



Process Mining

Generalized Alignment-Based Trace Clustering

Mathilde Boltenhagen¹

J.Carmona²

T.Chatain²

5th, June 2019

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

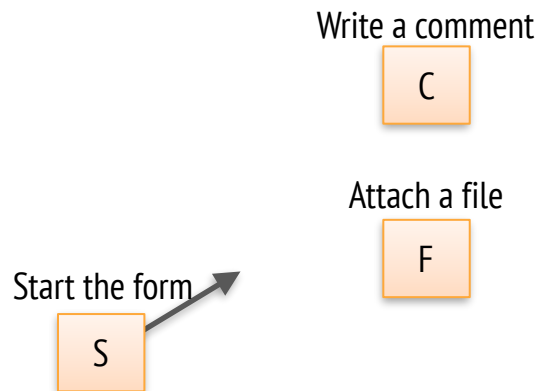
²Universitat Politècnica de Catalunya

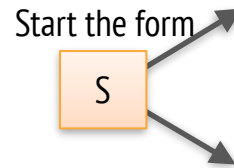




Start the form







Write a comment

C

Attach a file

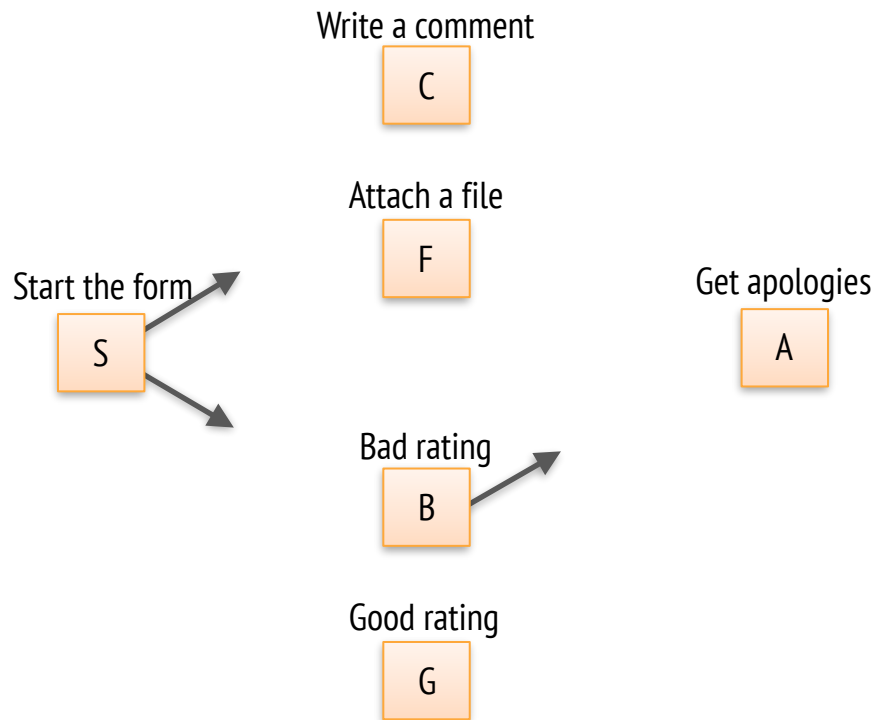
F

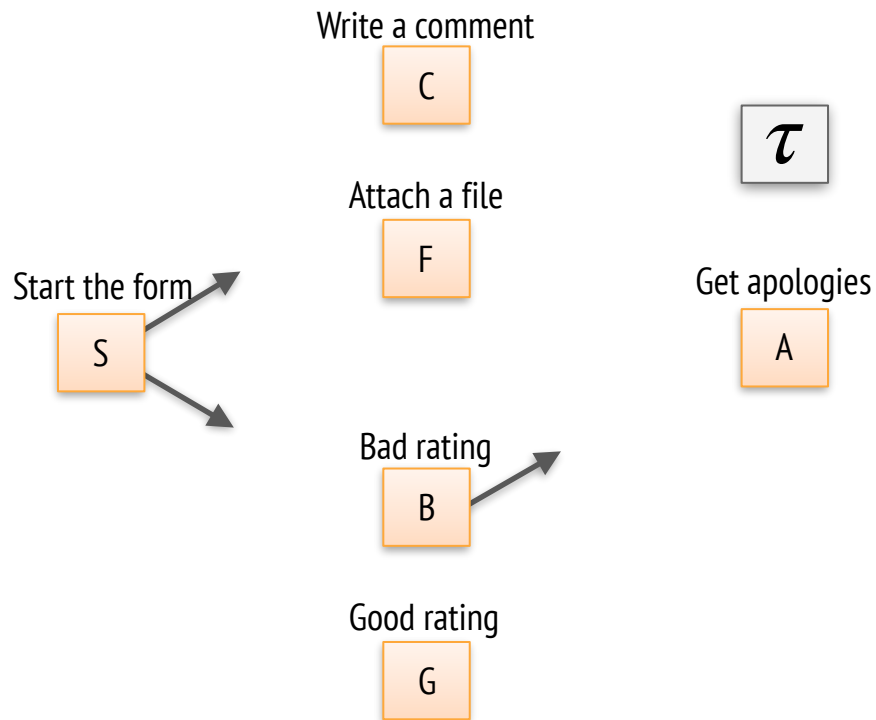
Bad rating

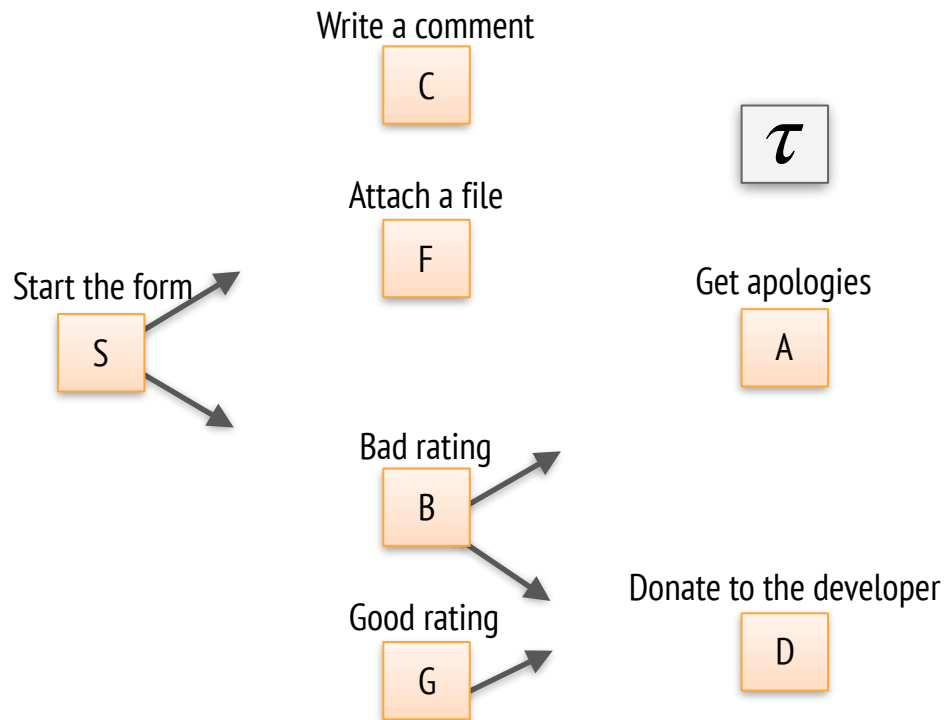
B

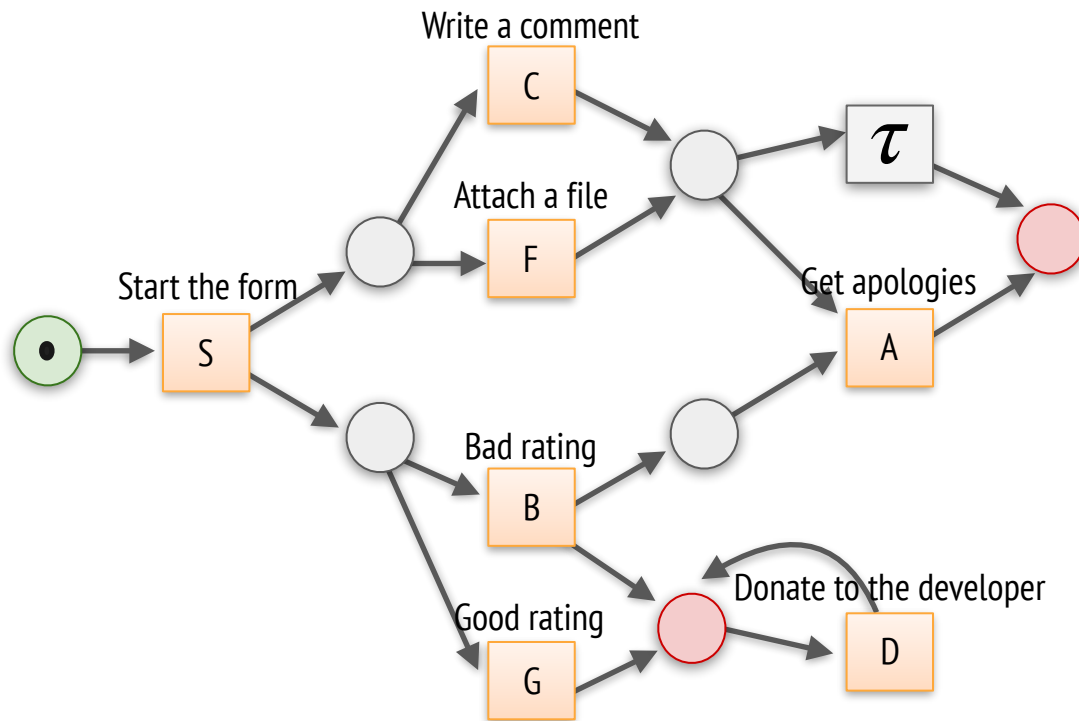
Good rating

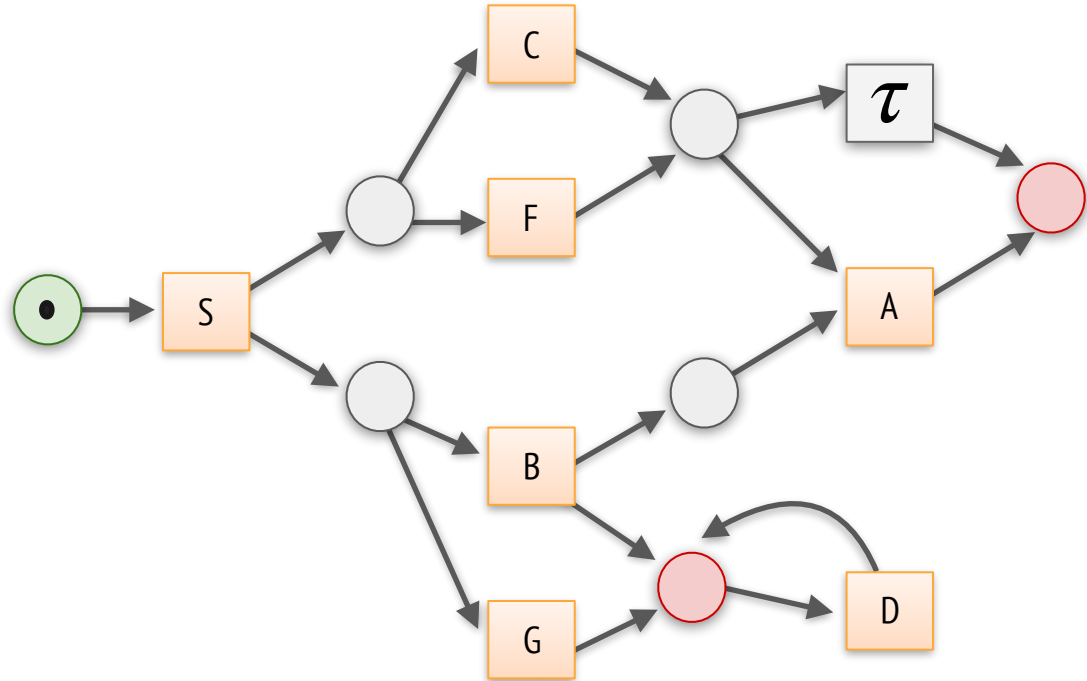
G











< Start the form, Give a comment, Good rating >
< Start the form, Give a comment, Good rating Donate to the developer>
 < Start the form, Bad rating, Send a File, Get Apologies >
