



# Temporal graph patterns by timed automata

Amir Aghasadeghi<sup>1</sup> · Jan Van den Bussche<sup>2</sup> · Julia Stoyanovich<sup>1</sup>

Received: 4 May 2022 / Revised: 23 January 2023 / Accepted: 19 March 2023  
© The Author(s) 2023

## Abstract

Temporal graphs represent graph evolution over time, and have been receiving considerable research attention. Work on expressing temporal graph patterns or discovering temporal motifs typically assumes relatively simple temporal constraints, such as journeys or, more generally, existential constraints, possibly with finite delays. In this paper we propose to use timed automata to express temporal constraints, leading to a general and powerful notion of temporal basic graph pattern (BGP). The new difficulty is the evaluation of the temporal constraint on a large set of matchings. An important benefit of timed automata is that they support an iterative state assignment, which can be useful for early detection of matches and pruning of non-matches. We introduce algorithms to retrieve all instances of a temporal BGP match in a graph, and present results of an extensive experimental evaluation, demonstrating interesting performance trade-offs. We show that an on-demand algorithm that processes total matchings incrementally over time is preferable when dealing with cyclic patterns on sparse graphs. On acyclic patterns or dense graphs, and when connectivity of partial matchings can be guaranteed, the best performance is achieved by maintaining partial matchings over time and allowing automaton evaluation to be fully incremental. The code and datasets used in our analysis are available at <http://github.com/amirpouya/TABGP>.

**Keywords** Temporal graphs · Property graphs · Graph query languages · Timed automata · dataflow systems

## 1 Introduction

Graph pattern matching, the problem of finding instances of one graph in a larger graph, has been extensively studied since the 1970s, and numerous algorithms have been proposed [14, 16, 21, 36, 50, 59]. Initially, work in this area focused on static graphs, in which information about changes in the graph is not recorded. Finding patterns in static graphs can be helpful for many important tasks, such as finding mutual interests among users in a social network. However, understanding how interests of social network users evolve over time, support for contact tracing, and many other research questions and applications require access to information about how a

graph changes over time. Consequently, the focus of research has shifted to pattern matching in *temporal graphs* for tasks such as finding temporal motifs, temporal journeys, and temporal shortest paths [26, 34, 48, 51, 54, 68, 73, 74].

Figure 1 gives our running example, showing an interaction graph, where each node represents an *employee* (node label *emp*), a *customer* (*cst*), or an *office* (*ofc*), and each edge represents either an email *message* (*msg*) or a *visit* (*visit*) from a source node to a target node. Each edge is associated with the set of timepoints when an interaction occurred. Such graphs with static nodes but dynamic edges that are active at multiple timepoints are commonly used to represent interaction networks [48, 72].

This research was supported in part by NSF Award No. 1916505.

✉ Julia Stoyanovich  
stoyanovich@nyu.edu

Amir Aghasadeghi  
amirpouya@nyu.edu

Jan Van den Bussche  
jan.vandenbussche@uhasselt.be

<sup>1</sup> New York University, New York, NY, USA

<sup>2</sup> Hasselt University, Hasselt, Belgium

**Example 1** Assume our temporal graph holds information about a publicly traded company. Suppose that employee  $v_1$  shared confidential information with their colleagues  $v_2$  and  $v_3$ , and that one of them subsequently shared this information with customer  $v_4$ , potentially constituting insider trading. Assuming that we have no access to the content of the messages, only to when they were sent, can we identify employees who may have leaked confidential information to  $v_4$ ?

Based on graph topology alone, both  $v_2$  and  $v_3$  could have been the source of the information leak to  $v_4$ . However, by considering the timepoints on the edges, we observe that there is no path from  $v_1$  to  $v_4$  that goes through  $v_3$  and visits the nodes in temporal order. We will represent this scenario with the following basic graph pattern (BGP), augmented with a temporal constraint:

$$v_1 \xrightarrow{y_1} x \xrightarrow{y_2} v_4 : \exists t_1 \in y_1.T, \exists t_2 \in y_2.T : t_1 \leq t_2$$

Here,  $v_1$  and  $v_4$  are node constants,  $x$  is a node variable,  $y_1$  and  $y_2$  are edge variables, and  $y_1.T$  and  $y_2.T$  refer to the sets of timepoints associated with edges  $y_1$  and  $y_2$ . The temporal constraint states that there must exist a pair of timepoints  $t_1$ , associated with edge  $y_1$ , and  $t_2$ , associated with edge  $y_2$ , such that  $t_1$  occurs before  $t_2$ . We refer to such combinations of BGPs and temporal constraints as *temporal BGPs*.

