

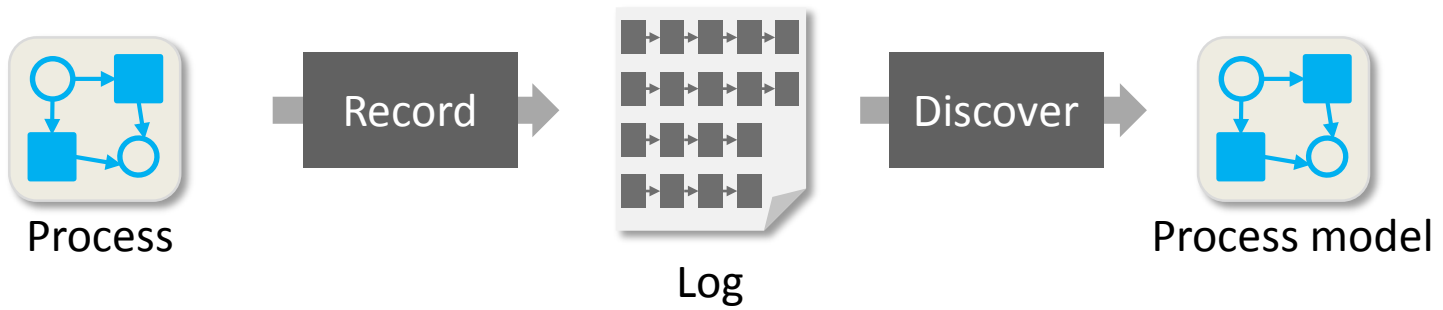
Finding Similar and Dissimilar Events for Preprocessing Logs and Improving Mining Results - *A Snapshot of Thesis and Tools*

Xixi Lu, Dirk Fahland, Wil van der Aalst
Eindhoven University of Technology

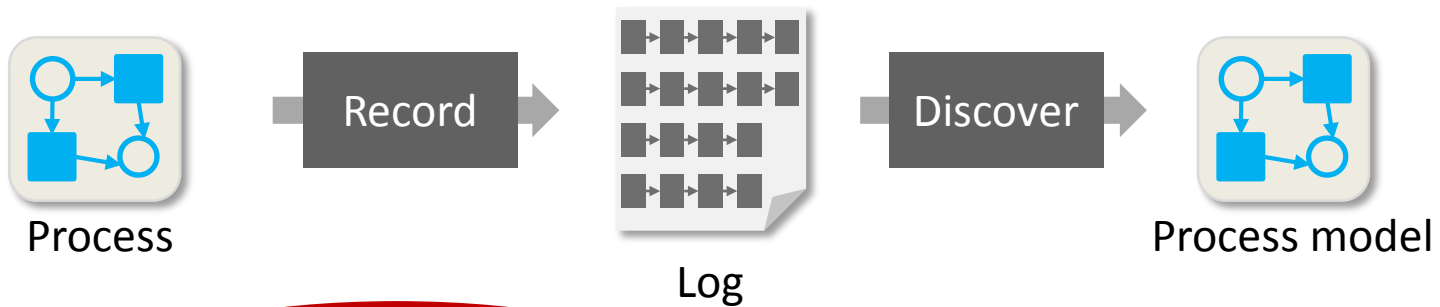
Agenda

- (5 min) I. Process discovery and challenges
- (5 min) II. Research problem
- (10 min) III. Approach and applications (Recall)
- (20 min) IV. Demo and Discussions

Process Discovery



Challenges in Flexible Environments

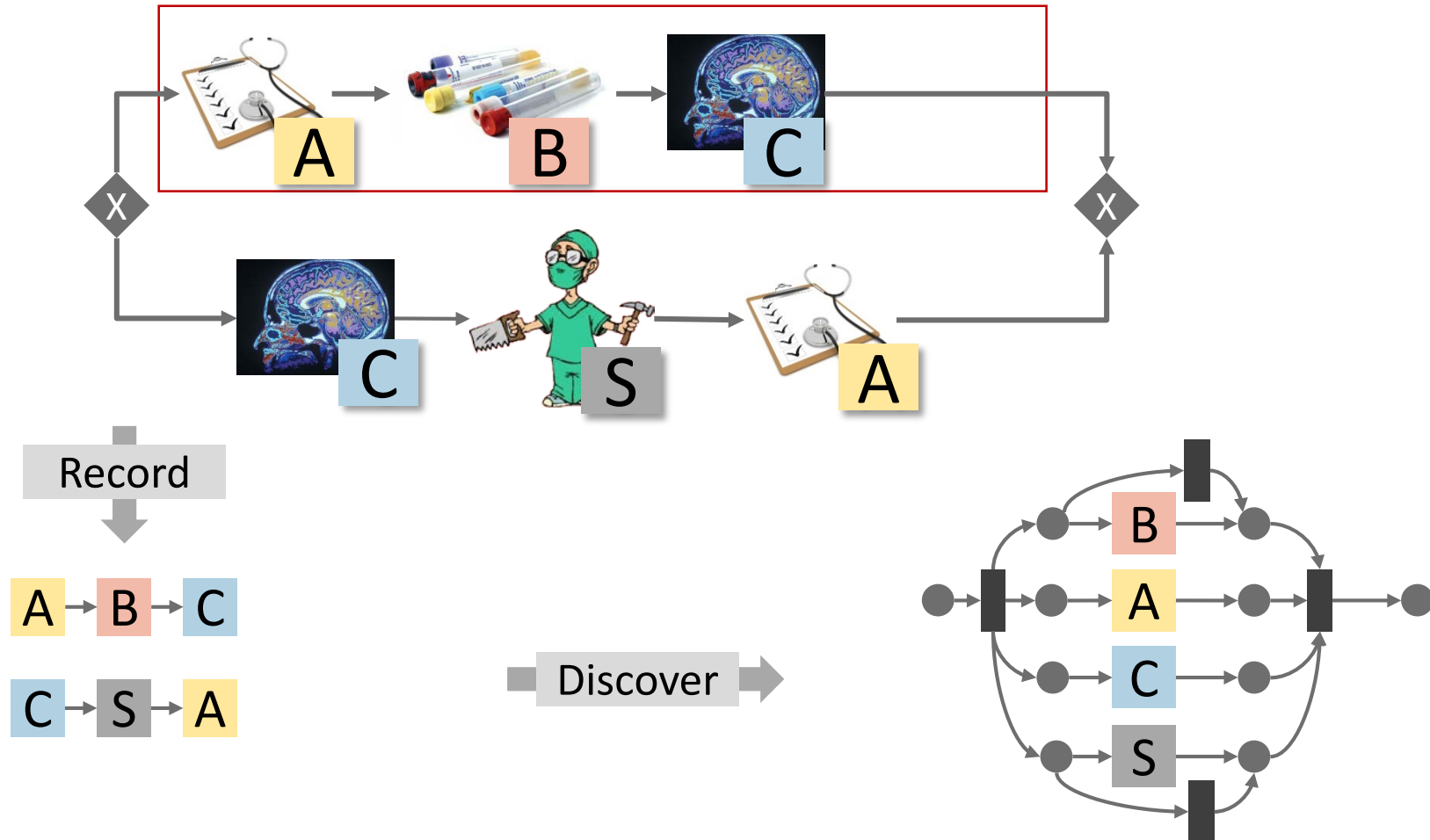


1. High variety
of behavior

2. Deviating
behavior

3. Duplicated
tasks

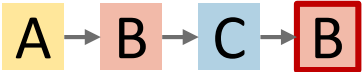
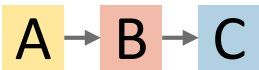
1. High Variety of Behavior