 < Start the form, Send a File, Send a File, Get Apologies >
 < Start the form, Send a File, Bad rating, Get Apologies >
< Start the form, Good rating, Send a File, Donate to the developer, Donate to the developer, Donate to the developer Donate to the developer >
 < Good rating, Give a comment, Start the form, Donate to the developer, Donate to the developer >
 < Start the form, Donate to the developer, Donate to the developer, Donate to the developer >

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

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

Data clustering

- Unsupervised algorithms
- No classes
- Number of clusters is unknown

To mine databases

(k-means, self-organizing map ...)

Data classification

- Supervised algorithms
- Labelled data/ known classes
- Number of classes is known

To respond to defined problems

(k-NN, SVM ...)

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 >

[Greco 2006] ; [Ferreira 2007]

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 >

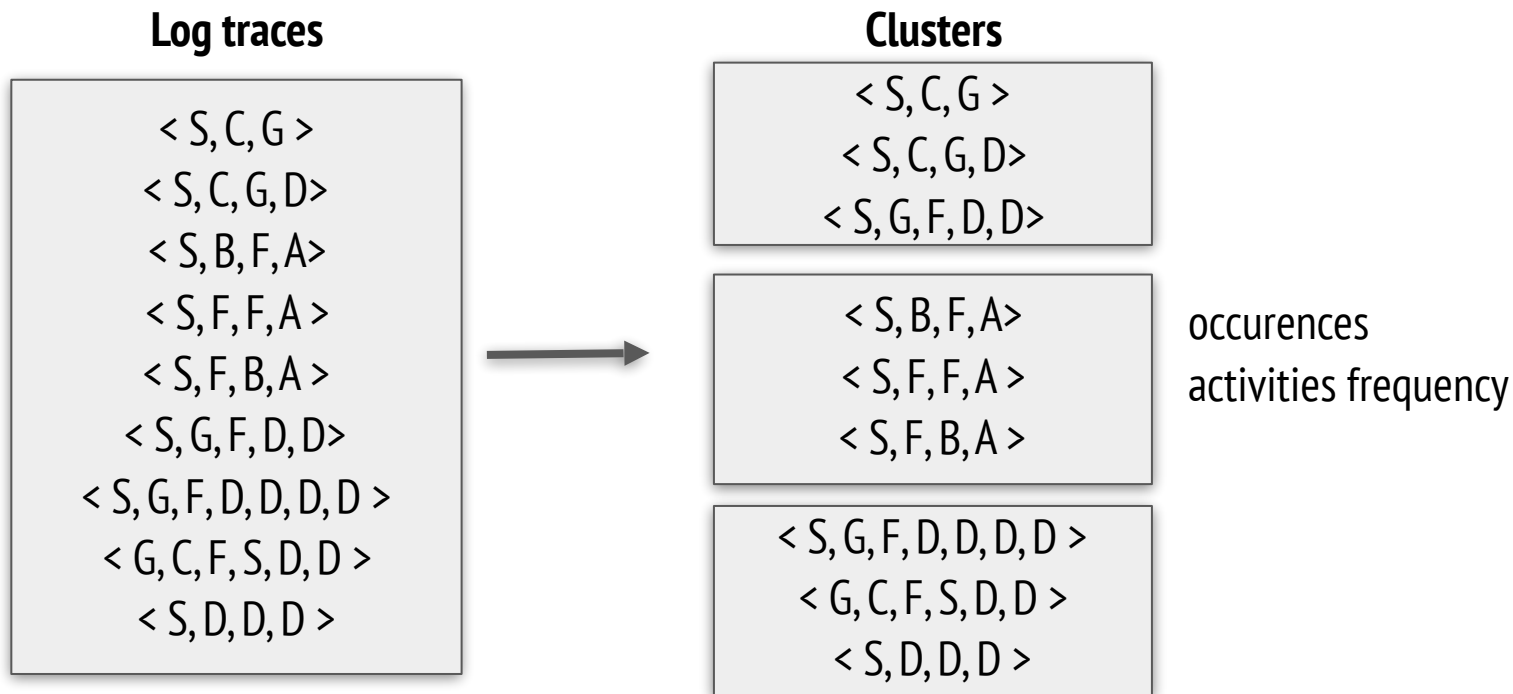
**Clusters**

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

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

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

[Greco 2006] ; [Ferreira 2007]



[Greco 2006] ; [Ferreira 2007]

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

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

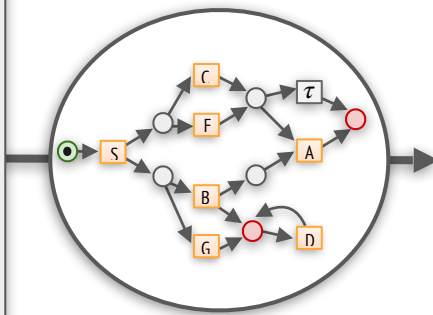
Log traces

< 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 2017]

