

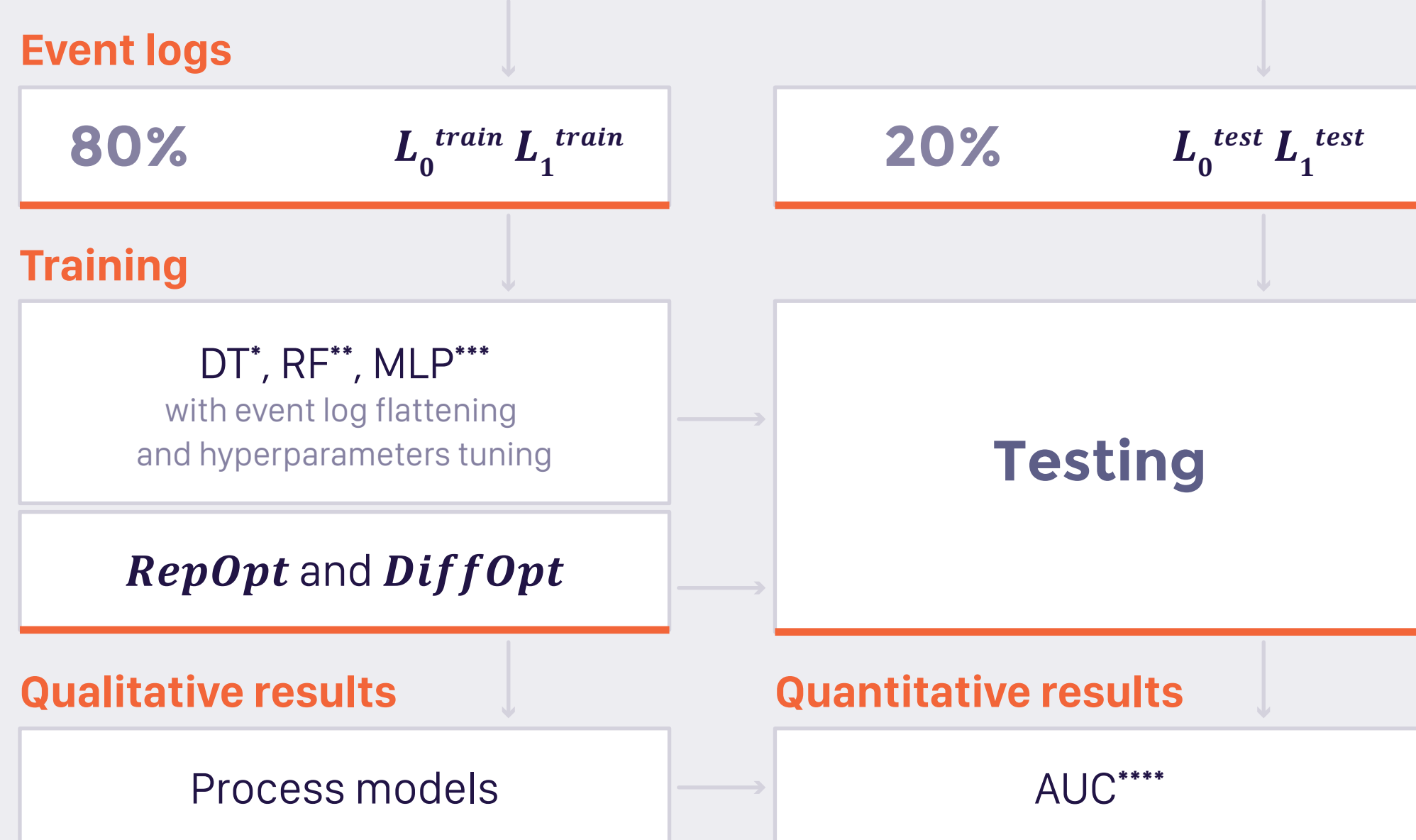
