

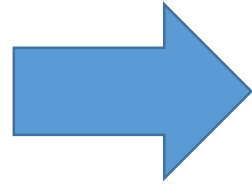
Database Process Mining

Alifah Syamsiyah
Boudewijn F. van Dongen
Wil M.P. van der Aalst

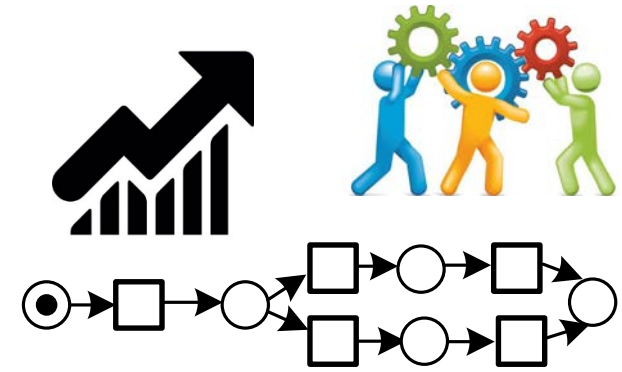
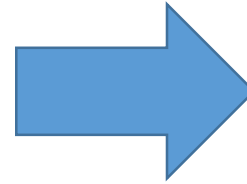
1



Event Data

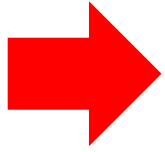


Process Mining



Discovery Results

Provide **instant discovery**
Enable **discovery based on time window**
Store **richer information for alignment**



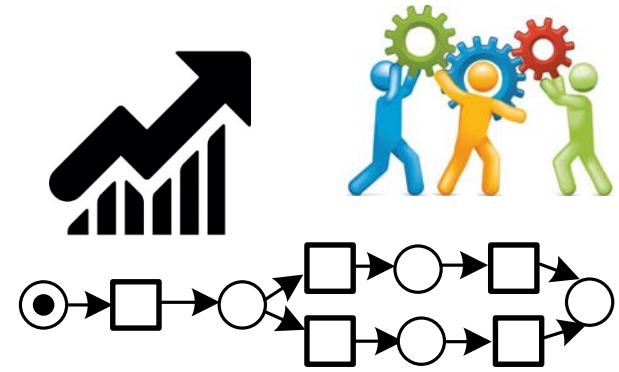
Instant Discovery

Discovery based on Time Window

Alignment

Experimental Results

Conclusion



Intermediate Structure

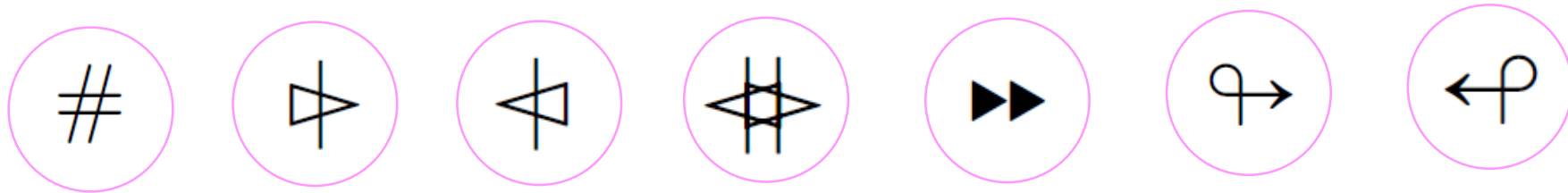
DFR

Declare relations

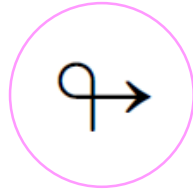
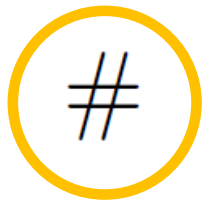
Prefix automaton



Incrementally update the intermediate structure



What are these relations?
How to incrementally updates these relations?



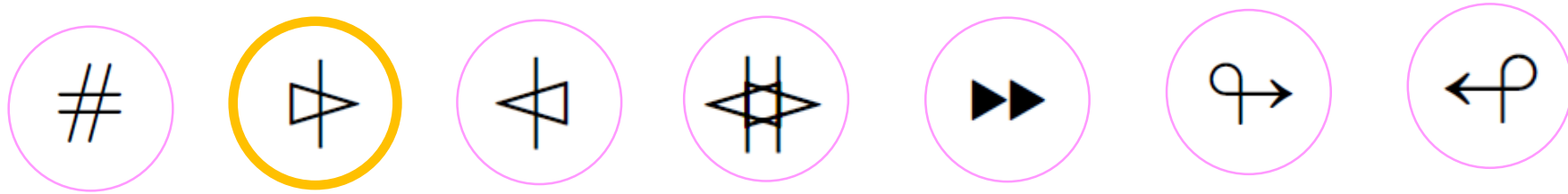
Occurrence $\#(a)$

Counts the occurrences of activity a

$L = \{ \langle a, b, a, a, c, a \rangle, \langle a, b, a, a, c, a, d \rangle \}$