Log traces

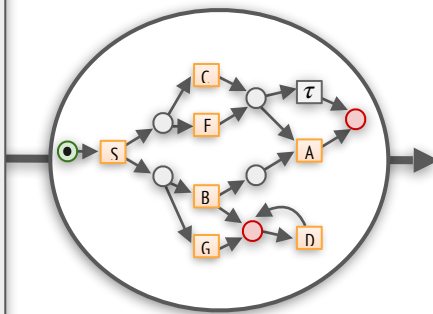
< S, C, G >
< S, C, G, D >
< S, B, F, A >
< S, F, F, A >
< S, F, B, A >
< S, G, F, D, D >
< S, G, F, D, D, D, D >
< G, C, F, S, D, D >
< S, D, D, D >



[Chatain 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 >

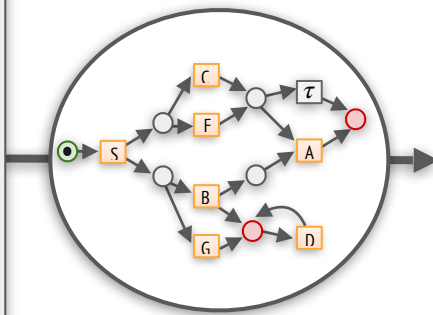
< S, G, F, D, D, D, D >

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

[Chatain 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 >

Centroids : runs

< S, C, τ , G >

< S, B, F, A >

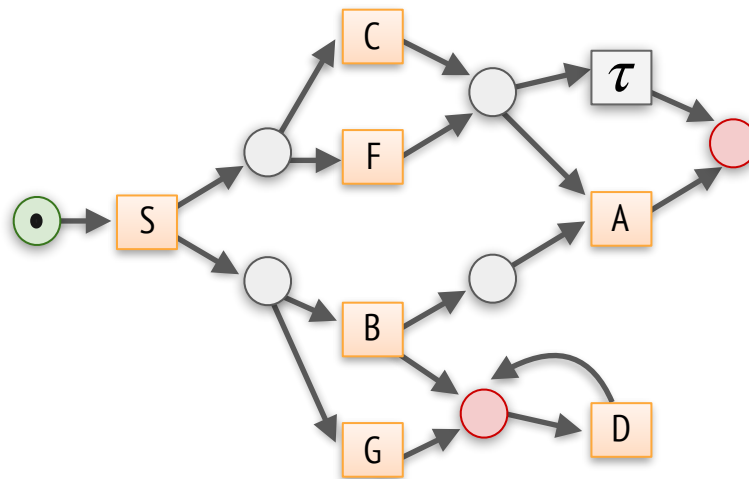
< S, F, B, A >

< S, G, F, D, D >

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

Non-clustered

[Chatain 2017]



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

[Chatain 2017]

N a process model, **L** a log :

For every trace $\sigma \in \mathbf{L}$,

Find $u \in \text{Runs}(\mathbf{N})$, *dist*(σ, u) is small

u is a centroid

N a process model, **L** a log :

For every trace $\sigma \in \mathbf{L}$,

Find $u \in \text{Runs}(\mathbf{N})$, *dist*(σ, u) is small

u is a centroid

Example :

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

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

[Chatain 2017]

N a process model, **L** a log :

For every trace $\sigma \in \mathbf{L}$,

Find $u \in \text{Runs}(\mathbf{N})$, *dist*(σ, u) is small

u is a centroid

Example :

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

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

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

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

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

[Chatain 2017]

N a process model, **L** a log :

For every trace $\sigma \in \mathbf{L}$,

Find $u \in \text{Runs}(\mathbf{N})$, *dist*(σ, u) is small

u is a centroid

Example :

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

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

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

dist(σ_1, u) = 0

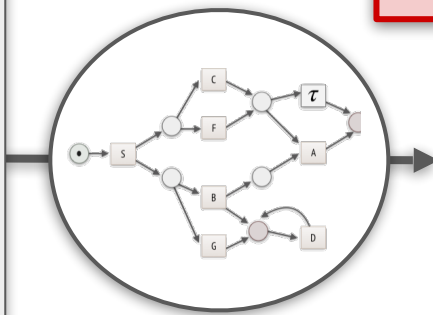
dist(σ_2, u) = 1

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

[Chatain 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

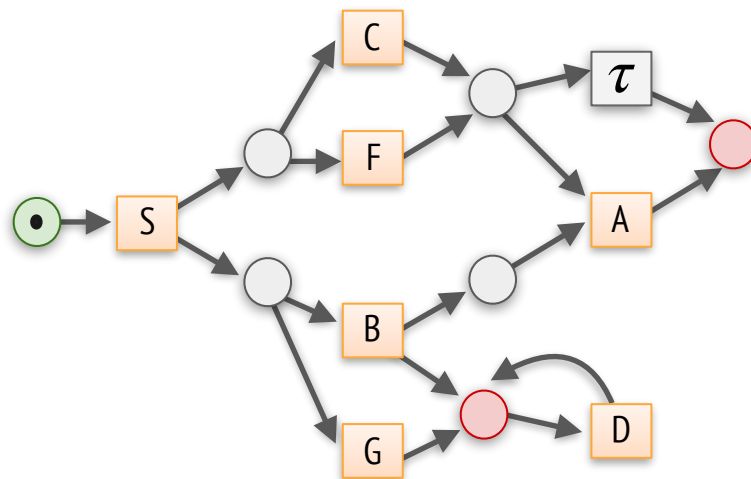
The use of full runs as centroids does not allow concurrency