The temporal constraint in the above example is *existential*: it requires the existence of timepoints where the edges from a BGP matchings are active, so that these timepoints satisfy some condition (in the example, a simple inequality).

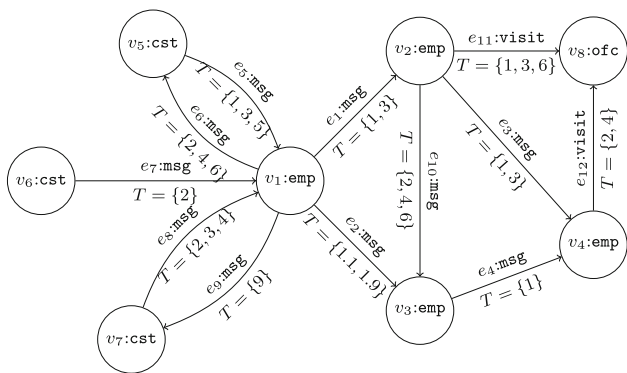


Fig. 1 Example of a temporal graph. Each edge is associated with a set of timepoints during which it is active

Existential constraints are typical in the literature on temporal graph pattern matching [26, 48, 54, 68, 73, 74]. Various forms of conditions, beyond inequalities on the timepoints, have been considered. For example, one may require that the timepoints belong to a common interval with a given start- and end-time (“temporal clique” [73]) or with a given length (“ $\delta$ -motifs” [48]), or one may specify lower and upper bounds on the gaps between the timepoints [74].

In this paper, our goal is to go beyond existential constraints. Indeed, many useful temporal constraints are *not* existential. We give two examples over temporal graphs such as the one in Fig. 1.

**Example 2** When monitoring communication patterns, we may want to look for extended interactions between customers and employees. Specifically, we are looking for matchings of the BGP shown in Fig. 2a, where edge variables  $y_1$  and  $y_2$  represent email messages exchanged by customer  $x_1$  and employee  $x_2$ . We impose the temporal constraint that  $y_1$  and  $y_2$  are active in an interleaved, alternating fashion: first  $y_1$  was active, then  $y_2$ , then  $y_1$  again, etc. This constraint is not existential. In Fig. 1, it is satisfied in the communication between  $v_5$  and  $v_1$ , but not between  $v_7$  and  $v_1$ .

**Example 3** In contact tracing, we may want to look for pairs of employees who have shared an office for a contiguous period of time with some minimal duration. We are looking for matchings of the BGP  $x_2 : emp \xrightarrow{y_1:visit} x_1 : ofc \xleftarrow{y_2:visit} x_3 : emp$ . As a temporal constraint, we impose that there exists a contiguous sequence of timepoints in the graph’s temporal domain, of duration at least, say, 3 time units, in which  $y_1$  and  $y_2$  were both active. This constraint is, again, not existential. In Fig. 1, the only matching of the BGP (involving employees  $v_2$  and  $v_4$  and office  $v_8$ ) does not satisfy the constraint; as a matter of fact,  $v_2$  and  $v_4$  were never active (i.e., at the office) at the same time!