$\#(a) = 8$

$\Sigma = \{a, b, c, d\}$



No Following $\triangleright(a,b)$

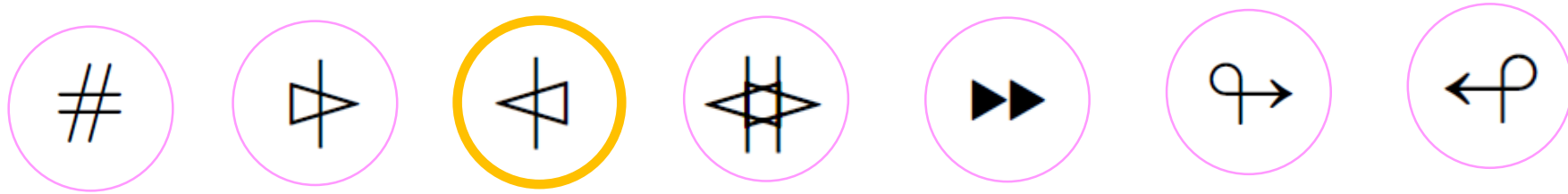
Counts the occurrences of activity a with no following b

$L = \{ \langle \boxed{a}, b, \boxed{a}, \boxed{a}, c, \boxed{a} \rangle, \langle a, b, a, a, c, \boxed{a}, d \rangle \}$

$$\triangleright(a,d) = 4$$

$$\triangleright(a,c) = 2$$

$\Sigma = \{a, b, c, d\}$



No Preceding $\triangleleft(a,b)$

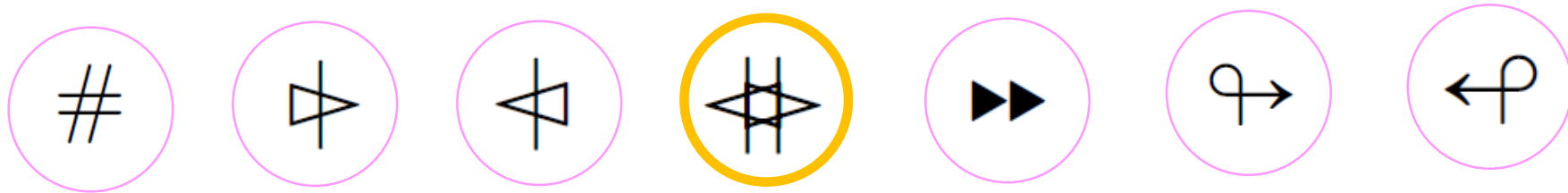
Counts the occurrences of activity a with no preceding b

$L = \{ \langle a, b, a, a, c, a \rangle, \langle a, b, a, a, c, a, d \rangle \}$

$$\triangleleft(a,b) = 2$$

$$\triangleleft(a,c) = 6$$

$\Sigma = \{a, b, c, d\}$



No Co-Occurring $\nabla (a,b)$

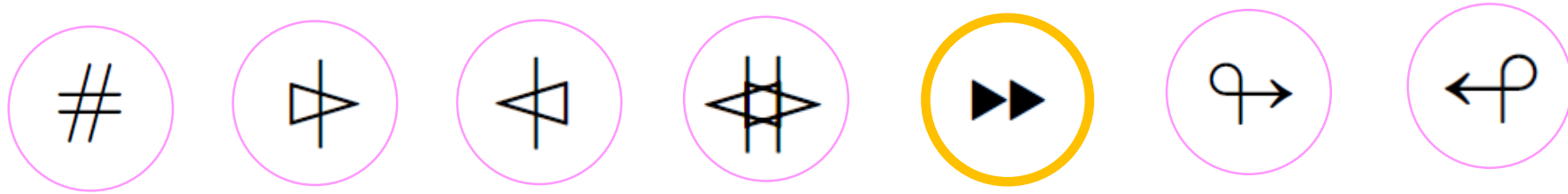
Counts the occurrences of activity a with no co-occurring b

$L = \{ \langle a, b, a, a, c, a \rangle, \langle a, b, a, a, c, a, d \rangle \}$