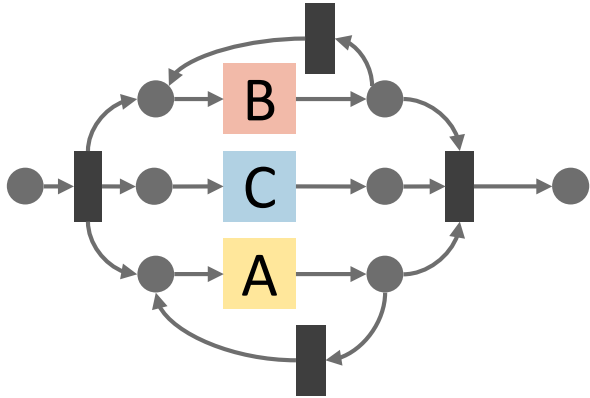
2. Deviating behavior



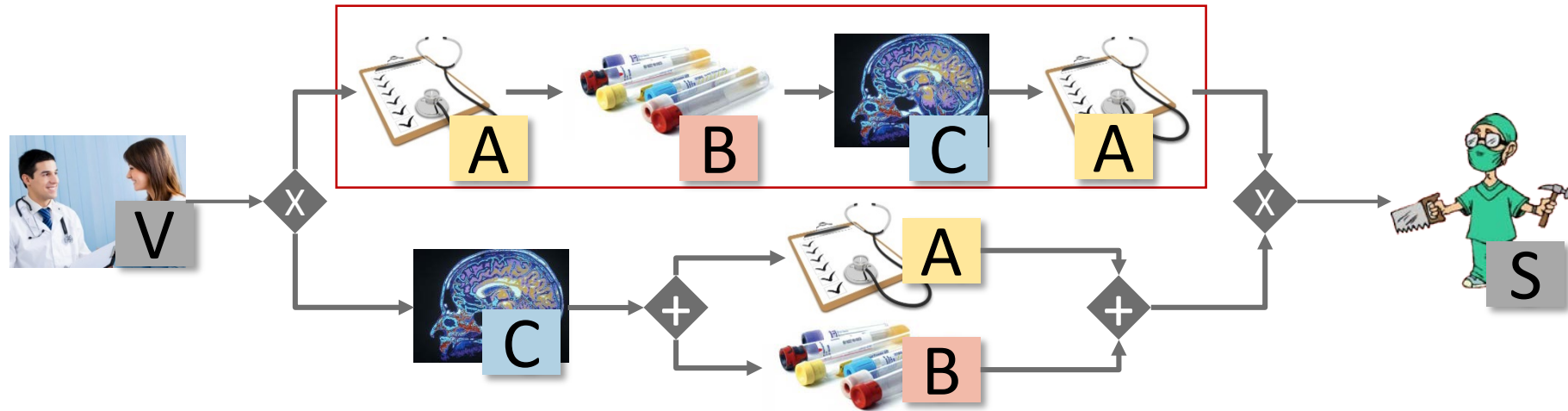
Record



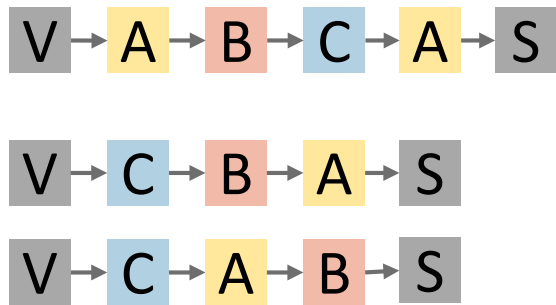
Discover



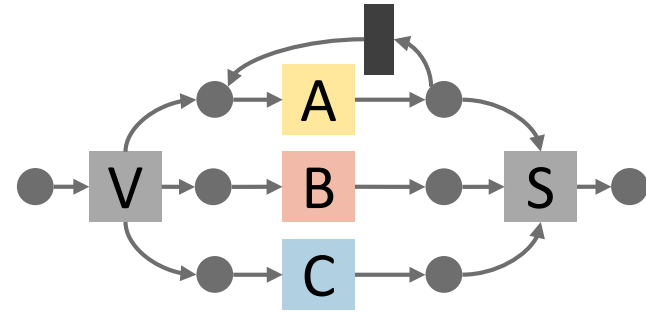
3. Duplicated Tasks



Record



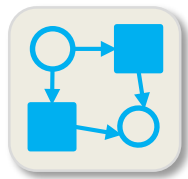
Discover



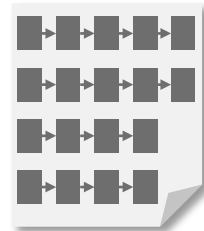
Challenges

System unknown!

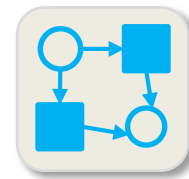
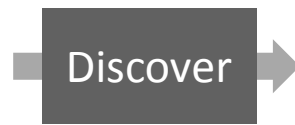
We don't know what is optimal!



Process



Log



Process model

1. High variety of behavior

2. Deviating behavior

3. Duplicated tasks

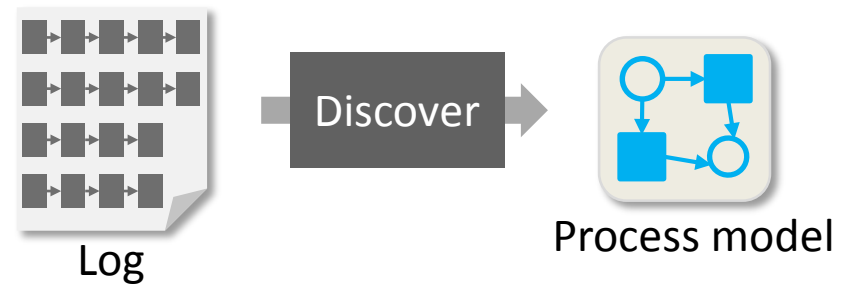
Cluster?

Filter?

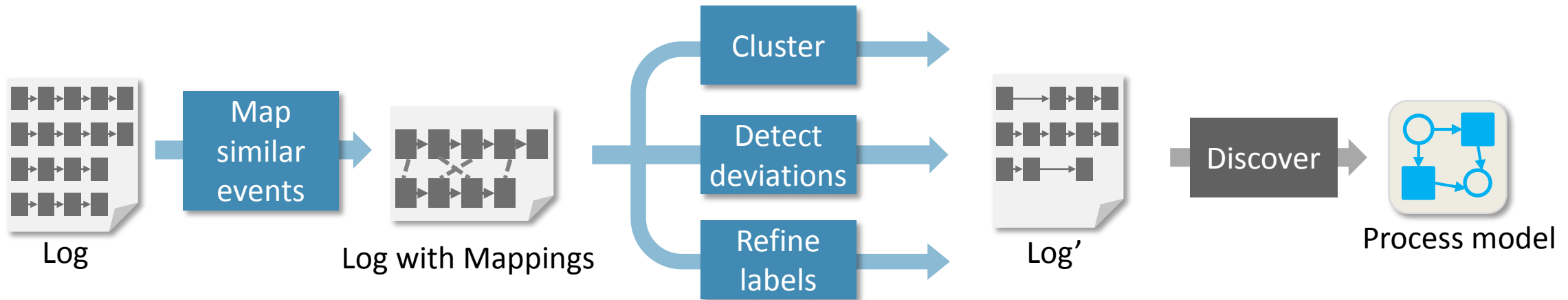
Yet Another Discovery?

Having dedicated algorithms?