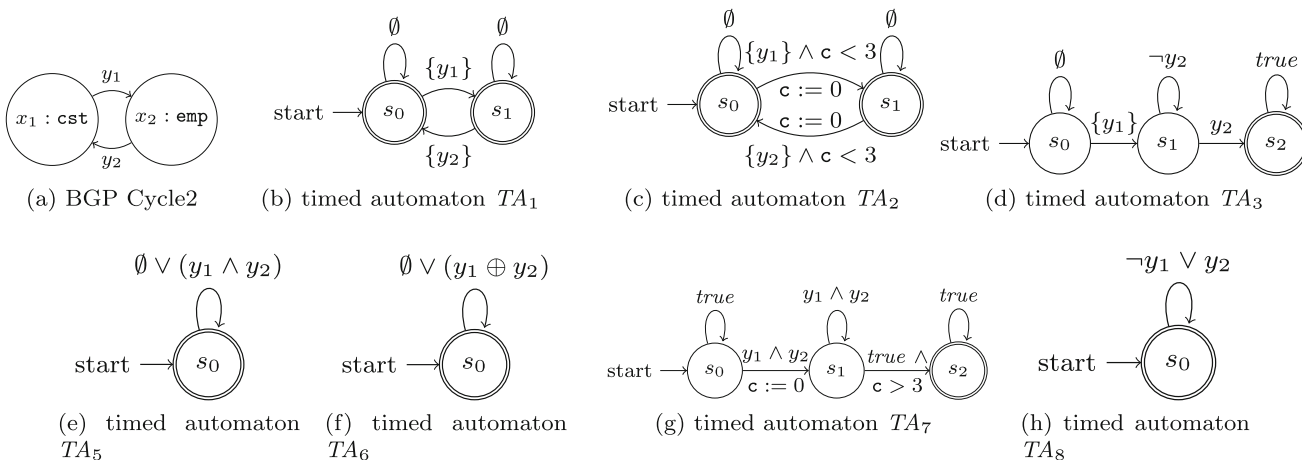


Fig. 2 Example of a temporal BGP and example timed automata explained in Example 4

Indeed, “ $y_1$  and  $y_2$  are never active at the same time” would be another natural example of a temporal constraint that is not existential.

In order to express temporal constraints (existential or not), we need a language. When the goal is the expression of possibly complex constraints, an obvious approach would be to use SQL. Indeed, any temporal graph can be naturally represented by three relations  $Node(vid, label)$ ,  $Edge(eid, vid_1, vid_2, label)$ , and  $Active(eid, time)$ , where  $vid$  and  $eid$  are node (vertex) and edge identifiers. The problem with this approach is that temporal constraints do not fit well in the SQL idiom. SQL is certainly expressive enough, but the resulting expressions tend to be complicated and hard to optimize. The alternating communication pattern from Example 2 would be expressed in SQL as follows:

```
QTA1: WITH matching AS
(SELECT E1.eid AS y1, E2.eid AS y2
 FROM edge E1, edge AS E2
 WHERE E1.dst = E2.src and E2.dst = E1.src),
Succ AS
(SELECT y1, y2, A1.eid AS e1, A2.eid AS e2
 FROM matching, active A1, active A2
 WHERE (A1.eid = y1 OR A1.eid = y2)
 AND (A2.eid = y1 OR A2.eid = y2)
 AND A1.time < A2.time AND NOT EXISTS
 (SELECT * FROM active A3
  WHERE (A3.eid = y1 or A3.eid = y2)
  AND A1.time < A3.time
  AND A3.time < A2.time))
SELECT * FROM matching M
WHERE NOT ( EXISTS (SELECT * FROM active
  WHERE eid = y2)
 AND NOT EXISTS (SELECT * FROM
  active WHERE eid = y1) )
AND ( NOT EXISTS (SELECT * FROM active A1,
  active A2
  WHERE A1.eid = y1 AND
  A2.eid = y2)
 OR (SELECT MIN(time) FROM active WHERE
  eid = y1) <
 (SELECT MIN(time) FROM active WHERE
  eid = y2) )
AND NOT EXISTS (SELECT * FROM Succ
  WHERE M.y1 = y1 AND M.y2 = y2
  AND e1 = e2)
```

Likewise, for the contiguous-duration office sharing pattern from Example 3:

```
QTA7: WITH matching AS
(SELECT E1.eid AS y1, E2.eid AS y2
 FROM edge E1, edge AS E2 WHERE E1.dst =
  E2.dst)
SELECT DISTINCT y1, y2
FROM matching, active A1, active A2, active
  B1, active B2
WHERE A1.eid = y1 AND A2.eid = y2 AND A1.time
  = A2.time
AND B1.eid = y1 AND B2.eid = y2 AND B1.time
  = B2.time
AND B1.time - A1.time > 3
```

```
AND NOT EXISTS
(SELECT * FROM active C
 WHERE A1.time < C.time AND C.time <
  B1.time
 AND NOT EXISTS (SELECT * FROM active
  C1, active C2
  WHERE C1.time = C.time
  AND C2.time = C.time
  AND C1.eid = y1 AND
  C2.eid = y2))
```

The hypothesis put forward in this paper is that specification formalisms used in fields such as complex event recognition [24] or verification of real-time systems [9] may be much more suitable for the expression of complex temporal constraints. In this paper, we specifically investigate the use of *timed automata* [2, 9].