< 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 >

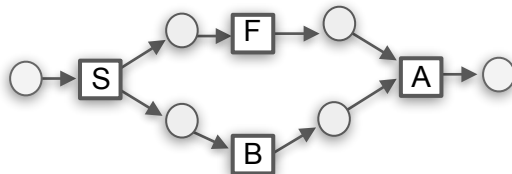
Centroids : runs

< S, B, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 Non-clustered

[Chatain 2017]



Example of process :



Linearizations of the process :

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

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

[Boltenhagen 2019]

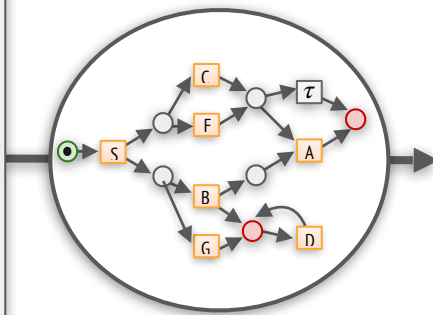
Log traces

< 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 >

[Boltenhagen 2019]

Log traces

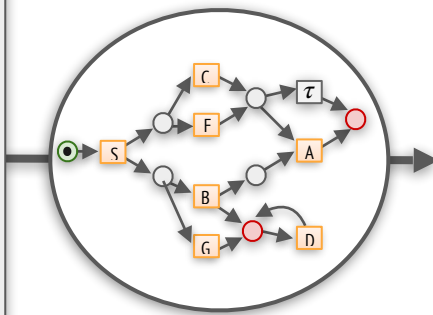
< S, C, G >
< S, C, G, D >
< S, B, F, A >
< S, F, F, A >
< S, F, B, A >
< S, G, F, D, D >
< S, G, F, D, D, D, D >
< G, C, F, S, D, D >
< S, D, D, D >



[Boltenhagen 2019]

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 >

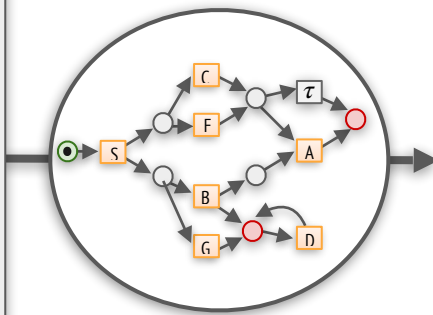
< S, G, F, D, D, D, D >

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

[Boltenhagen 2019]

Log traces

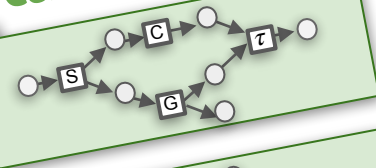
< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >



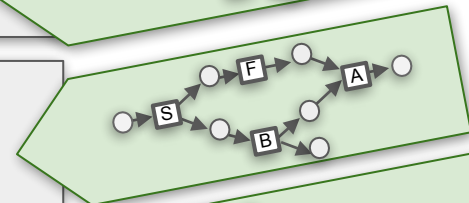
Clusters

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

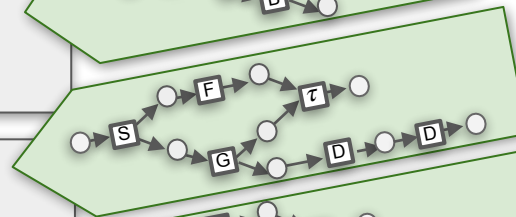
Centroids : processes



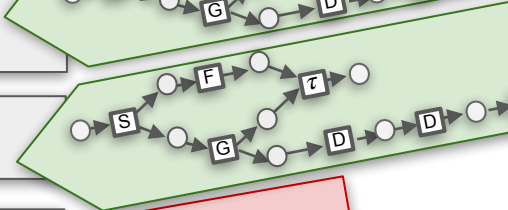
< S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >



< S, G, F, D, D >



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



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

Non-clustered

N a process model, **L** a log :

For every trace $\sigma \in \mathbf{L}$,

Find $\mathbf{P} \in \text{Processes}(\mathbf{N})$, *dist*(σ, \mathbf{P}) is small

\mathbf{P} is a centroid

N a process model, **L** a log :

For every trace $\sigma \in \mathbf{L}$,

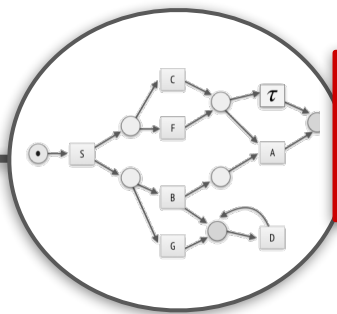
Find $\mathbf{P} \in \text{Processes}(\mathbf{N})$, *dist*(σ, \mathbf{P}) is small

\mathbf{P} is a centroid

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

Log traces

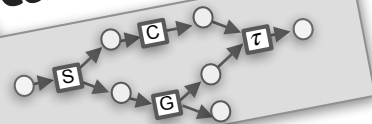
< S, C, G >
 < S, C, G, D >
 < S, B, F, A >
 < S, F, F, A >
 < S, F, B, A >
 < S, G, F, D, D >
 < S, G, F, D, D, D, D >
 < G, C, F, S, D, D >
 < S, D, D, D >