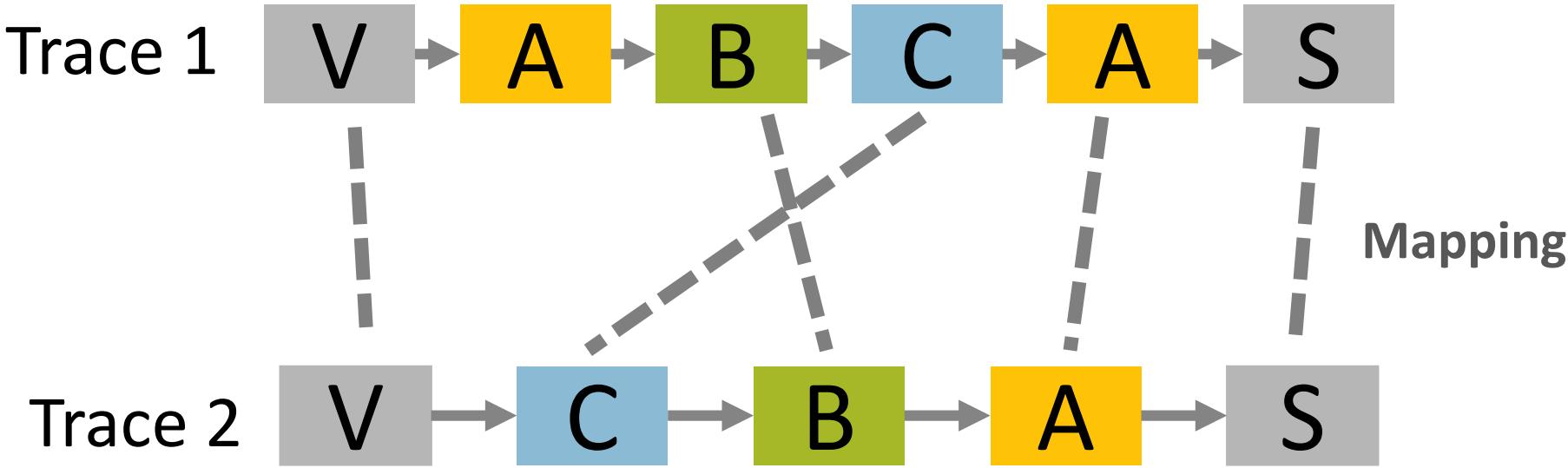
Challenges



Research Problem and Approach

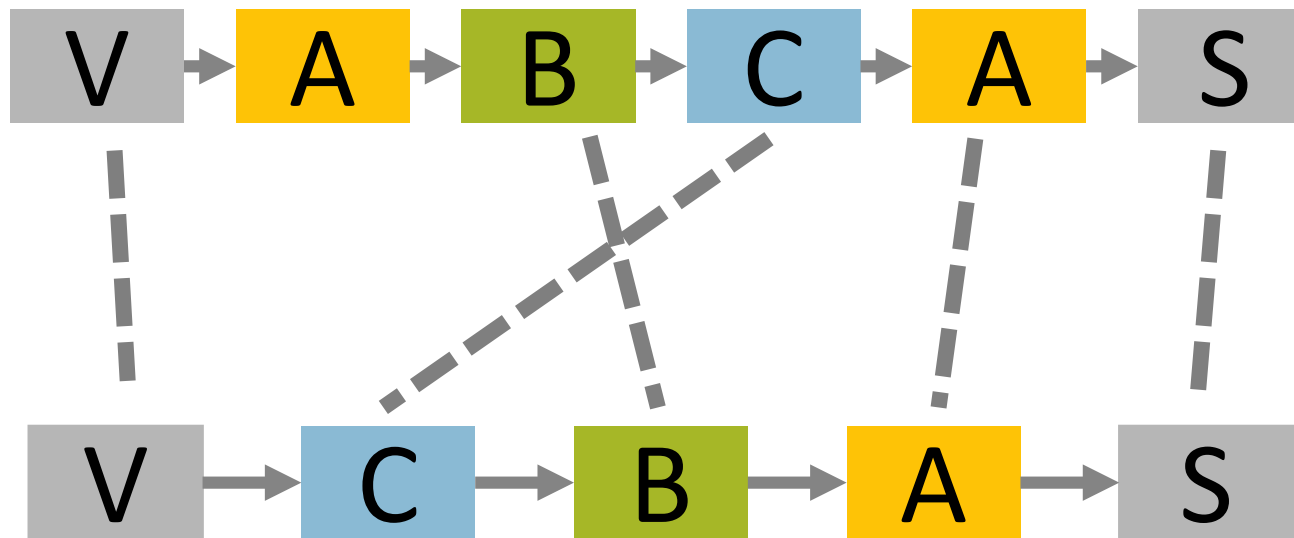


Mapping Between Events



Quantify Similarity of Mapped Events

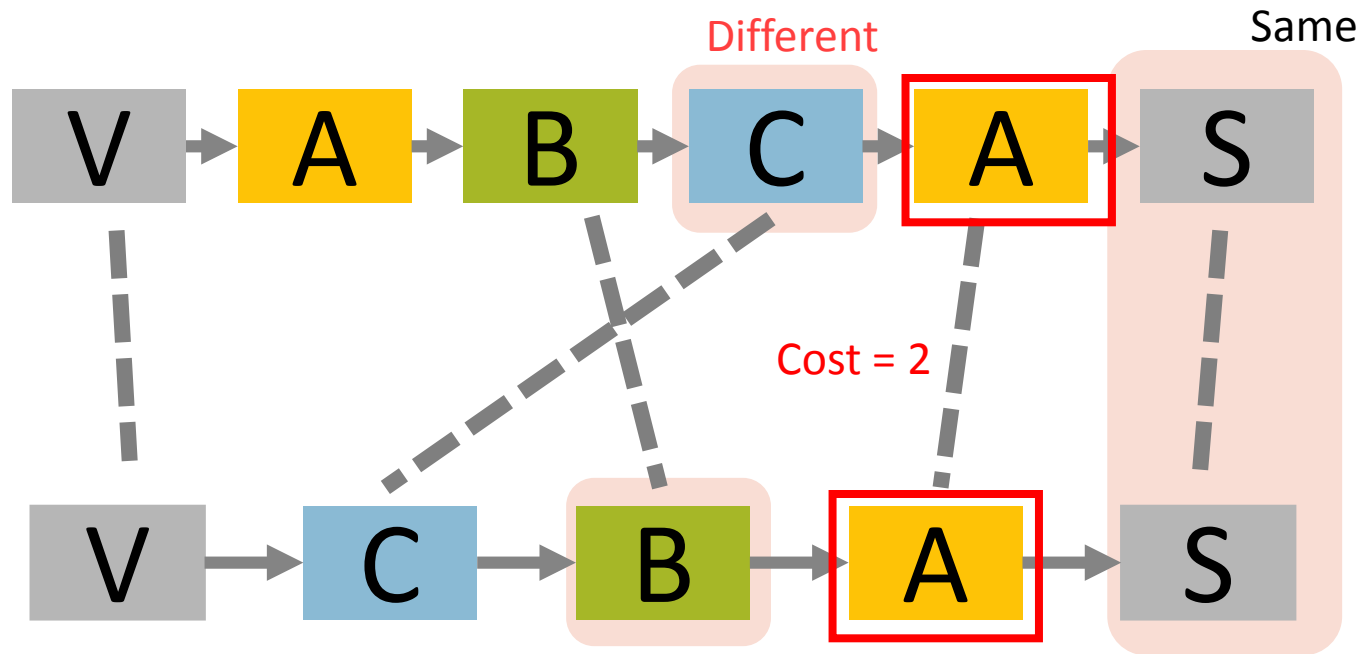
... based on “structural context of events”



Quantify Similarity of Mapped Events

... based on “structural context of events”

(1) Differences in neighbors

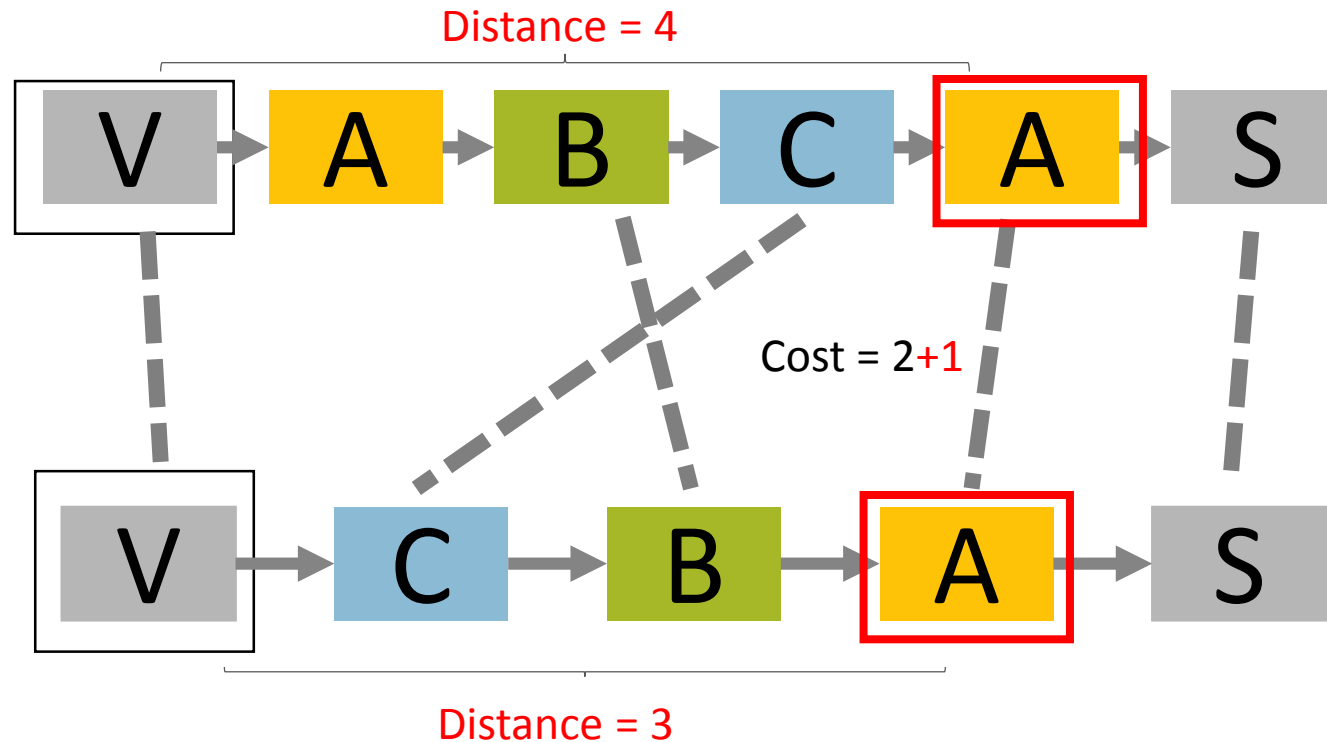


Quantify Similarity of Mapped Events

... based on “structural context of events”

