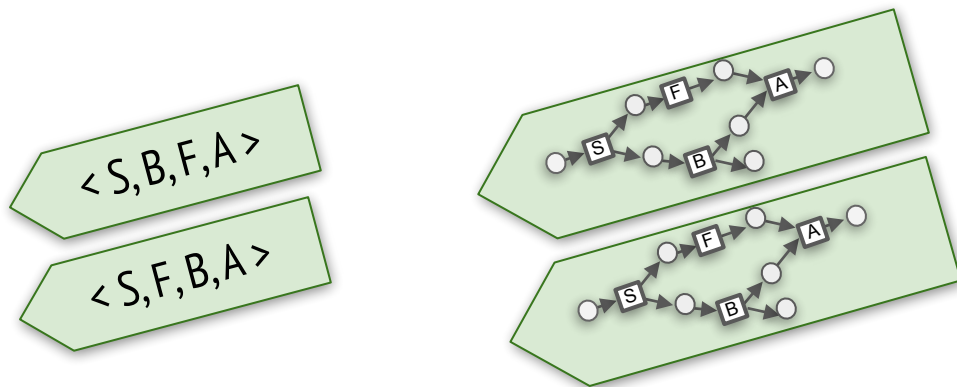
Full runs centroids



Processes centroids



Subnets centroids



> Relating clustering by generalizing centroids

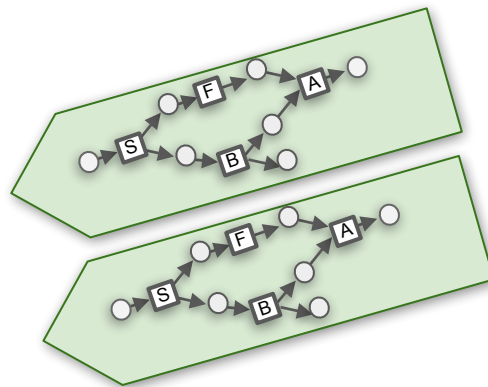
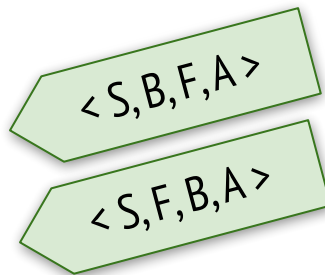
Full runs centroids



Processes centroids



Subnets centroids



$$\Phi(C) = 0$$

- > Tool DARK SIDER*
- > SAT formulas
- > Optimal clusterings

```
|T|: 18
|P|: 14
|L|: 9
clusters : full runs
nb variables : 75438
nb constrains : 165569
total time :5.099369
d(c) : 2
delta(c) : 16
n(c) : 3
ĉ(c) : 0.666667

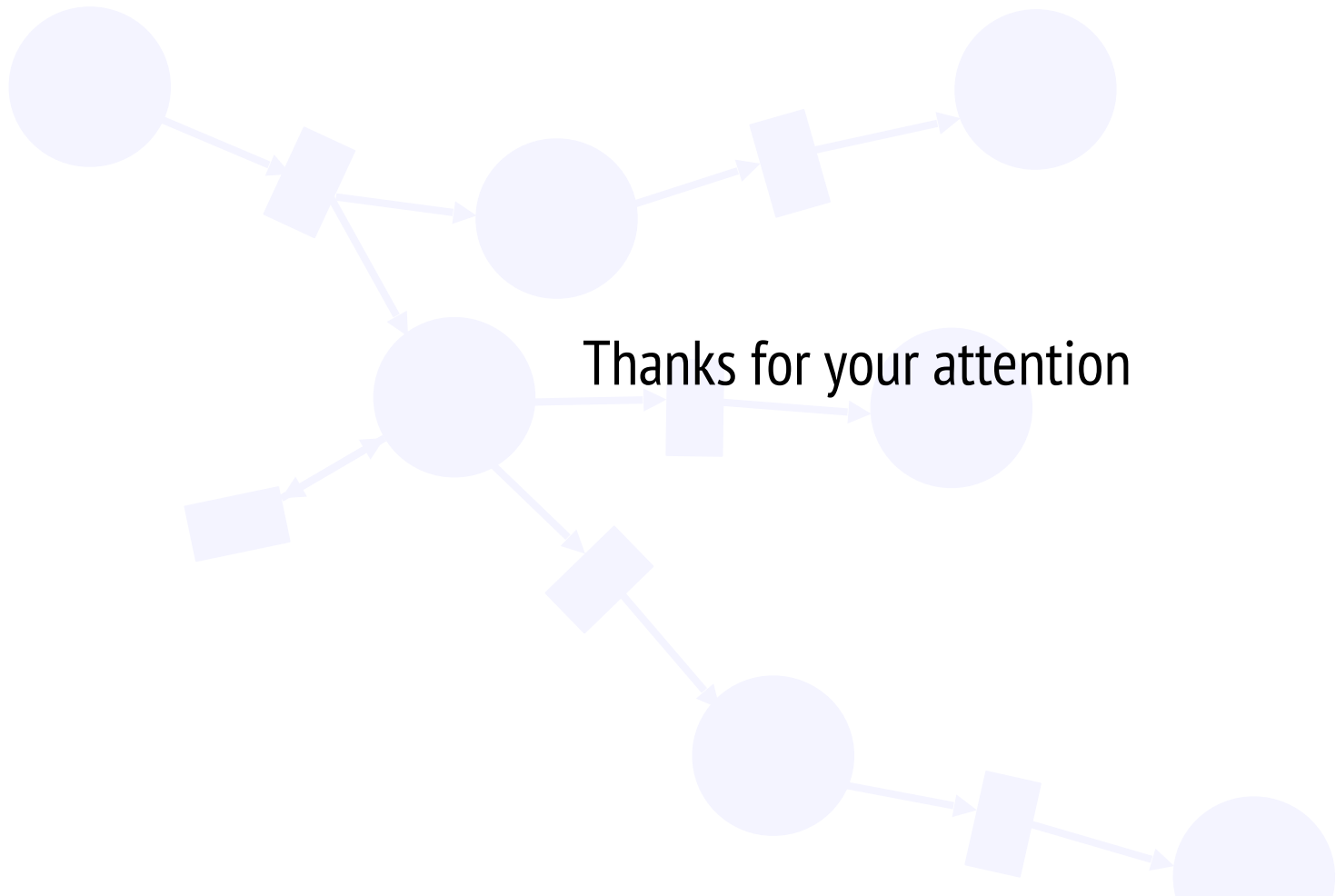
(A C E X3 I J X1 F D B , [A C E X3 I J X1 F D B  A C E X3 I J X1 F D B  ],0.222222)
(A C E L N M F D B , [A C E L N M F D B  A C E L O M F D B  ],0.222222)
(A C E G H F D B , [A C E G H F D B  A C E G H F F B  ],0.222222)
(nc, [ A C E X3 I J K I J X1 F D B  A C E X3 I J K I J K I J X1 F D B  A A C E G H F D B  ],0.333333)
```

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

- > 2 novel clustering methods for Process Mining problems
- > Quality criteria
- > Tool and experimentation

Limits : execution time

Future works : non-clustered traces





Generalized Alignment-Based Trace Clustering for Process Behavior

Mathilde Boltenhagen¹, Thomas Chatain¹, Josep Carmona²

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

²Universitat Politècnica de Catalunya