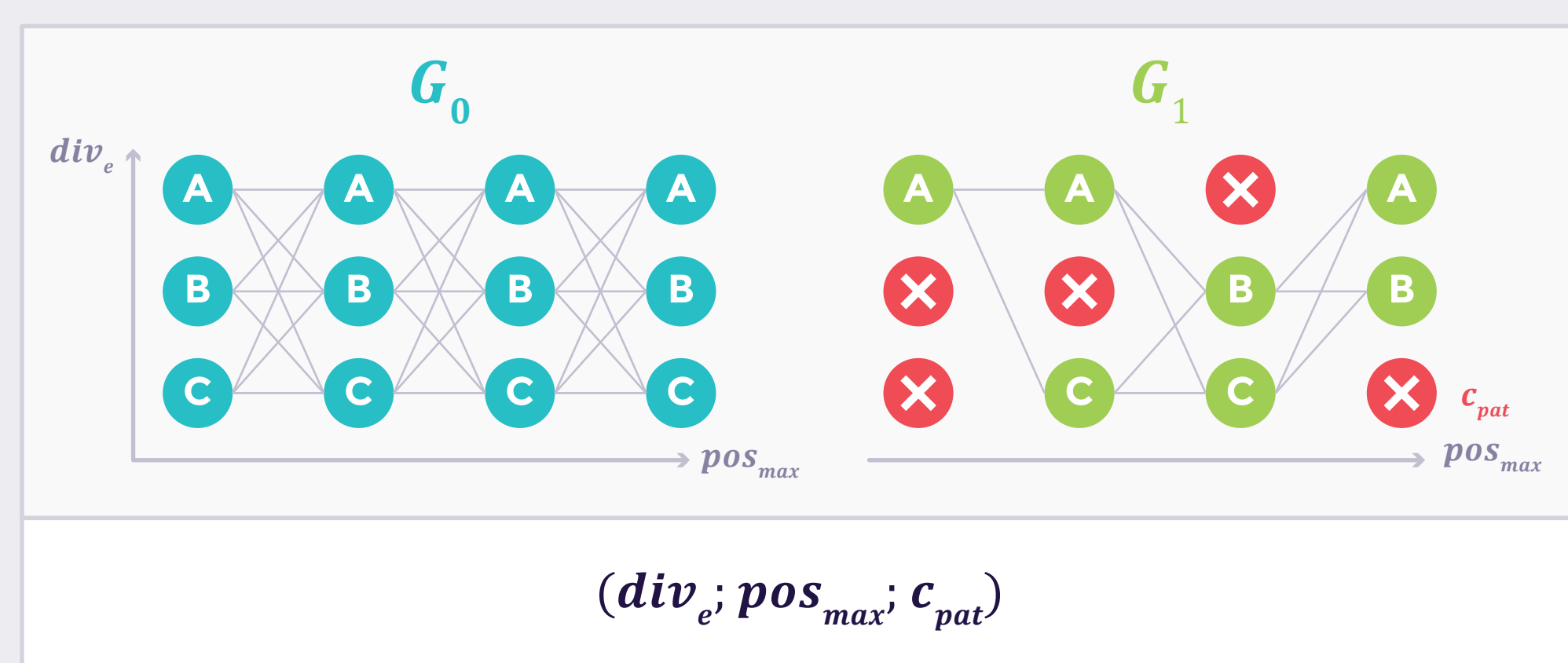
## Introduction