(1) Differences in neighbors

(2) Differences in structure

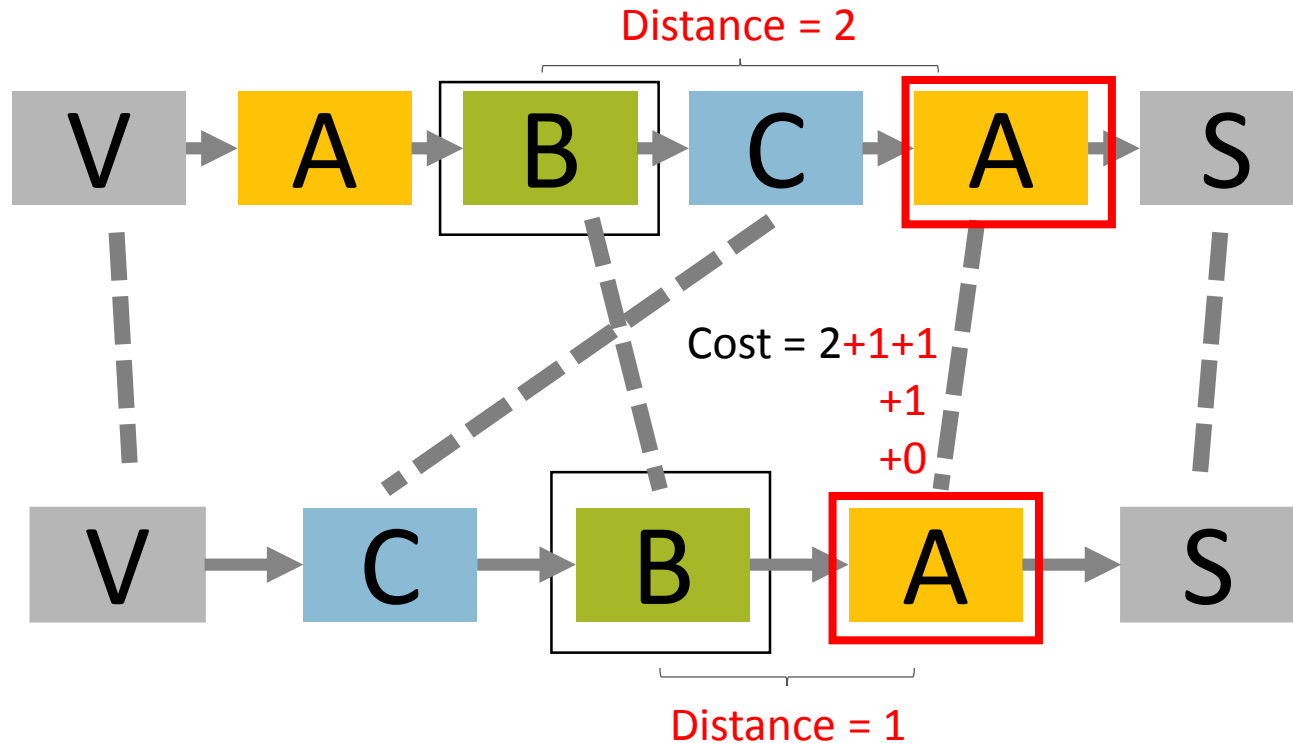


Quantify Similarity of Mapped Events

... based on “structural context of events”

(1) Differences in neighbors

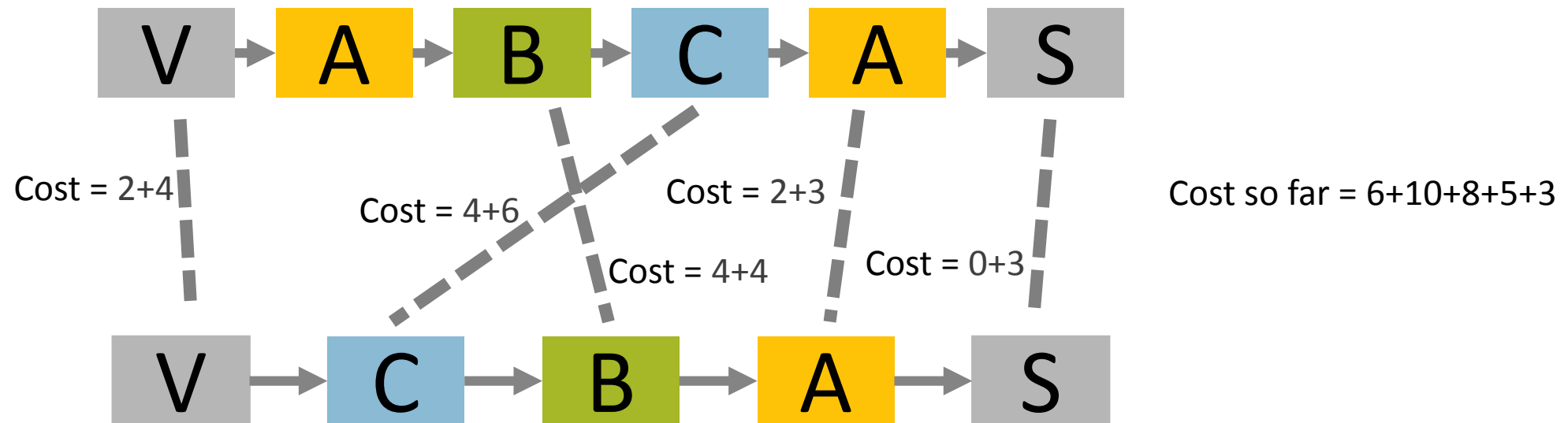
(2) Differences in structure



Quantify Similarity of Mapped Events

... based on “structural context of events”

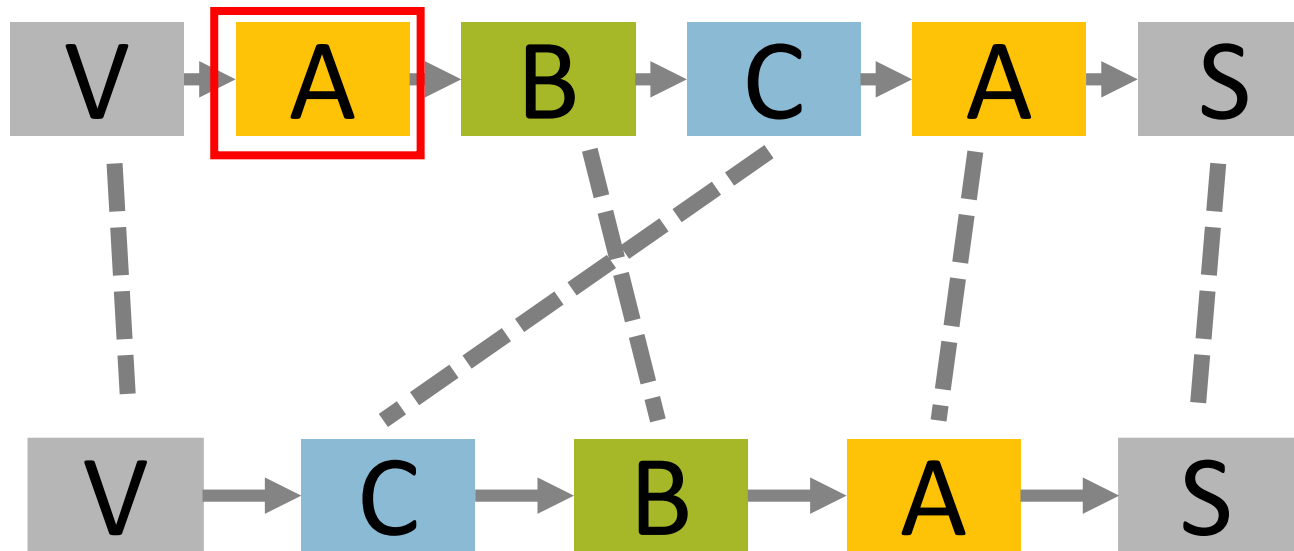
- (1) Differences in neighbors
- (2) Differences in structure



Quantify Similarity of Mapped Events

... based on “structural context of events”

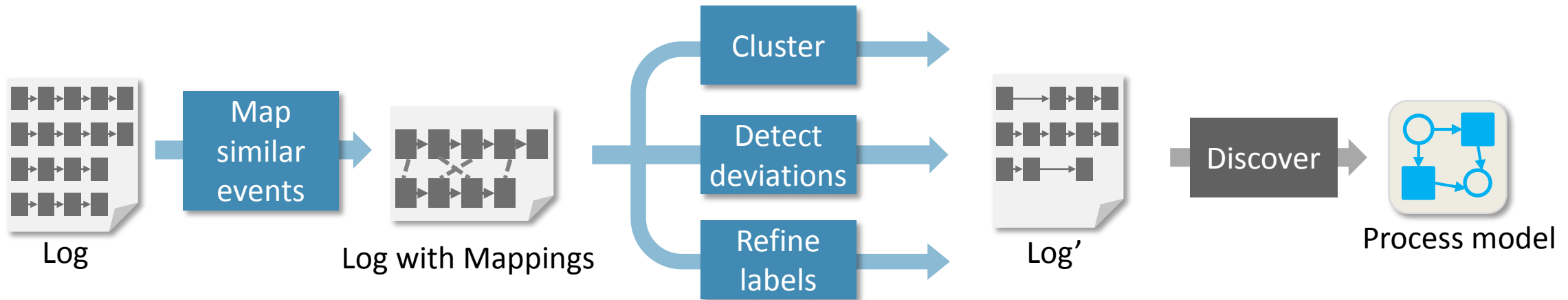
- (1) Differences in neighbors
- (2) Differences in structure
- (3) #Dissimilar events