Clusters

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

Centroids : processes



The use of processes as centroids does not allow loops

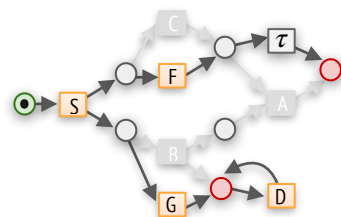
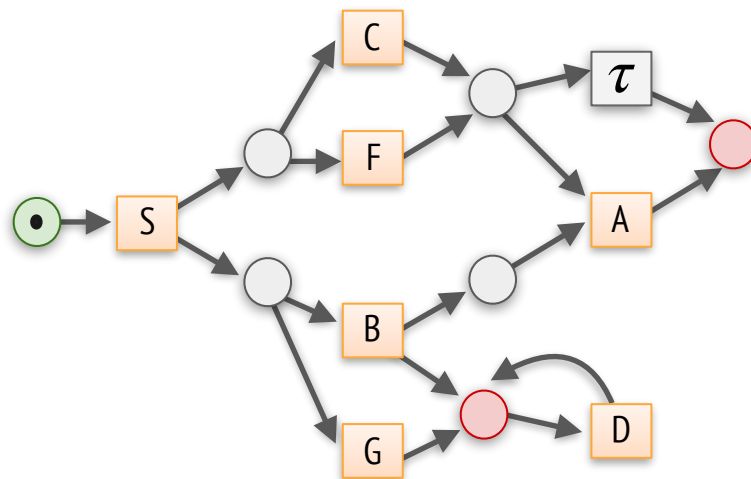
< S, G, F, D, D >

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

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

< S, D, D, D >

Non-clustered



Example of subnet :

[Boltenhagen 2019]

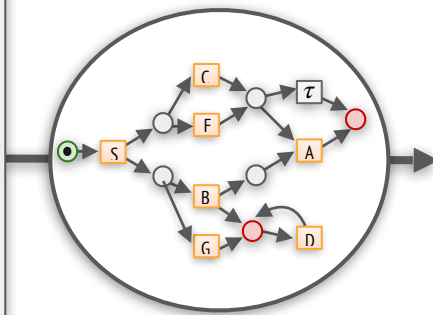
Log traces

< 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 >

[Boltenhagen 2019]

Log traces

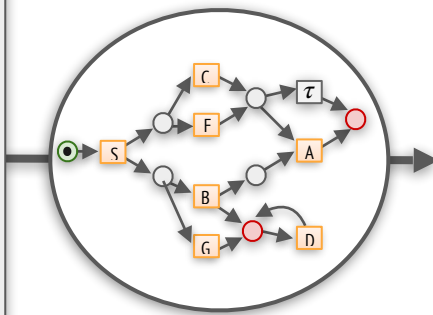
< S, C, G >
< S, C, G, D >
< S, B, F, A >
< S, F, F, A >
< S, F, B, A >
< S, G, F, D, D >
< S, G, F, D, D, D, D >
< G, C, F, S, D, D >
< S, D, D, D >



[Boltenhagen 2019]

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 >

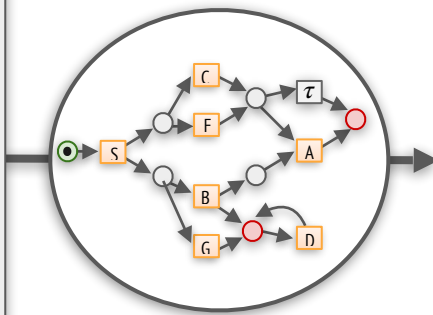
< S, G, F, D, D >
< S, G, F, D, D, D, D >

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

[Boltenhagen 2019]

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

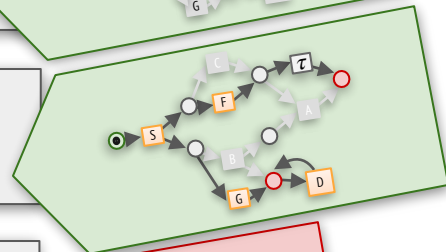
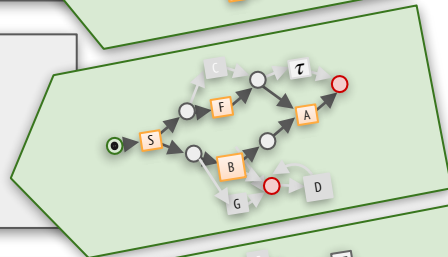
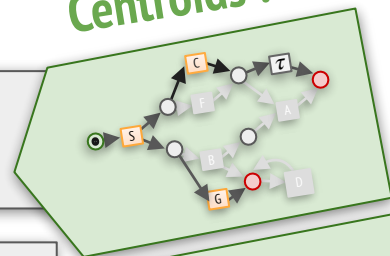
< S, C, G >
 < S, C, G, D >

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

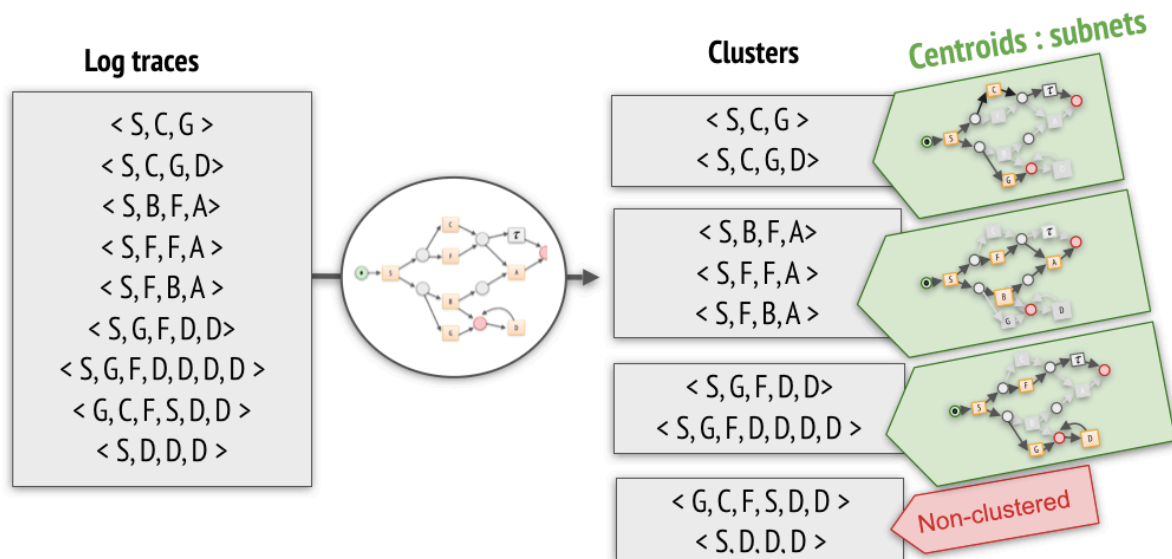
< S, G, F, D, D >
 < S, G, F, D, D, D, D >