**Example 4** Figure 2 shows various examples of timed automata that can be applied to matchings of a BGP with two edge variables  $y_1$  and  $y_2$ , such as the BGPs considered in Examples 2 and 3. One can think of the automaton as running over the snapshots of the temporal graph. A matching is accepted if there is a run such that, after seeing the last snapshot, the automaton is in an accepting state. The edge variables serve as Boolean conditions on the transitions of the automaton. When the edge matched to  $y_1$  ( $y_2$ ) is active in the current snapshot, the Boolean variable  $y_1$  ( $y_2$ ), is true. We use  $\emptyset$  as an abbreviation for  $\neg y_1 \wedge \neg y_2$ , and  $\{y_1\}$  for  $y_1 \wedge \neg y_2$  (and similarly  $\{y_2\}$ ).

The alternation constraint of Example 2 is expressed by  $TA_1$ .  $TA_2$  is similar but additionally requires that each message gets a reply within 3 time units (a clock  $c$  is used for this purpose). The contiguous-duration constraint of Example 3 is expressed by  $TA_7$ , also using a clock. The constraint “ $y_1$  and  $y_2$  are never active together” is expressed by  $TA_6$ ; the opposite constraint “ $y_1$  and  $y_2$  are always active together” by  $TA_5$ . Likewise,  $TA_8$  expresses that  $y_2$  is active whenever  $y_1$  is (in SQL, this constraint would correspond to a set containment join [8]). Finally,  $TA_3$  expresses that  $y_1$  has been active strictly before the first time  $y_2$  becomes active. Also, existential constraints such as the one from Example 1 are readily expressible by timed automata (see Sect. 3).

Timed automata offer not only a good balance between expressiveness and simplicity. A temporal constraint expressed by a timed automaton can also be processed efficiently, as the iterative state assignment mechanism allows early acceptance and early rejection of matchings. In this paper, we will introduce three algorithms for the evaluation of temporal BGPs with timed automata as temporal constraints. The first is a baseline algorithm intended for offline processing when the complete history of graph evolution is available at the time of execution. The second is an on-demand algorithm that supports online query processing when the temporal graph arrives as a stream. The third is a

partial-match algorithm that speeds up processing by sharing computation between multiple matches.

We will present an implementation of these algorithms in a dataflow framework, and will analyze performance trade-offs induced by the properties of the temporal BGP and of the underlying temporal graph. We will also compare performance with main-memory SQL systems, and will observe that temporal BGPs with temporal constraints that are not existential can be impractical when expressed and processed as SQL queries.

## 2 Temporal graphs and temporal graph patterns

We begin by recalling the standard notions of graph and graph pattern used in graph databases [4, 63]. Assume some vocabulary  $L$  of labels. We define:

**Definition 1** (*Graph*) A graph is a tuple  $(N, E, \rho, \lambda)$ , where:

- $N$  and  $E$  are disjoint sets of *nodes* and *edges*, respectively;
- $\rho : E \rightarrow (N \times N)$  indicates, for each edge, its source and destination nodes; and
- $\lambda : N \cup E \rightarrow L$  assigns a label to every node and edge.

Figure 1 gives an example of a graph, with  $N = \{v_1, \dots, v_7\}$ ,  $E = \{e_1, \dots, e_9\}$ ,  $\lambda(v_1) = \lambda(v_2) = \lambda(v_3) = \text{emp}$ ,  $\lambda(v_4) = \dots = \lambda(v_7) = \text{cst}$ , and  $\lambda(e_1) = \dots = \lambda(e_9) = \text{msg}$ . In this graph,  $\rho(e_5) = (v_5, v_1)$ .

Next, recall the conventional notion of basic graph pattern (BGP).

**Definition 2** (*Basic graph pattern (BGP)*) A BGP is a tuple  $(C, X, Y, \rho, \lambda)$ , where:

- $C, X$  and  $Y$  are pairwise disjoint finite sets of *node constants*, *node variables*, and *edge variables*, respectively;
- $\rho : Y \rightarrow (C \cup X) \times (C \cup X)$  indicates, for each edge variable, its source and destination, which can be a node constant or a node variable; and
- $\lambda : X \cup Y \rightarrow L$  is a partial function, assigning a label from  $L$  to some of the variables.