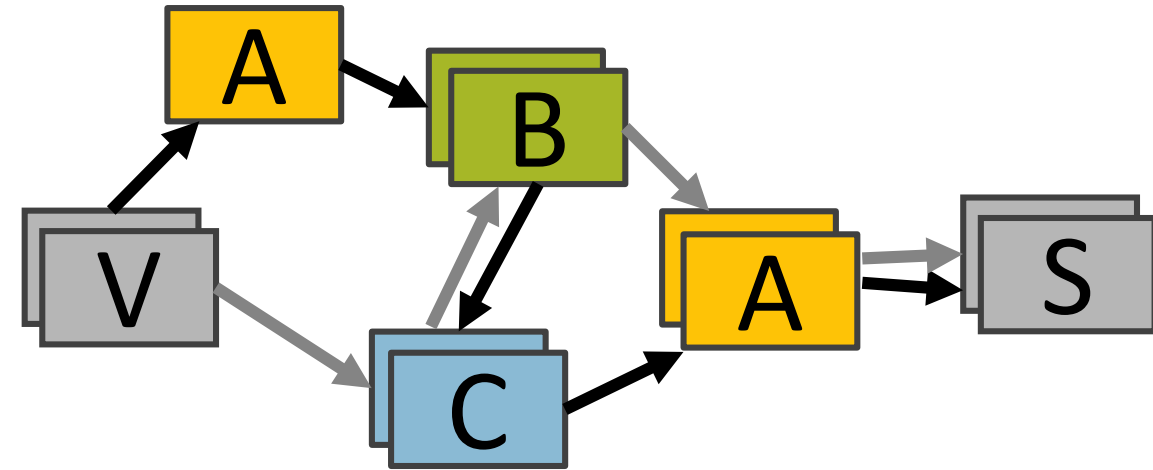
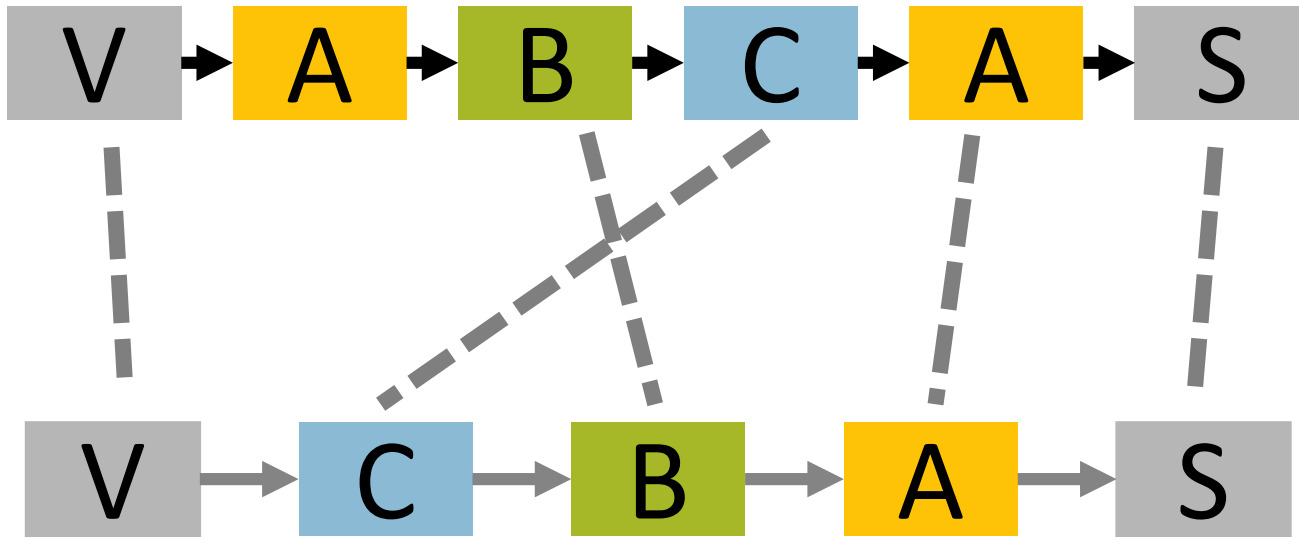
Total Cost = 6+10+8+5+3 +1

There is an algorithm to compute an optimal mapping :

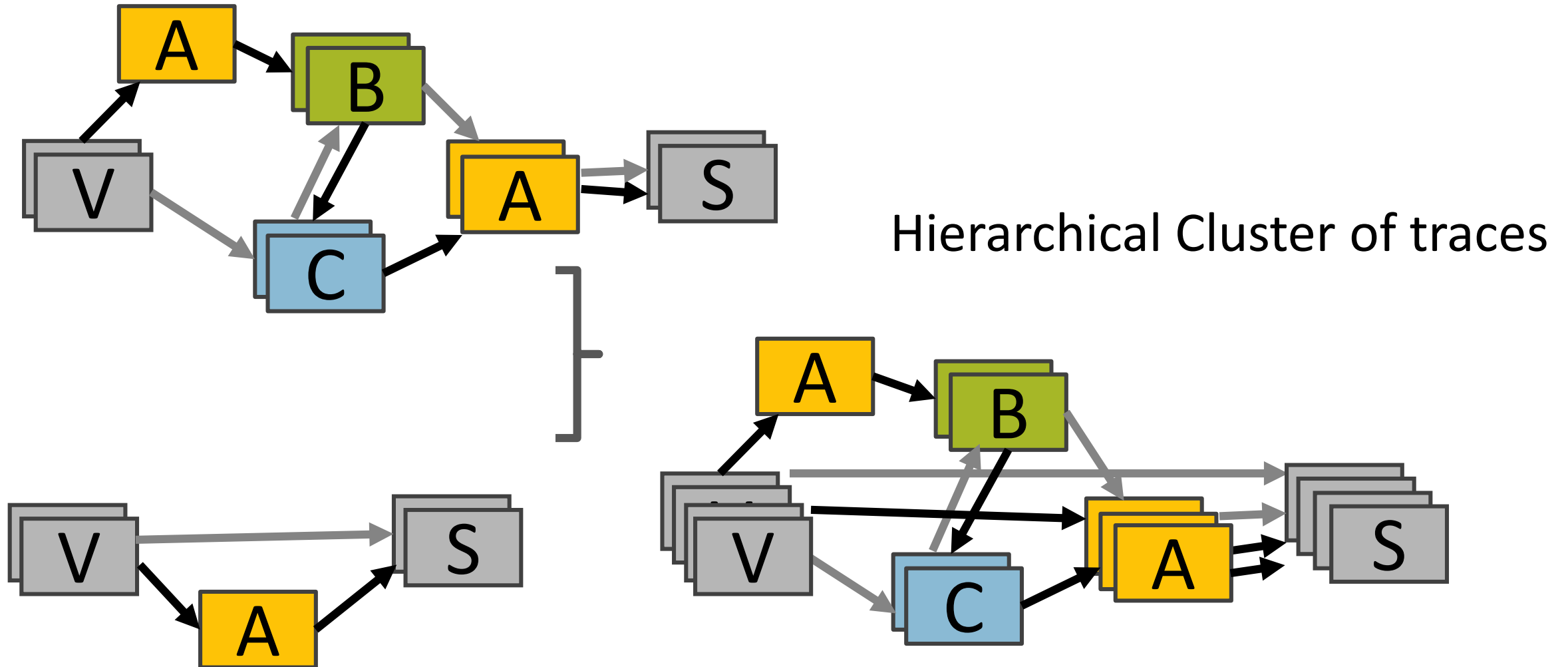
Proposing Approach



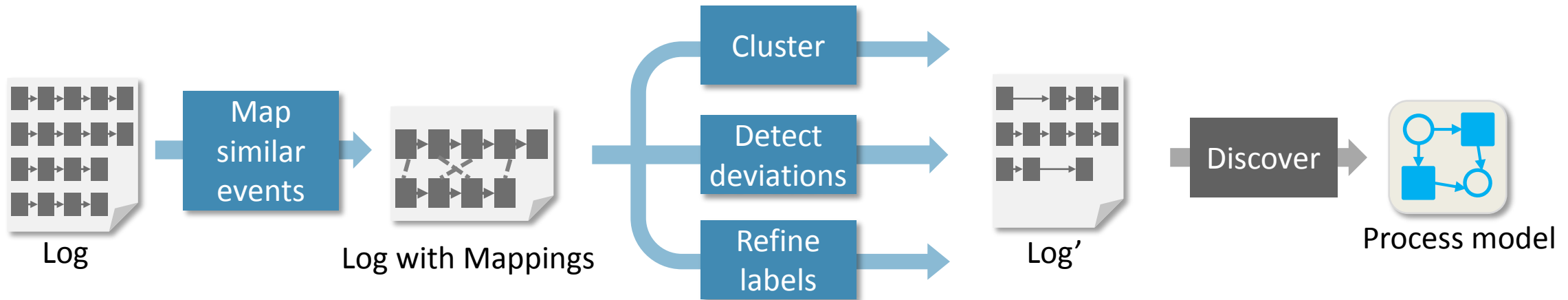
Cluster Traces Using Mappings and Fusion