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

Centroids : subnets

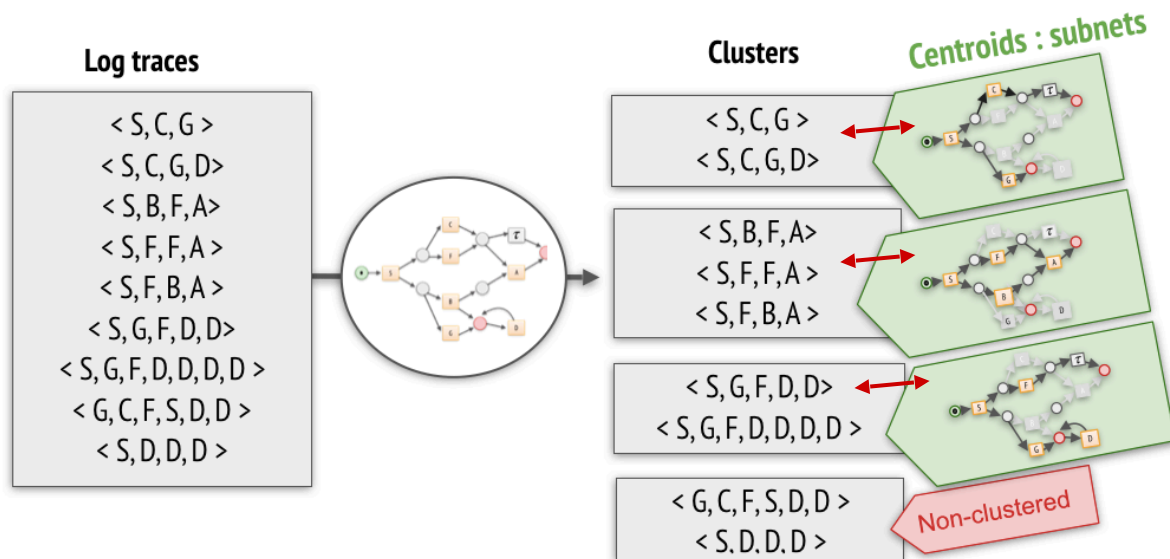


Non-clustered



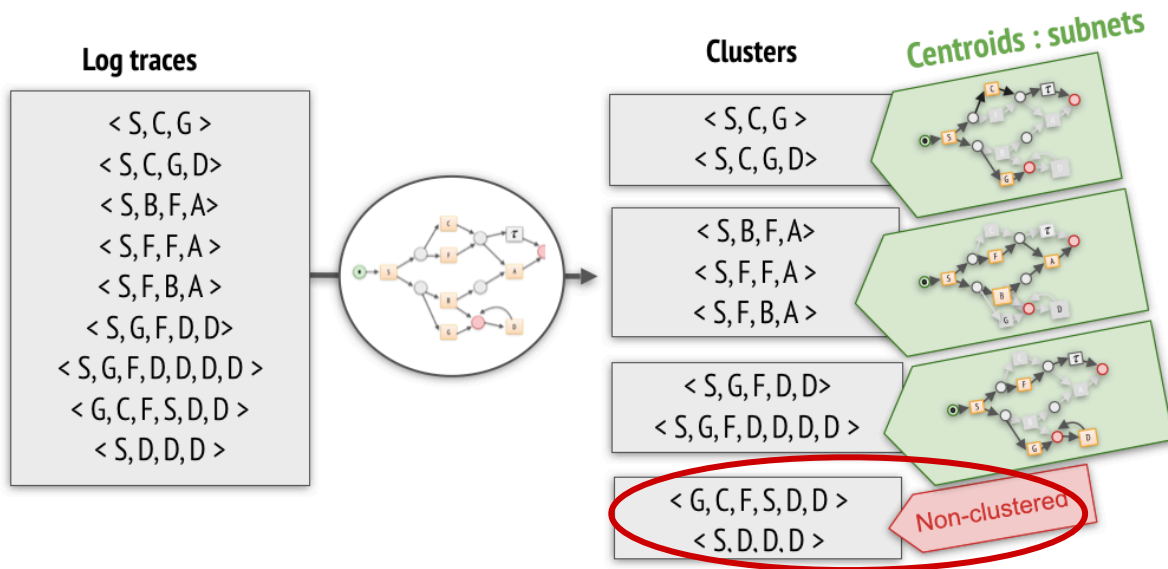
[Boltenhagen 2019]