$\nabla (a,c) = 0$

$\nabla (a,d) = 4$

$\Sigma = \{a, b, c, d\}$



Directly Follows Relation ▶▶ (a,b)

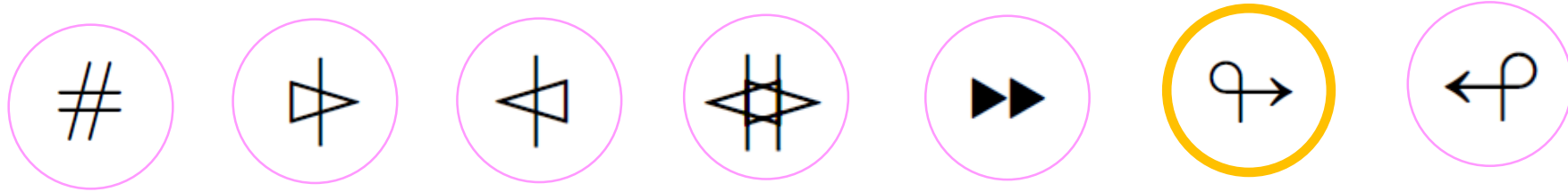
Counts the number of times activity a is directly followed by activity b

$L = \{ \langle \boxed{a}, \boxed{b}, a, \boxed{a}, \boxed{c}, a \rangle, \langle \boxed{a}, \boxed{b}, a, \boxed{a}, \boxed{c}, a, d \rangle \}$

▶▶ (a,b) = 2

▶▶ (a,c) = 2

$\Sigma = \{a, b, c, d\}$



Repetition $\varrho \rightarrow (a,b)$

Counts how many times after an occurrence of a , a repeats until the first occurrence of b .
 If no b occurs after a , then the repetitions after a are not counted.

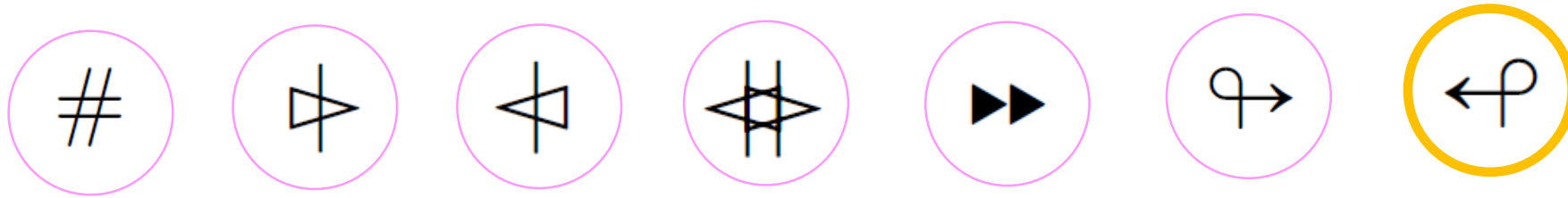
$L = \{ \langle a, b, \boxed{a}, \boxed{a}, c, a \rangle, \langle a, b, \boxed{a}, \boxed{a}, c, \boxed{a}, d \rangle \}$

$\Sigma = \{a, b, c, d\}$

$$\varrho \rightarrow (a,c) = 4$$

$$\varrho \rightarrow (a,b) = 0$$

$$\varrho \rightarrow (a,d) = 3$$



Repetition (Backwards) $\leftarrow \rho(a,b)$

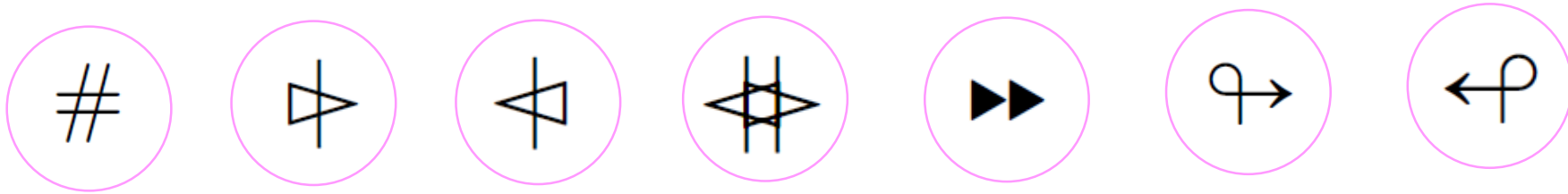
Is similar to $\rho \rightarrow$ but reading the trace backwards

$$L = \{ \langle a, b, \boxed{a}, \boxed{a}, c, a \rangle, \\ \langle a, b, \boxed{a}, \boxed{a}, c, a, d \rangle \}$$