Cluster Traces Using Mappings and Fusion



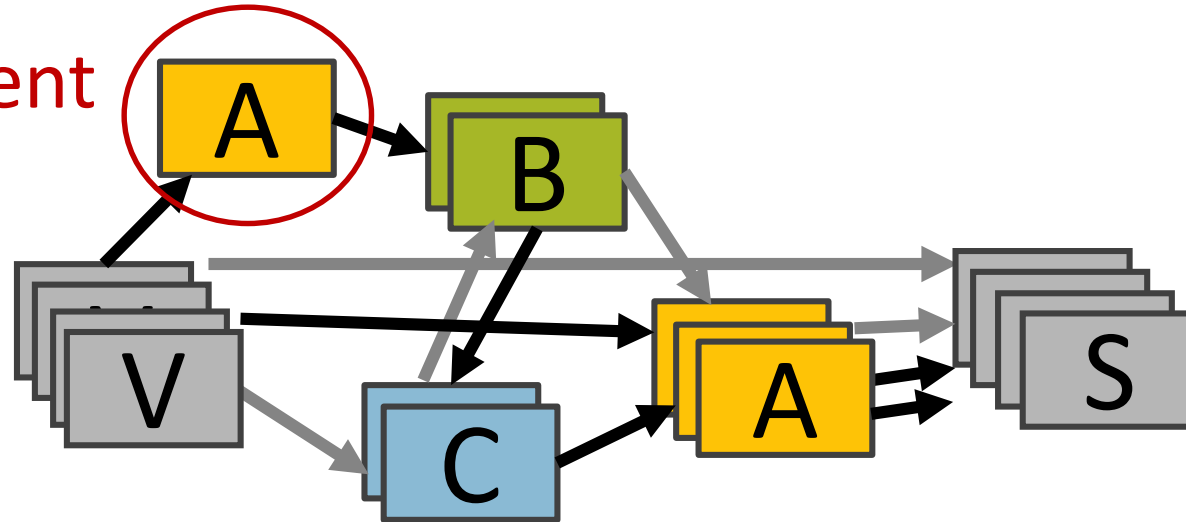
Proposing Approach



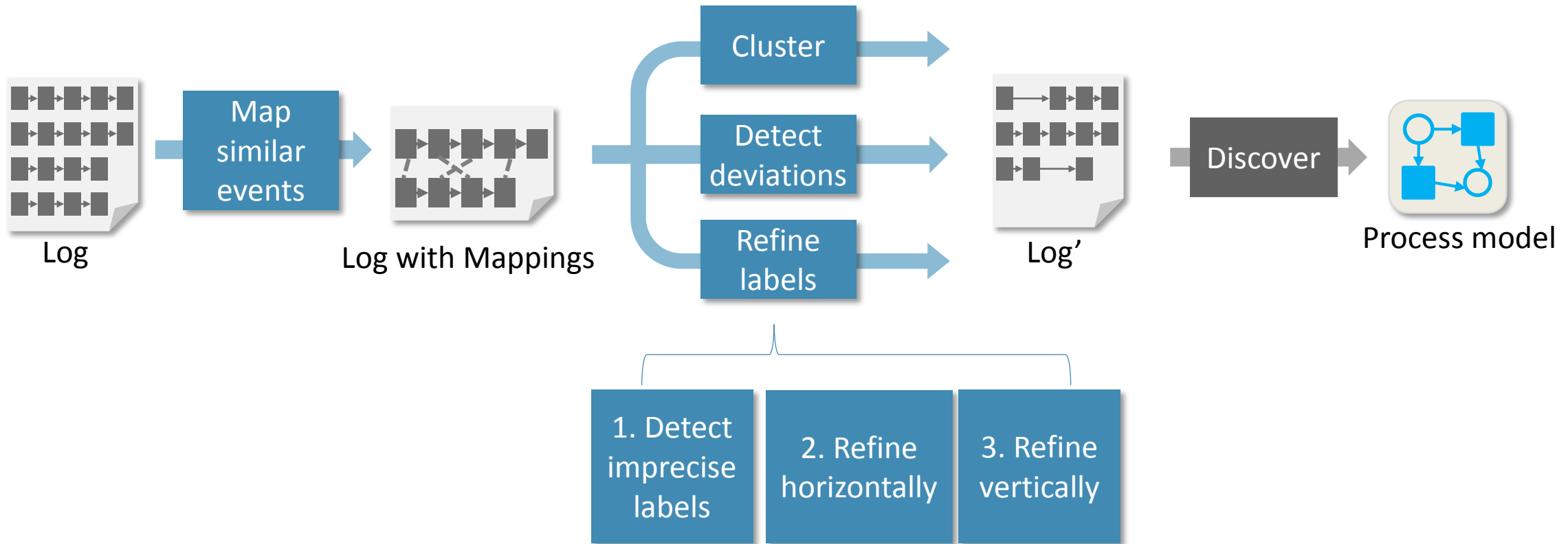
Detect Deviation Using Mappings

Threshold: “Deviating event” if $< T\%$ of cases

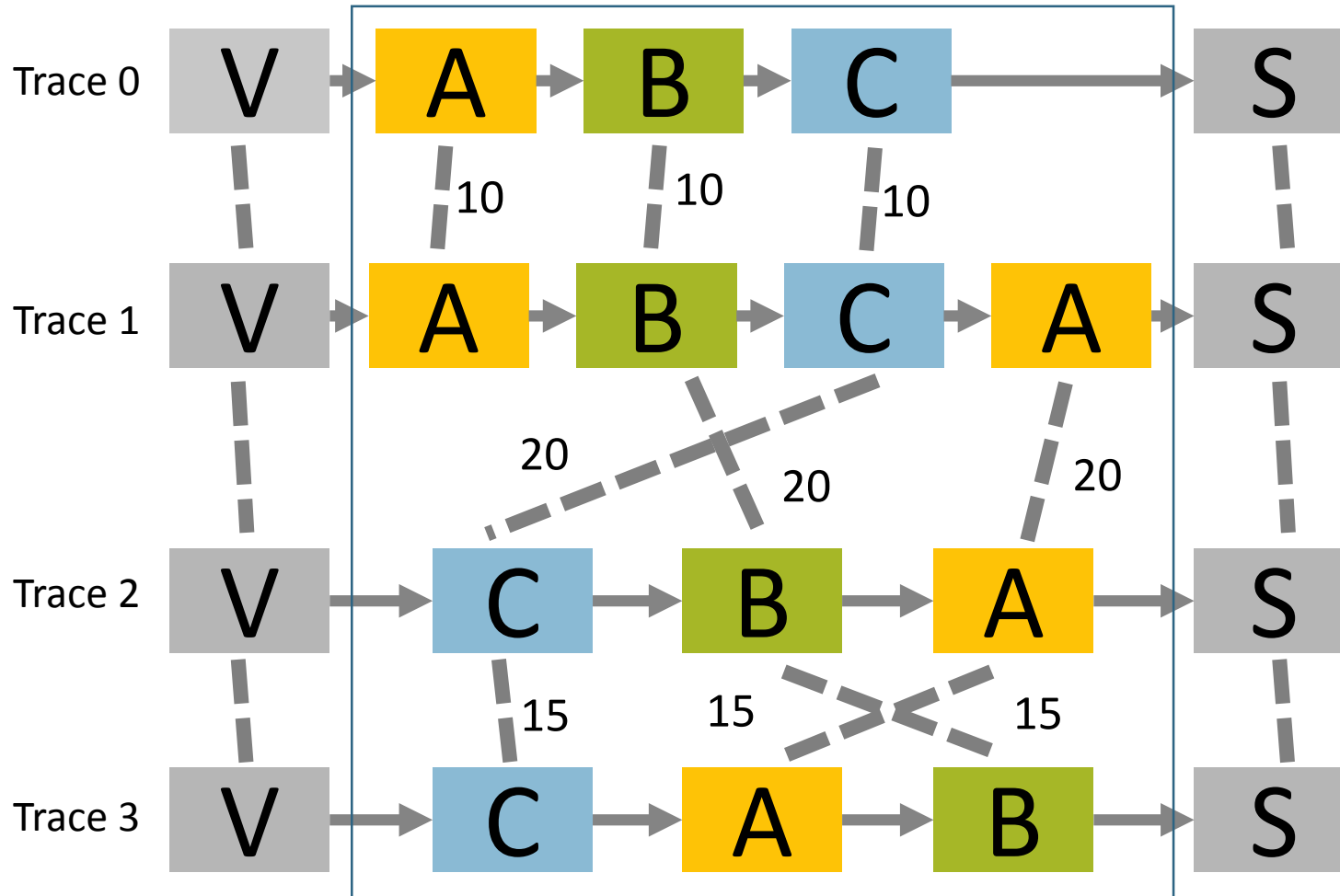
Deviating event
for $T = 30\%$