Event logs are a widespread type of data structure carrying information of time and ordering of events. As the complexity increases when time-dependent processes are considered, human understanding of predictive models is a lever for acceptability and practical deployment. We present here a new binary classification algorithm, created for event log data. Moreover, the proposed algorithm provides transparency by producing a process model to explain training results and future predictions. Transparency of predictive models is a current challenge, particularly in healthcare where deep learning has become state of the art<sup>[1]</sup>.

## Computational experiment

### Generation of traces with two graphs $G_0$ and $G_1$

After choosing a size configuration ( $pos_{max}$  for the length and  $div_e$  for the diversity), we create two graphs  $G_0$  and  $G_1$  with  $pos_{max}$  identical layers composed of  $div_e$  different nodes. Then, a proportion  $c_{pat}$  of shared patterns in  $G_1$  is deleted. Traces are then created by randomly crossing the graphs, forming event logs  $L_0$  and  $L_1$  (10 per parameters combination).



# Process model-based classification for event log data

## Proposed methodology

### Presentation

The problem addressed here is the training of a binary classifier on event log data of two classes  $L^{train}=(L_0^{train}, L_1^{train})$  of class 0 and 1.

The replayability  $R(PsM, \sigma) \in [0, 1]$  is a measure which quantifies the ability of a process model  $PsM$  to represent a trace<sup>[2]</sup>. The idea behind the proposed framework is the construction of a process model which well represents traces from event log  $L_1^{train}$  while less representing traces from  $L_0^{train}$ . Thus, it extracts discriminative patterns from  $L_1^{train}$ .

Two optimization functions **RepOpt** and **DiffOpt** can be used to construct the model, with a metaheuristic to find optimal process model.



### Event log and process model definitions

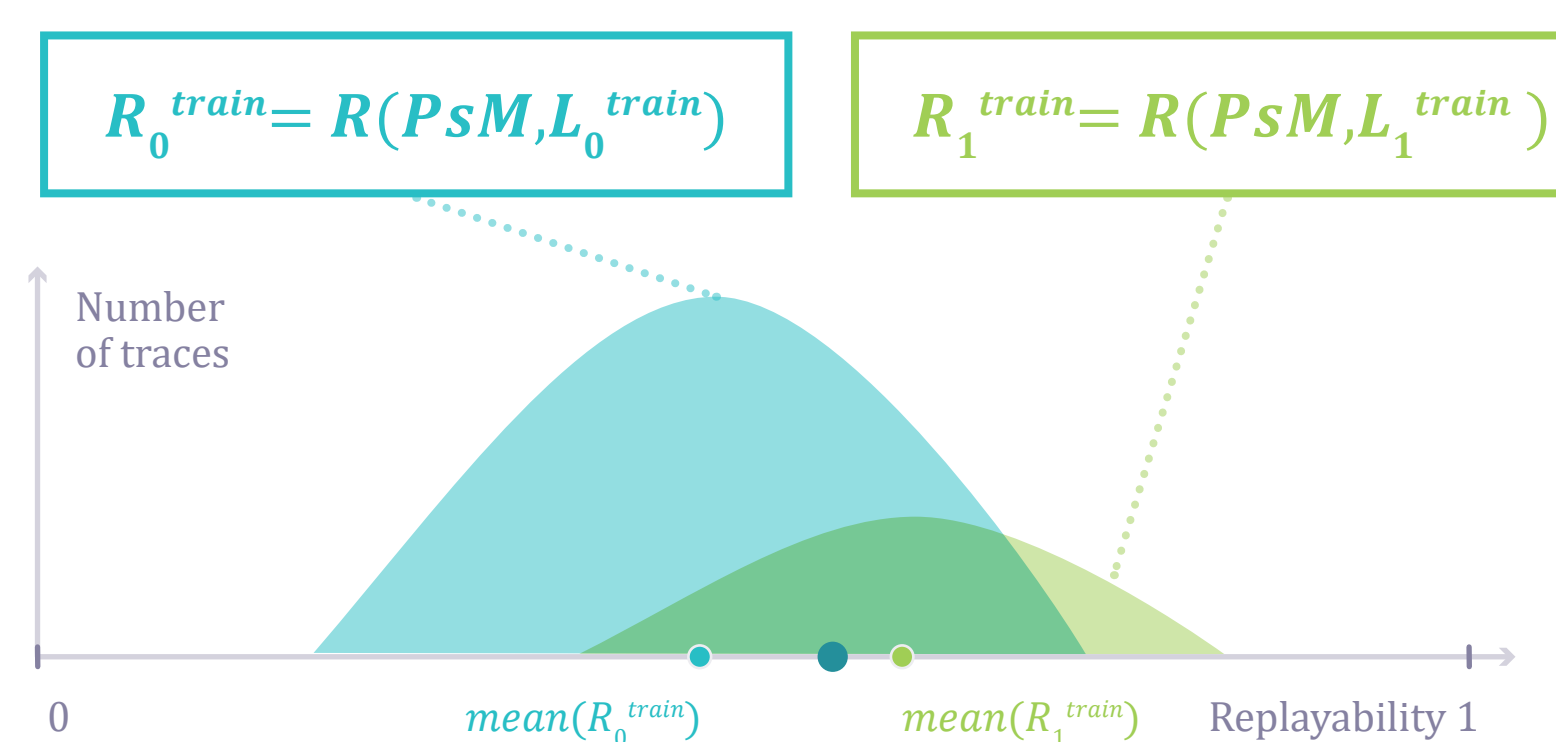
An **event log**  $L=\{\sigma_1, \sigma_2, \dots, \sigma_n\}$  regroups data in traces, each **trace**  $\sigma_i = e_1 e_2 \dots e_m$  being an ordered list of events, each **event**  $e_k=(a, t)$  having a **label**  $a \in A$  and a **time-stamp**  $t$ .