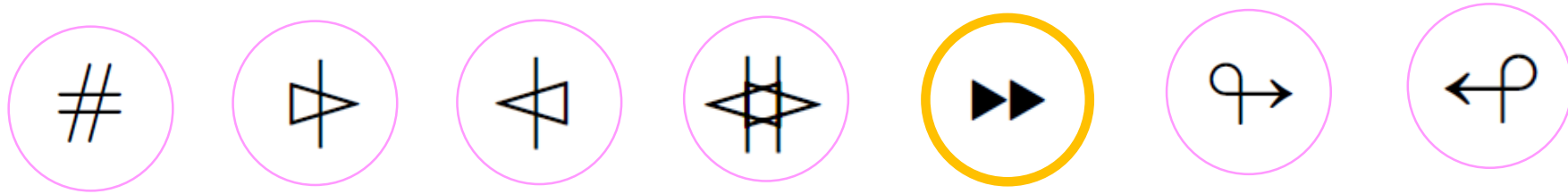
$$\leftarrow \rho(a,c) = 0$$

$$\leftarrow \rho(a,b) = 4$$

$$\Sigma = \{a, b, c, d\}$$

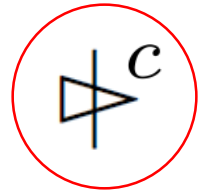
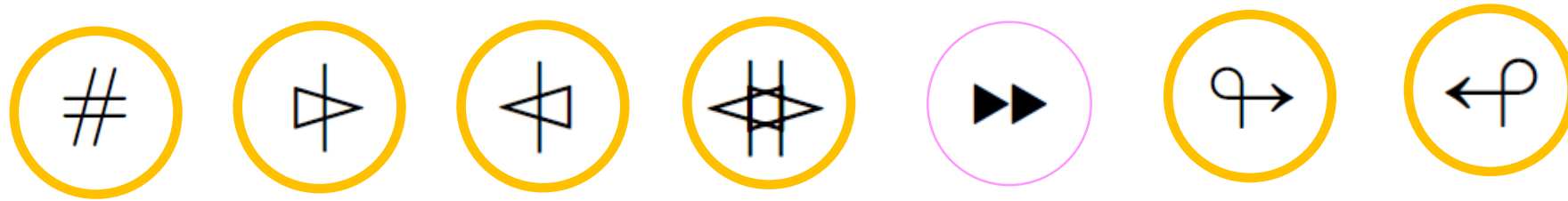


What are these relations? 
How to incrementally updates these relations?



$L = \{ \langle a, b, a, a, c, a \rangle \}$

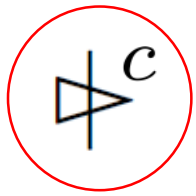
Activity 1	Activity 2	Frequency
a	b	1
b	a	1
a	a	1
a	c	1
c	a	1



Controller function $\triangleright^c(a,b)$

If $a \neq b$: counts the occurrences of activity a with no following b in trace t .

If $a = b$: counts the occurrence of a in trace t .



$\Sigma = \{a, b, c, d\}$

For a new activity x :

- **add 1** to the row of x
- **reset** the column of x except (x,x)