Proposing Approach



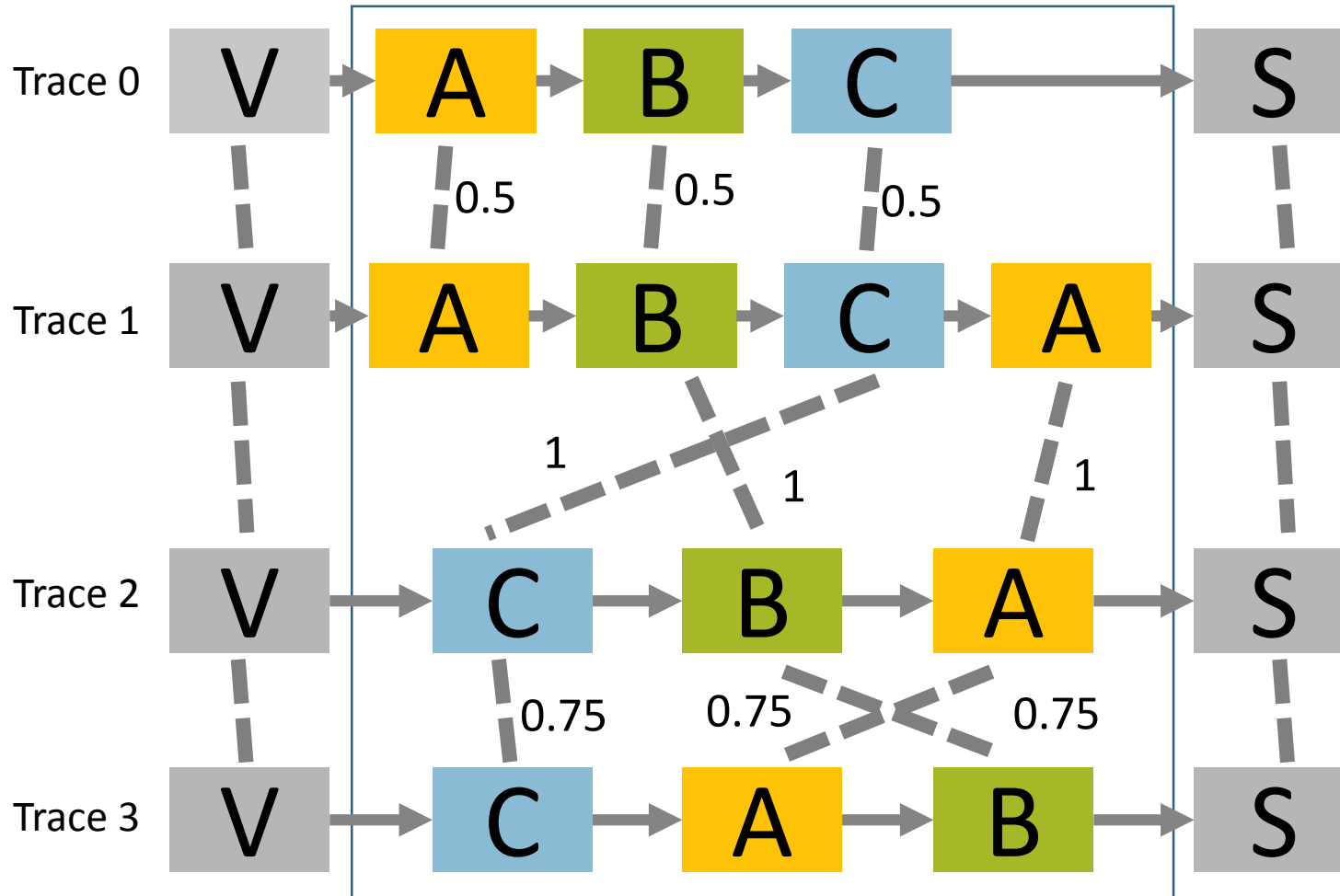
Refine Labels Horizontally



a) Normalize costs w.r.t. maximal cost seen in the log

Imprecise label candidates {A, B, C}

Refine Labels Horizontally



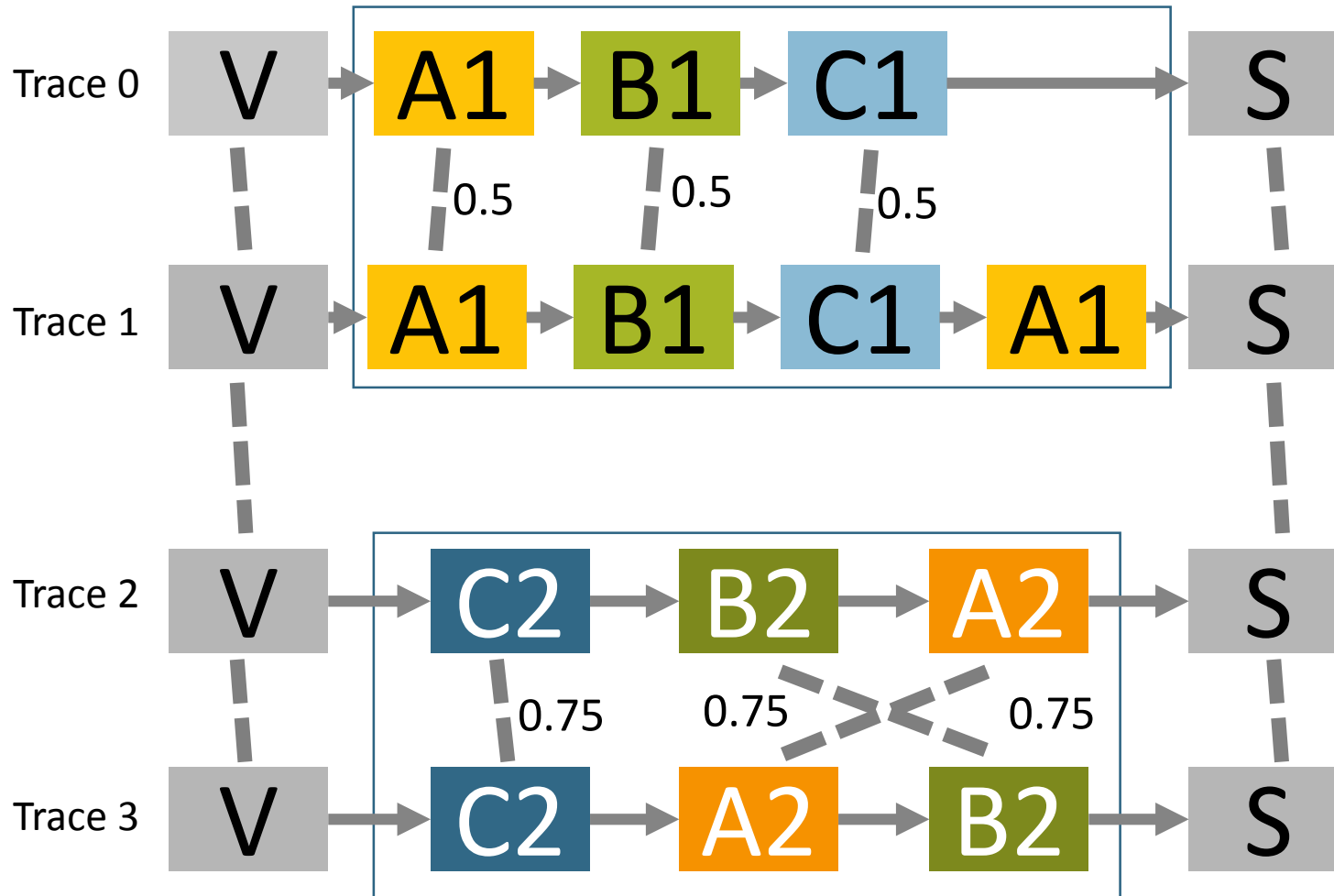
a) Normalize costs w.r.t. maximal cost seen in the log