A **process model**  $PsM=(N, E, L, P)$  is defined as a four-tuple with a set of **nodes**  $N$  and **edges**  $E$ , a **label function**  $L$  and a **position function**  $P$ .  $L$  and  $P$  map each node  $n \in N$  respectively to a label  $a \in A$  and a position  $p \in \mathbb{N}^*$ . Each edge links a node to another one of **strictly higher position**.

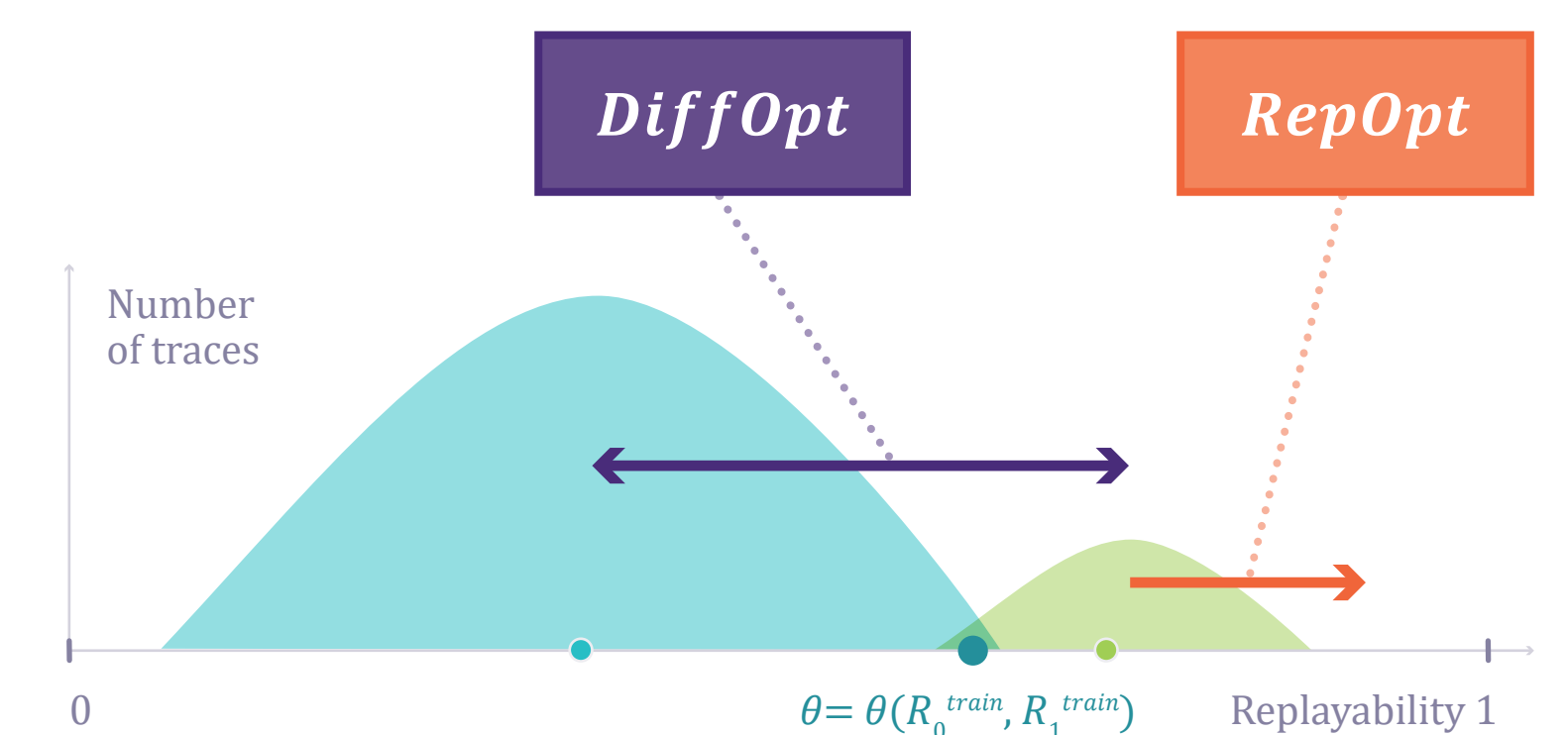
### The procedure of class prediction of a trace $\sigma$