	a	b	c	d
a				
b				
c				
d				

$L = \{ \langle a, b, a, a, c, a, d \rangle \}$

Occurrence

(a) = 4

(c) = 1

(b) = 1

(d) = 1

∇^c	a	b	c	d
a				
b				
c				
d				

$L = \{ \langle a, b, a, a, c, a, d \rangle \}$

No Following

$$\nabla (a,b) = 2$$

$$\nabla (a,c) = 0$$

∇^c	a	b	c	d
a				
b				
c				
d				

$$L = \{ \langle a, b, a, a, c, a, d \rangle \}$$

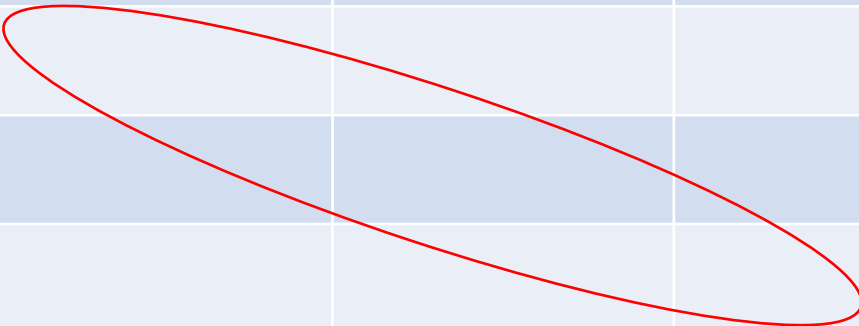
No Preceding

$$\nabla_{L'}(a,b) = 1$$

$$\nabla_{L'}(a,d) = 1$$

$$\nabla_{L'}(a,c) = 1$$

∇_c	a	b	c	d
a				
b				
c				
d				



$$L = \{ \langle a, \quad \rangle \}$$

No Co-Occurrence

$$\#_{L'}(a,b) = 0$$

$$\#_{L'}(a,d) = 1$$

$$\#_{L'}(a,c) = 1$$

$\forall c$		a	b	c	d
a					
b					
c					
d					

$$L = \{ \langle a, b, \quad \rangle \}$$

Repetition

$$\varphi \rightarrow_{L'}(a,c) = 2$$

3 occurrence of a with no following c

=

After an occurrence of a, a repeats twice until the first occurrence of c

$\forall c$	a	b	c	d
a				
b				
c				
d				

$$L = \{ \langle a, b, a, a, c, \quad \rangle \}$$

Repetition (Backwards)