> Distance between traces and their centroid

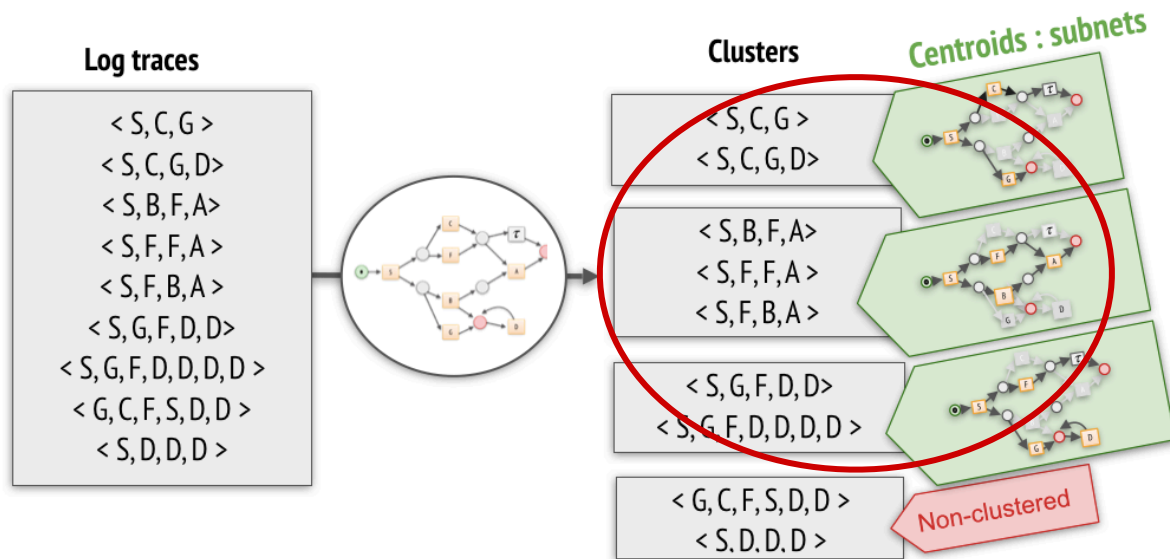


[Boltenhagen 2019]

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

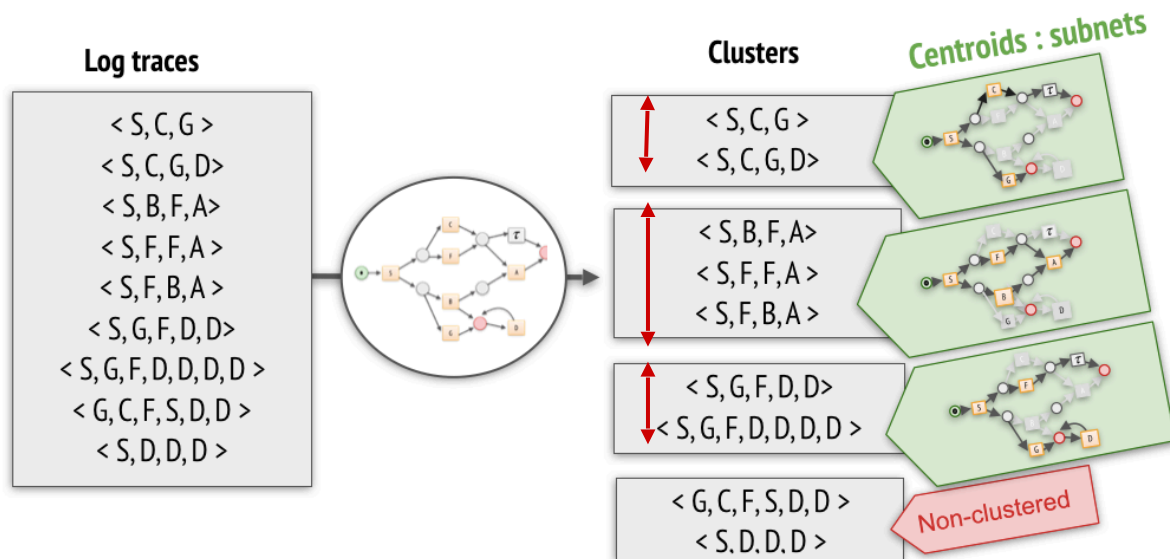


- > Distance between traces and their centroid
- > Number of non-clustered traces
- > Number of clusters



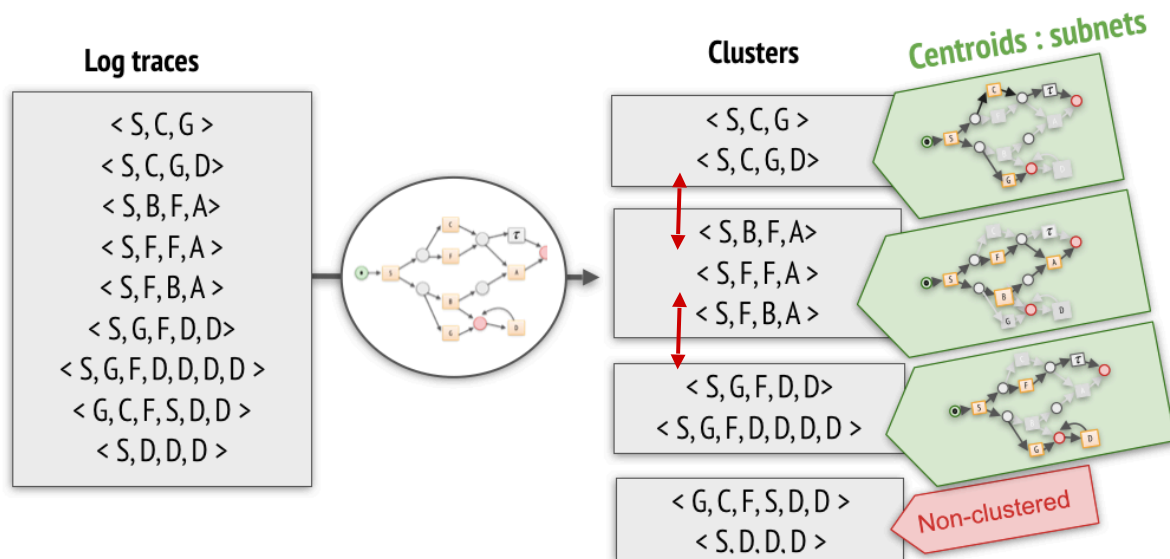
[Boltenhagen 2019]

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



[Boltenhagen 2019]

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



[Boltenhagen 2019]

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

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

How can one repair an
existing process model ?

At least a better
fitness

How can one repair an
existing process model ?

Not too far from
the original model

....