The fundamental task related to BGPs is to find all matchings in a graph, defined as follows:

**Definition 3** (*Matching*) A *partial matching* of a BGP  $(C, X, Y, \rho_P, \lambda_P)$  in a graph  $G = (N, E, \rho, \lambda)$  is a function  $\mu : Z \rightarrow N \cup E$  satisfying the following conditions:

- $Z$ , the domain of  $\mu$ , is a subset of  $X \cup Y$ .
- $\mu(Z \cap X) \subseteq N$  and  $\mu(Z \cap Y) \subseteq E$ .

- Let  $y$  be an edge variable in  $Z$  and let  $\rho_P(y) = (x_1, x_2)$ . Then, for  $i = 1, 2$ , if  $x_i$  is a node variable, then  $x_i \in Z$ . Moreover,  $\rho(\mu(y)) = (\mu(x_1), \mu(x_2))$ , where we agree that  $\mu(c) = c$  for any node constant  $c$ .
- For every  $z \in Z$  for which  $\lambda_P(z)$  is defined, we have  $\lambda(\mu(z)) = \lambda_P(z)$ .

If  $Z$  equals  $X \cup Y$  then  $\mu$  is called a (total) *matching*.

Consider the BGP in Fig. 2a. Evaluating it over the graph in Fig. 1 yields 7 partial matchings:  $v_1 \xrightarrow{e_6} v_5$ ,  $v_1 \xrightarrow{e_9} v_7$ ,  $v_2 \xrightarrow{e_3} v_4$ ,  $v_3 \xrightarrow{e_4} v_4$ ,  $v_5 \xrightarrow{e_5} v_1$ ,  $v_6 \xrightarrow{e_7} v_1$ ,  $v_7 \xrightarrow{e_8} v_1$ , and 2 total matchings:  $v_5 \xrightarrow{e_5} v_1 \xrightarrow{e_6} v_5$  and  $v_7 \xrightarrow{e_8} v_1 \xrightarrow{e_9} v_7$  as total matchings.

**Remark 1** Our semantics of matching is based on the semantics of standard conjunctive queries, so does not require matchings to be injective. An alternative approach in graph matching is to search for exact copies of subgraphs (graph isomorphism). Many of our general ideas also apply to such a semantics. For example, our baseline algorithm (Sect. 4.2) can be used under injective semantics. However, once more optimized query processing strategies are involved, injectivity can make a difference. For example, classical minimization (i.e., finding redundant joins) becomes much more involved [58]. On the other hand, particular algorithms exploiting graph isomorphism can be devised under injective semantics.

We now present the notion of a *temporal graph* in which edges are associated with sets of timepoints, while nodes persist over time. Extending our work to temporal property graphs in which both nodes and edges are associated with temporal information, and where the properties of nodes and edges can change over time [41], is an interesting direction for further research. We assume that *timepoints* are strictly positive real numbers and define:

**Definition 4** (*Temporal graph*) A temporal graph is a pair  $(G, \xi)$ , where  $G$  is a graph and  $\xi$  assigns a finite set of timepoints to each edge of  $G$ . When  $e$  is an edge and  $t \in \xi(e)$ , we say that  $e$  is *active* at time  $t$ .

In the temporal graph in Fig. 1,  $\rho(e_5) = (v_5, v_1)$  and  $\xi(e_5) = \{1, 3, 5\}$ , indicating that  $v_5$  messaged  $v_1$  at the listed timepoints.

To extend the notion of matching to temporal graphs, we enrich BGPs with temporal constraints, defined as follows.

**Definition 5** (*Temporal variables, assignments, and constraints*) Let  $V$  be a set of temporal variables. A *temporal assignment*  $\alpha$  on  $V$  is a function that assigns a finite set of timepoints to every variable in  $V$ . A *temporal constraint* over  $V$  is a set of temporal assignments on  $V$ . This set is typically infinite. When a temporal assignment  $\alpha$  belongs to a temporal constraint  $\Gamma$ , we also say that  $\alpha$  *satisfies*  $\Gamma$ .