$$\leftarrow \rho_{L'}(a,b) = 2$$

First occurrence

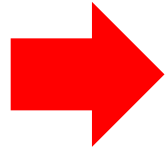
1st repetition

2nd repetition

$\forall c$	a	b	c	d
a				
b				
c				
d				

$$L = \{ \langle a, b, a, a, c, a, \rangle \}$$

Instant Discovery



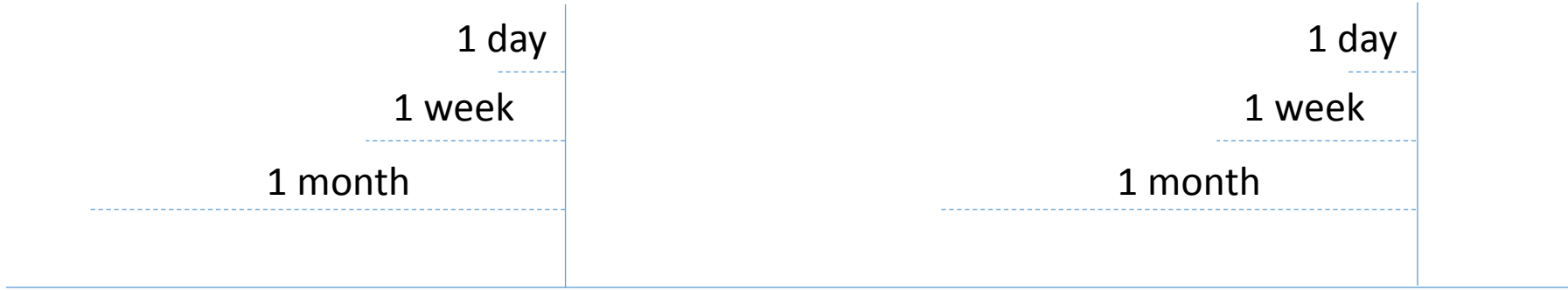
Discovery based on Time Window

Alignment

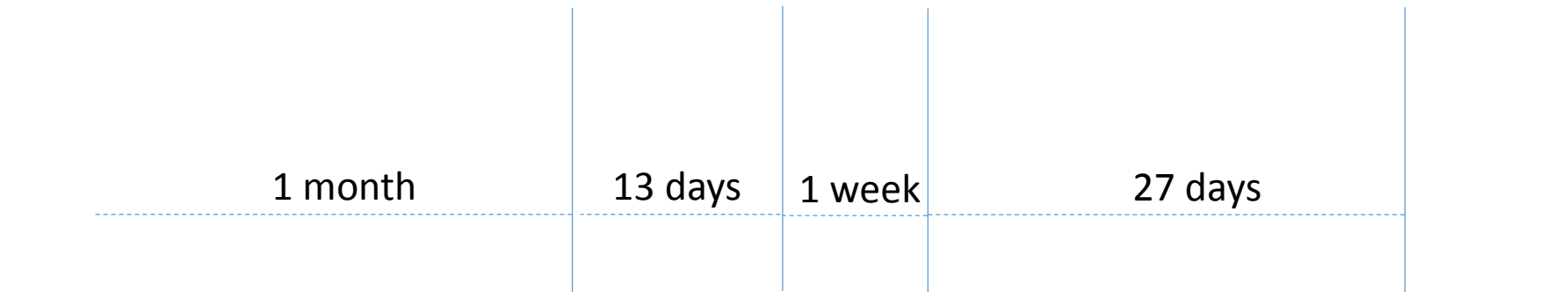
Experimental Results

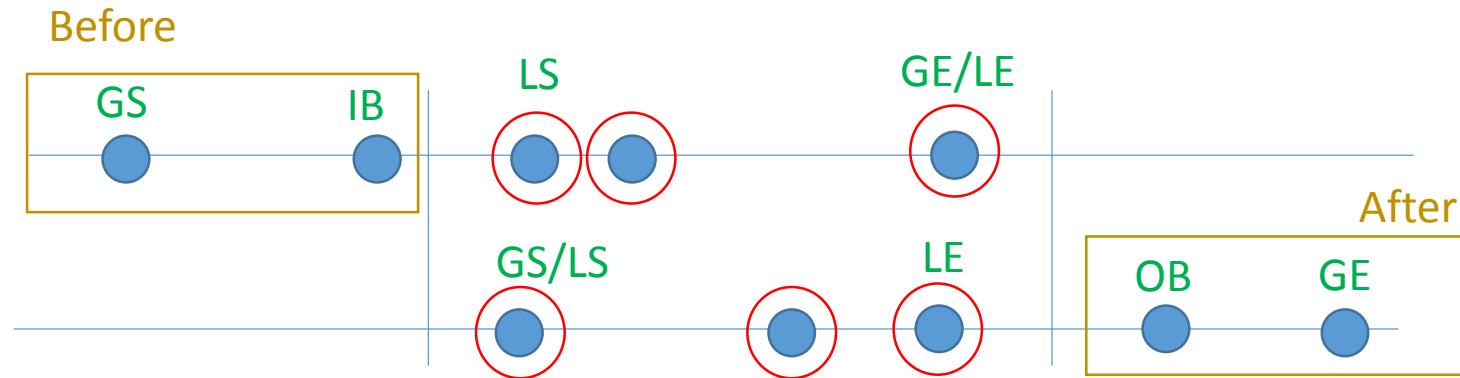
Conclusion

Sliding Window



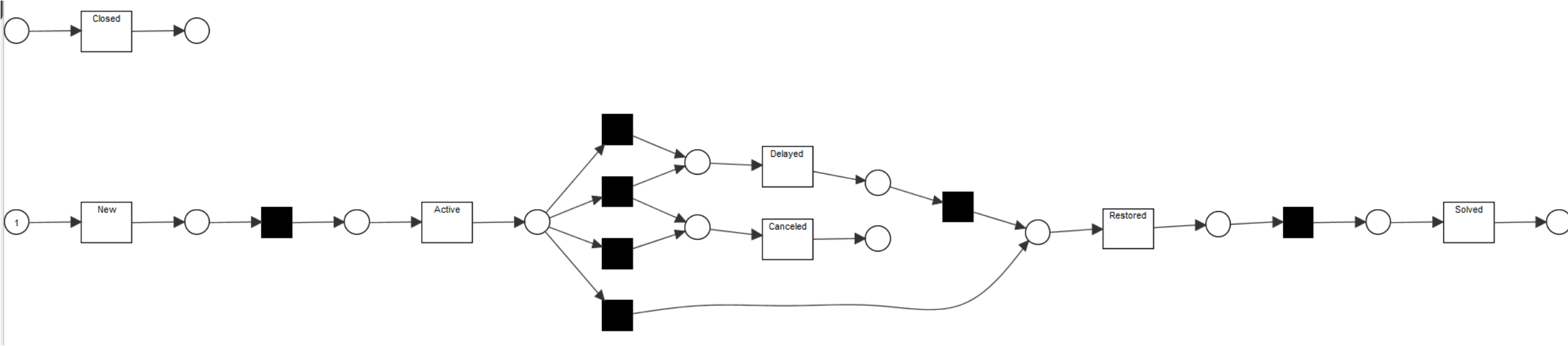
Chained Window





- = event always included
- IB = input boundary events
- OB = output boundary events
- GS = global start events
- GE = global end events
- LS = local start events
- LE = local end events

Telefonica Process Model of the last 11 days



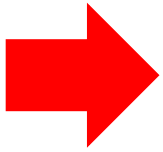
Instant Discovery

Discovery based on Time Window

Alignment

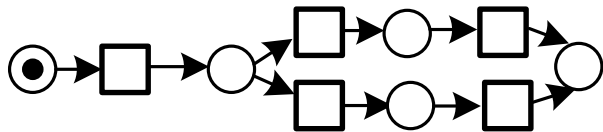
Experimental Results

Conclusion



Types of Alignment

	Optimal Alignment	Prefix Alignment	Partial Alignment
Initial Marking	Fixed for all traces	Fixed for all traces	Different from one trace to the others
Final Marking	Fixed for all traces	Different from one trace to the others	Different from one trace to the others



a		b c d		e
a b		d		f
		a d		c e

Motivation

- The existing technique is not able to give reliable fitness value for partial traces. Incomplete traces, for example, are treated as deviations.
- In a steady state process, we want to get stable fitness values in any window time
- In a big log, we want to have fitness value without taking into account all data
- We want to get an alarm as soon as a deviation happens (don't need to wait until a process finishes)