- 1 Train  $PsM$  using **RepOpt** or **DiffOpt** on  $L^{train}=(L_0^{train}, L_1^{train})$  to obtain replayability distributions  $R_0^{train}=R(PsM, L_0^{train})$  and  $R_1^{train}=R(PsM, L_1^{train})$

#### Overlapped replayability distributions



#### Distinct after optimization



- 2 Set a threshold of separation  $\theta = \theta(R_0^{train}, R_1^{train})$ , e.g. using Gini impurity
- 3 Compute the replayability  $r_\sigma = R(PsM, \sigma)$
- 4 Predict the class of the trace  $\sigma$  by comparison between  $r_\sigma$  and  $\theta$

## Quantitative results

Results show that the proposed method using the optimization equation **DiffOpt** generally outperforms **RepOpt** and other methods.

Moreover, the AUC<sup>\*\*\*\*</sup> variability on the 10 replications per parameter combinations is lower for **DiffOpt**.

The performance gaps between **DiffOpt** and other methods, as well as performance variabilities increase with event logs complexity (from top to bottom in the table).

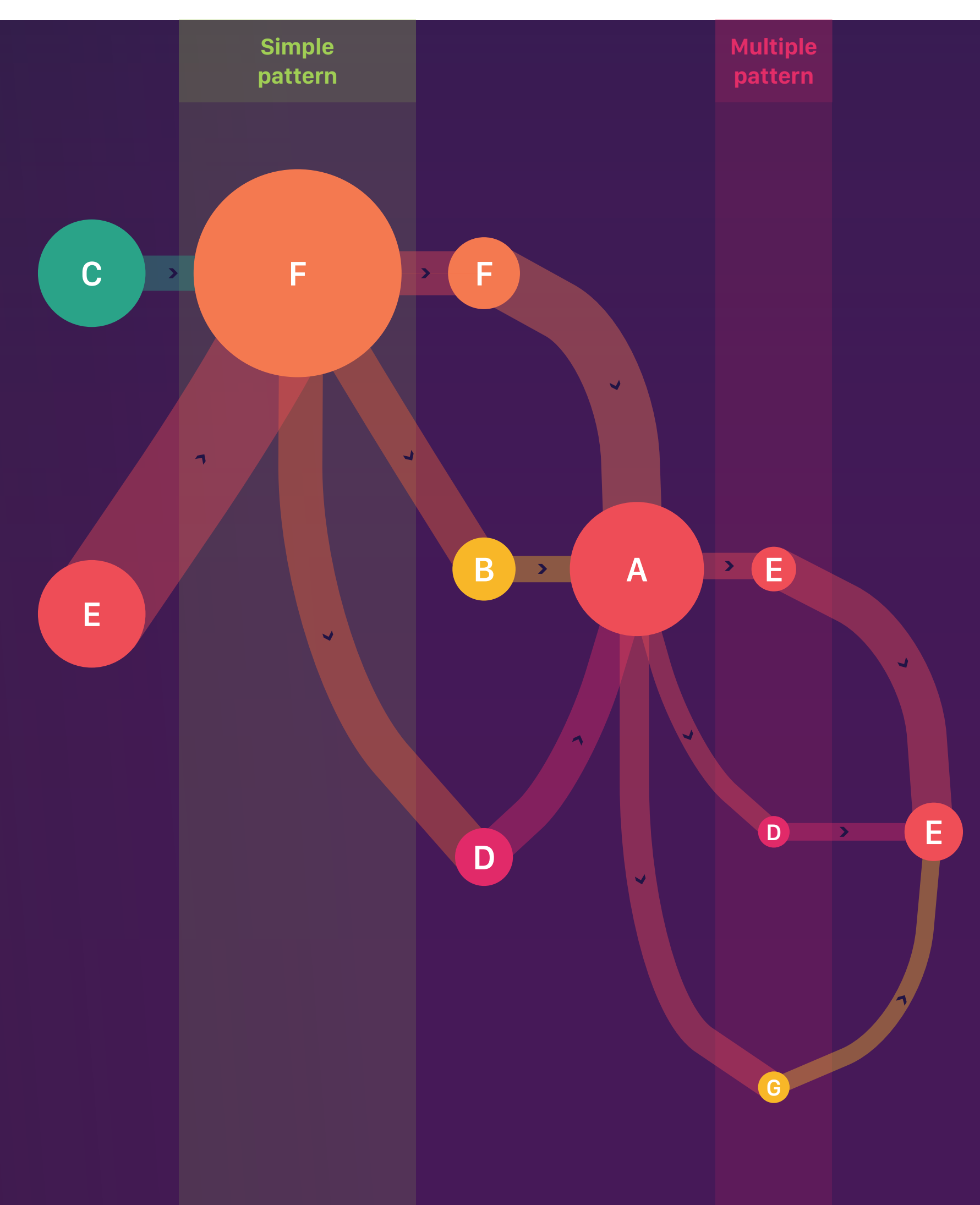
$c_{pat}$	$div_e$	$p_{max}$	DT*		RF**		MLP***		RepOpt		DiffOpt	
			AVG	STD	AVG	STD	AVG	STD	AVG	STD	AVG	STD
0.9	10	10	0.96	0.01	0.96	0.01	0.99	0.01	0.99	0.01	1	0
		25	0.95	0.01	0.95	0.01	1	0	0.99	0.02	0.99	0.02
		50	0.96	0.01	0.96	0.01	1	0	0.99	0.02	0.98	0.01
	50	10	0.95	0.02	0.95	0.02	0.97	0.02	0.98	0.01	1	0
		25	0.95	0.02	0.95	0.02	0.97	0.01	0.97	0.01	0.99	0
		50	0.95	0.03	0.95	0.03	0.98	0.02	0.97	0.01	0.99	0
0.75	10	10	0.95	0.01	0.95	0.01	0.96	0.01	0.98	0.01	0.99	0