b) Set variant threshold, say, 0.8

c) Remove edges if cost > variant threshold

Imprecise label candidates {A, B, C}

Refine Labels Horizontally

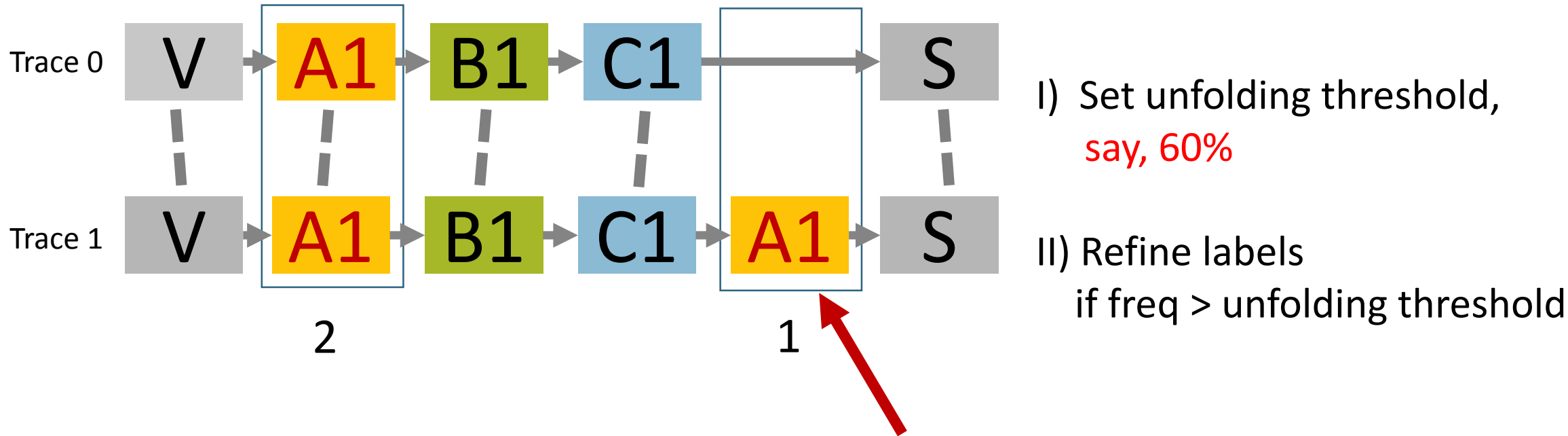


a) Normalize costs w.r.t. maximal cost seen in the log

b) Set variant threshold say, 0.8

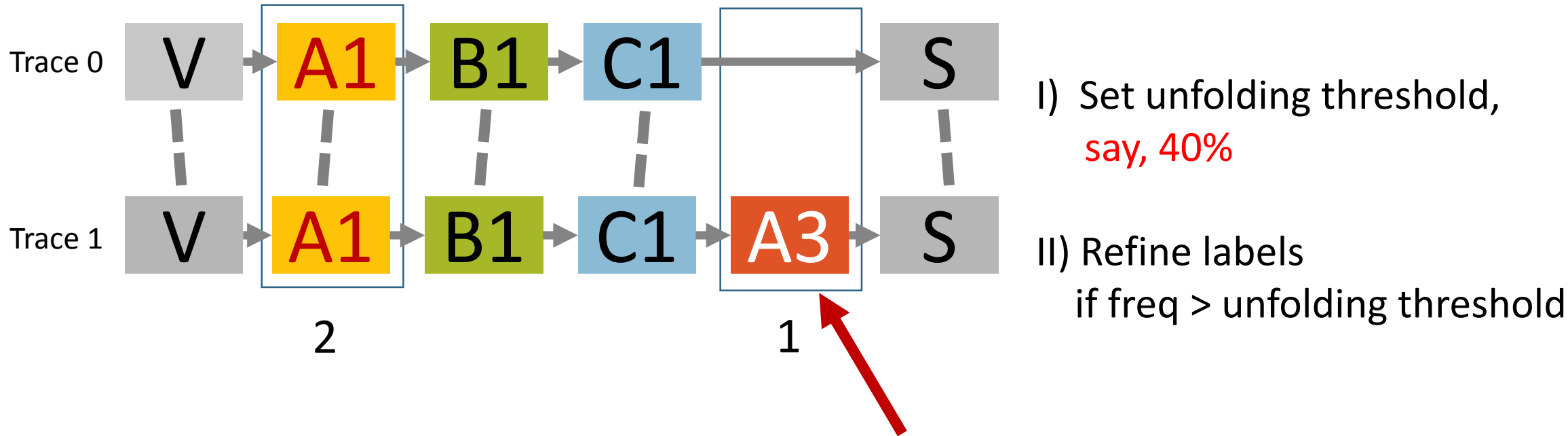
c) Remove edges if cost > variant threshold

Refine Labels Vertically



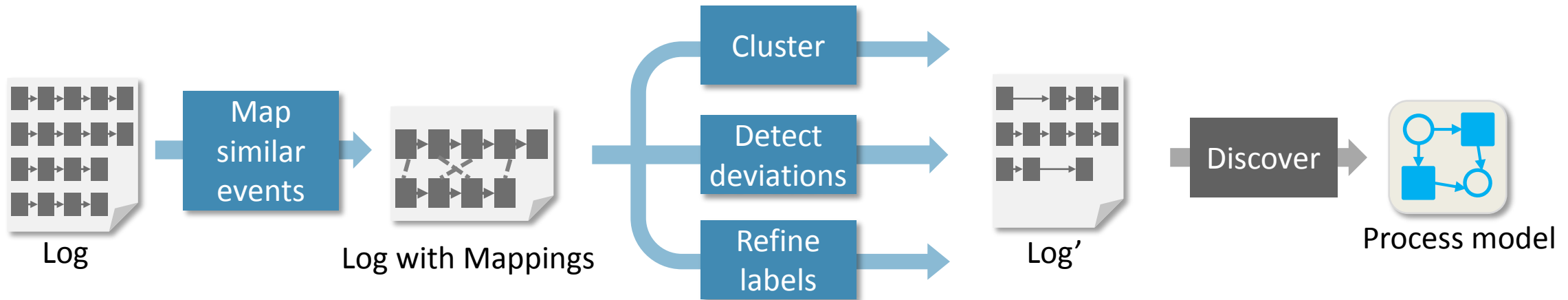
Is this A1 a 2nd iteration of a loop?

Refine Labels Vertically



Or is this a different task?

Proposing Approach



IV. Demo and Discussion

- Plugin “Log to Model Explorer”

The screenshot displays the ProM 6 RepEGraphs interface. On the left, a tree view under 'RepEGraphs' shows a 'Root' folder containing several clusters of traces, such as 'Cluster 1994 (53 trcs)' and 'Cluster 1992 (52 trcs)'. A central panel shows a process model with various activities like 'Create Fine', 'Send Fine', 'Insert Fine Notification', 'Appeal to Judge', 'Add penalty', 'Send Appeal to Prefecture', 'Receive Result Appeal from Prefecture', 'Notify Result Appeal to Offender', 'Send for Credit Collection', and 'Payment'. The model is annotated with execution traces. On the right, a 'Cluster' panel lists 'Imprecise Label Candidates' and provides controls for refining labels, including sliders for 'Set Max. Cost' and 'Set Unfolding Threshold', and buttons for 'Export model' and 'Export log'. A 'Set visualization type' dropdown is set to 'IM'. The interface is titled 'ProM UITopia' and 'ProM 6'.

Cluster Traces

Refine Labels

Discover and Visualize

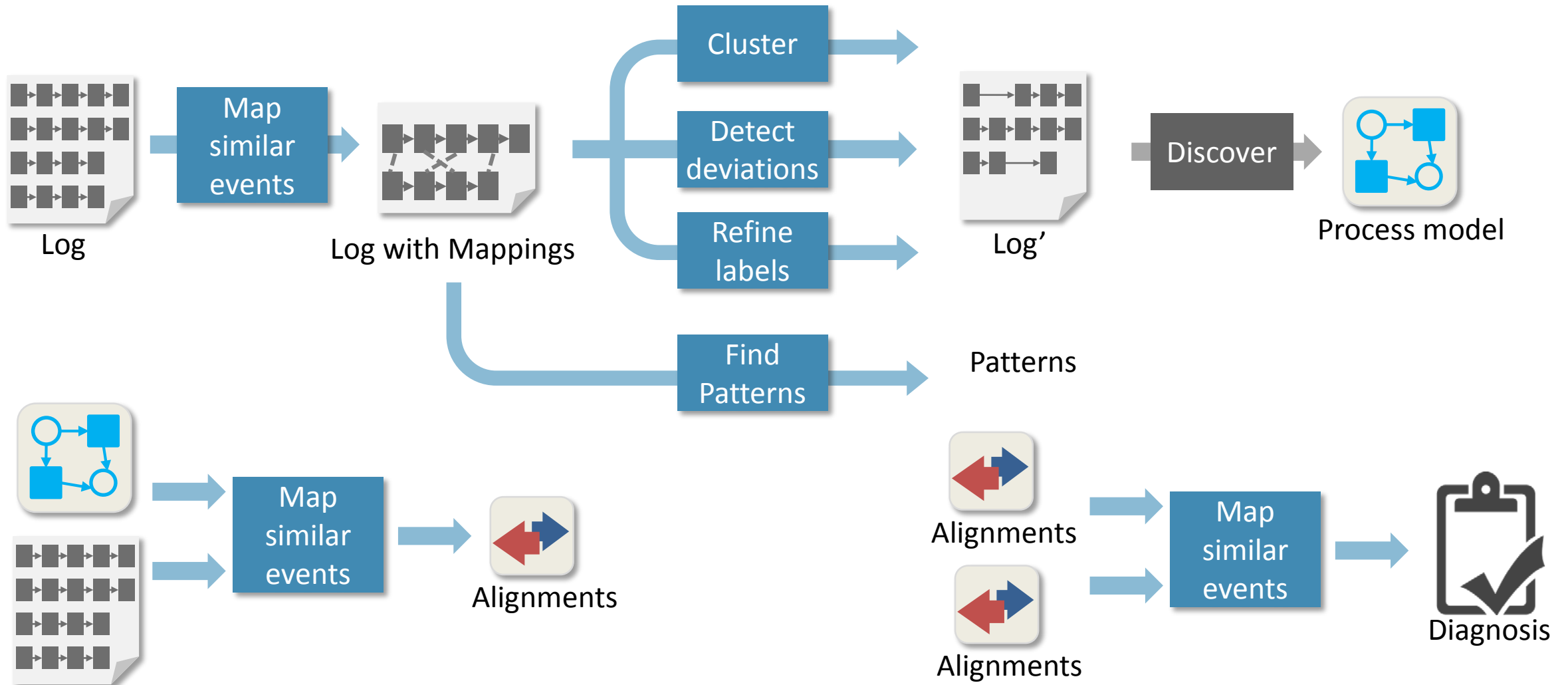
Detect Deviations

Limitations?

- Implementation inefficient
- Do not know what is happening in the background
- Can not handle large loops and large parallel branches

- Clustering
 - Do not know why these clusters
- Filtering
 - Do not support finding missing events (yet)
- Refining labels
 - ...

The Bigger Picture?



Questions and Feedback?