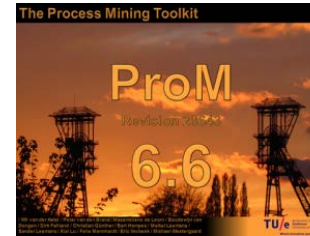
Store events, marking,
and last sequence

Sequence of events

Last marking



Current marking



Calculate fitness

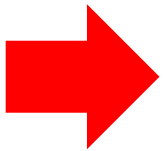
Instant Discovery

Discovery based on Time Window

Alignment

Experimental Results

Conclusion



Experimental Setup



BPI Challenge 2017 (31509 traces, 1202267 events)

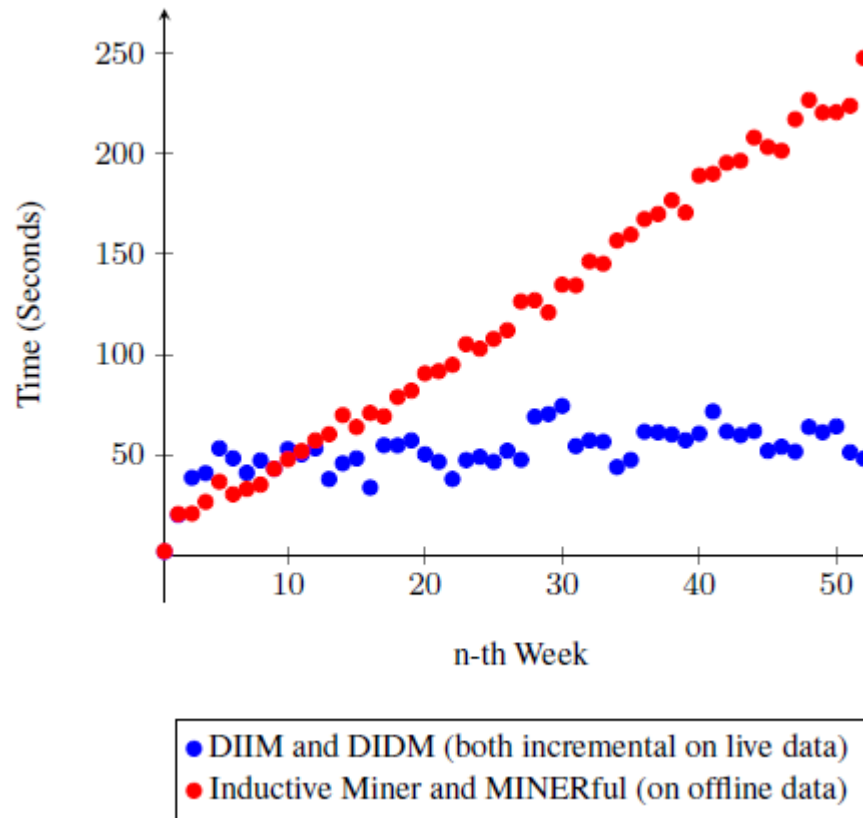
Weekly report on the last working day (every Friday)