For the moment, this is a purely semantic definition of temporal constraints; in Sect. 3 we will present how such constraints may be specified using timed automata.

If we have a matching  $\mu$  from a BGP in a graph  $G$ , and we consider a temporal graph  $(G, \xi)$  based on  $G$ , we automatically obtain a temporal assignment on the edge variables of the BGP. Indeed, each edge variable is matched to an edge in  $G$ , and we take the set of timepoints of that edge. Thus, edge variables serve as temporal variables, and we arrive at the following definition:

**Definition 6** (*Temporal BGP, matching*) A temporal BGP is a pair  $(P, \Gamma)$  where  $P$  is a BGP and  $\Gamma$  is a temporal constraint over  $Y$  (the edge variables of  $P$ ).

Let  $(G, \xi)$  be a temporal graph. Given a matching  $\mu$  of  $P$  in  $G$ , we can consider the temporal assignment  $\alpha_\mu$  on  $Y$  defined by

$$\alpha_\mu(y) := \xi(\mu(y)) \quad \text{for } y \in Y.$$

Now a *matching* of the temporal BGP  $(P, \Gamma)$  in the temporal graph  $(G, \xi)$  is any matching  $\mu$  of  $P$  in  $G$  such that  $\alpha_\mu$  satisfies  $\Gamma$ .

In the next section, we describe how timed automata such as that in Fig. 2b can be used to represent and enforce such constraints.

### 3 Expressing temporal constraints

Our conception of a temporal BGP, as a standard BGP  $P$  equipped with a temporal constraint  $\Gamma$  on the edge variables of  $P$ , leaves open how  $\Gamma$  is specified. We pursue the idea to use *timed automata*, an established formalism for expressing temporal constraints in the area of verification [2, 9]. Timed automata are often interpreted over infinite words, but here we will use them on finite words.

*Timed automata* A timed automaton over a finite set  $Y$  of variables is an extension of the standard notion of non-deterministic finite automata (NFA), over the alphabet  $\Sigma = 2^Y$  (the set of subsets of  $Y$ ). Recall that an NFA specifies a finite set of states: an initial state, a set of final states, and a set of transitions of the form  $(s_1, \theta, s_2)$ , where  $s_1$  and  $s_2$  are states and  $\theta$  is a Boolean formula over  $Y$ . The automaton reads a word over  $\Sigma$ , starting in the initial state. Whenever the automaton is in a current state  $s_1$ , the next letter to be read is  $a$ , and there exists a transition  $(s_1, \theta, s_2)$  such that  $a$  satisfies  $\theta$ , the automaton can change state to  $s_2$  and move to the next letter. If, after reading the last letter, the automaton is in a final state, the run accepts. If there is no suitable transition at some point, or if the last state is not final, then the run fails. A word is accepted if there exists an accepting run.

The extra feature added by timed automata to the standard NFA apparatus is a finite set  $C$  of *clocks*, which can be used to measure time gaps between successive letters in a *timed* word (to be defined momentarily). Transitions are of the extended form

$$(s_1, \theta, \delta, R, s_2), \tag{*}$$

where  $s_1, \theta$  and  $s_2$  are as in NFAs;  $\delta$  is a Boolean combination of *clock conditions*; and  $R$  is a subset of  $C$ . Here, by a clock condition, we mean a condition of the form  $c \leq g$  or  $c \geq g$ , where  $c$  is a clock and  $g$  is a real number constant representing a time gap.

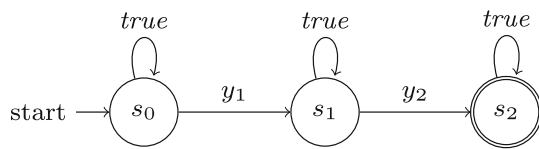
As just mentioned, a timed automaton works over timed words. A timed word over an alphabet  $\Sigma$  is a sequence of the form

$$(t_1, a_1) \dots (t_n, a_n),$$

where each  $a_i \in \Sigma$ , and  $t_1 < \dots < t_n$  are timepoints. When the automaton is started on the timed word, all clocks are initially set to  $t_0 := 0$ . For  $i = 1, \dots, n$  the automaton runs as follows. Upon reading position  $(t_i, a_i)$ , every clock has increased by  $t_i - t_{i-1}$ . Now the automaton can take a transition  $(*)$  as above on condition that the current state is  $s_1$  and  $a_i$  satisfies  $\theta$ , as before; and, moreover, the current valuation of the clocks satisfies  $\delta$ . If this is so, the automaton can change state to  $s_2$ , move to the next position in the timed word, and must reset to zero all clocks in  $R$ . As with NFAs, a run is accepting if it ends in a final state, and a timed word is accepted if there exists an accepting run.