		25	0.92	0.05	0.92	0.05	0.97	0.01	0.98	0.02	0.99	0
		50	0.9	0.07	0.9	0.07	0.97	0.02	0.97	0.01	0.99	0
	50	10	0.88	0.05	0.88	0.05	0.94	0.03	0.95	0.05	0.97	0.02
		25	0.89	0.04	0.9	0.04	0.95	0.04	0.95	0.06	0.96	0.02
		50	0.85	0.06	0.86	0.06	0.93	0.04	0.94	0.04	0.91	0.06
0.5	10	10	0.88	0.03	0.88	0.03	0.86	0.06	0.9	0.03	0.95	0.02
		25	0.87	0.04	0.85	0.06	0.87	0.05	0.91	0.04	0.94	0.01
		50	0.88	0.03	0.86	0.06	0.85	0.08	0.87	0.03	0.94	0.02
	100	10	0.77	0.1	0.78	0.06	0.86	0.03	0.87	0.05	0.93	0.02
		25	0.65	0.06	0.64	0.07	0.8	0.11	0.81	0.04	0.92	0.02
		50	0.64	0.07	0.64	0.06	0.85	0.05	0.72	0.05	0.86	0.07

## Qualitative results

A process model is used to discriminate the two classes using replayability. As a result, discriminative patterns extracted during the training procedure can be visualized.

The graph here highlights the specific patterns of  $G_1$  which are absent of  $G_0$  (for configuration  $c_{pat}=0.9$ ,  $div_e=10$ , and  $pos_{max}=10$ ).

Circles represent nodes of the model, and flux from circles represents edges. The size of nodes and edges are proportional to the number of traces represented.



## Conclusion

Across this study, we proposed a new binary classification algorithm for event log data, based on process model optimization. Quantitative and qualitative results show the competitiveness and the transparency of the method. Future research will focus on the integration of more information into the model to increase prediction performance for other types of discriminative patterns.