In each report, both **procedural and declarative models** are mined

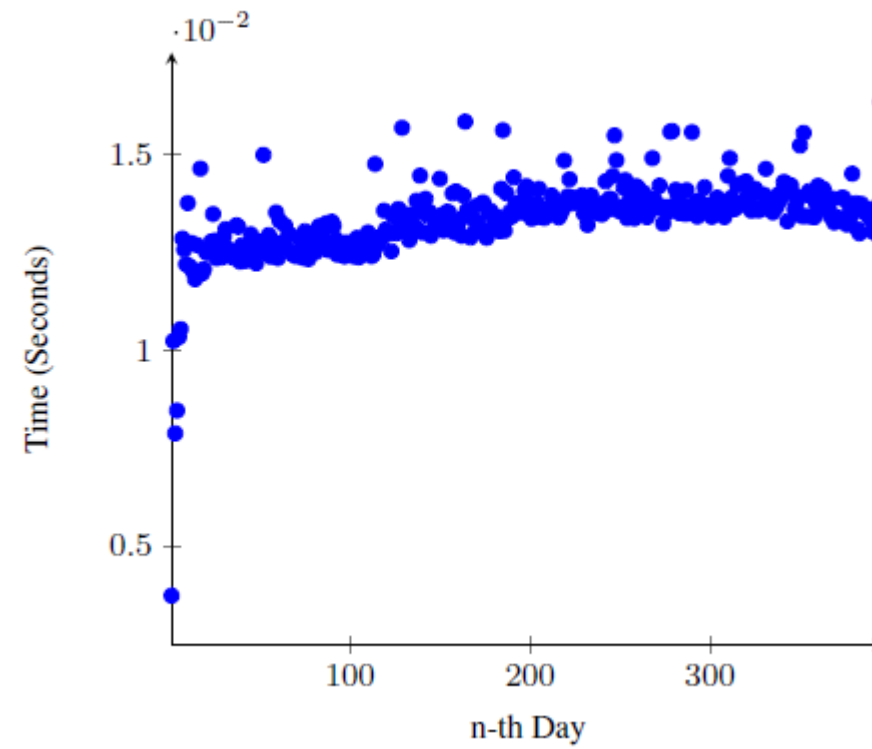
Implemented proposed approach:

- **DIIM** (Database-Incremental Inductive Miner)
- **DIDM** (Database-Incremental Declare Miner)

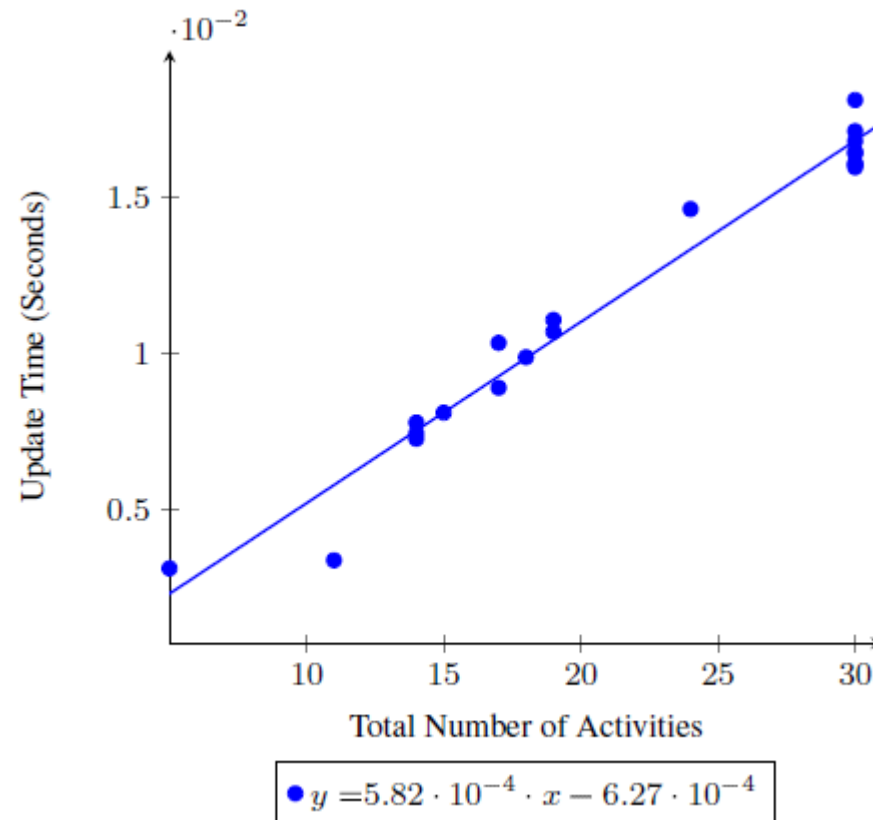
Recurrent Process Discovery Using DIIM and DIDM vs Traditional Inductive Miner and MINERful



Average Update Per Event



The Influences of Number of Activities to Update Time



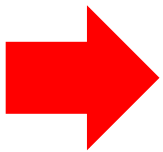
Instant Discovery

Discovery based on Time Window

Alignment

Experimental Results

Conclusion



Conclusion



- We can reduce analysis time by storing and computing intermediate structures in a persistent storage (e.g. database) and in an incremental way.
- We showed that our incremental implementations is constant for updating per event, while the retrieval and mining times are independent of the size of the underlying data.
- The use of database in process mining enables discovery based on time window and stores richer information for alignment.

Future Works



- Updating the intermediate structures **in batches of events**
- Keep the intermediate structures **live under removal of events.**