*Using timed automata to express temporal BGP constraints* A timed automaton defines the set of timed words that it accepts. But how does it define, as announced, a temporal constraint over  $Y$ , which is not a set of timed words, but a set of temporal assignments? This is simple once we realize that a temporal assignment over  $Y$ , in the context of a temporal graph  $H = (G, \xi)$ , is nothing but a timed word over  $\Sigma = 2^Y$ . We can see this as follows. Let  $T = \bigcup\{\xi(e) \mid e \in E\}$  be the set of all distinct timepoints used in  $H$ ; we also refer to  $T$  as the *temporal domain* of  $H$ . Let  $T = \{t_1, \dots, t_n\}$ , ordered as  $t_1 < \dots < t_n$ . We can then view any temporal assignment  $\alpha : Y \rightarrow T$  as the timed word  $(t_1, a_1) \dots (t_n, a_n)$ , where  $a_i = \{y \in Y \mid t_i \in \alpha(y)\}$ .

NFAs are a special case of timed automata without any clocks, and this special case is already useful for expressing temporal constraints. For example, consider the NFA in Fig. 3 which expresses the existential constraint  $\exists t_1 \in y_1 \exists t_2 \in y_2 : t_1 < t_2$  from Example 1. Suppose we instead want to express the existential constraint  $\exists t_1 \in y_1 \exists t_2 \in y_2 : t_1 < t_2 - 7$ . We can do this by introducing a clock  $c$ . When we see  $y_1$ , we reset the clock ( $c := 0$ ). Then, when we see  $y_2$ , we check



**Fig. 3**  $TA_e$ : an existential constraint

that the clock has progressed beyond the desired 7 time units ( $c > 7$ ).

Of course, as already argued in the Introduction, timed automata are much more powerful than mere existential constraints.

**Checking acceptance** The algorithm for checking whether an input timed word  $(t_1, a_1) \dots (t_n, a_n)$  is accepted by a timed automaton extends the analogous algorithm for NFAs [28]. We maintain a set  $S$  of current possible configurations. Each configuration is a pair  $(s, \nu)$ , where  $s$  is a state of the automaton, and  $\nu$  is a real-valued map defined on  $C$  (the set of clocks). We initialize  $S$  as the singleton  $S_0 := \{(s_0, \nu)\}$ , where  $s_0$  is the initial state and  $\nu(c) = 0$  for each  $c \in C$ .

For each  $i = 1, \dots, n$ , we now update  $S_{i-1}$  to  $S_i$  as follows. Recall that  $t_0 = 0$ .

1. Initialize  $S_i := \emptyset$ .
2. For each  $(s, \nu) \in S_{i-1}$ , find all applicable transitions as follows.
  - (a) Let  $\nu'$  be the map obtained by  $\nu$  by updating  $\nu'(c) := \nu(c) + (t_i - t_{i-1})$  for each clock  $c$ .
  - (b) Consider each transition  $(s_1, \theta, \delta, R, s_2)$  where  $s_1 = s$ , and test whether  $a_i$  satisfies  $\theta$  and  $\nu'$  satisfies  $\delta$ . If the test succeeds, let  $\nu''$  be the same as  $\nu'$  except that we reset  $\nu''(c) := 0$  for each  $c \in R$ . Now add  $(s_2, \nu'')$  to  $S_i$ .

If, in the end,  $S_n$  contains a final state, the timed word is accepted. This procedure runs, for any fixed automaton, in time linear in  $n$ .

**Expressive power of timed automata** Since their inception in 1994, timed automata were rapidly adopted in the real-time verification community, but were ignored in research on temporal query languages. Timed automata can be used to express a rich set of temporal constraints that is now considered standard [9] and remains a topic of active research.<sup>1</sup>

Part of the reason for the popularity of timed automata is that they provide a natural and concise syntax for the expression of temporal constraints that can be directly used by programmers of complex queries on temporal graphs. Equivalent logical formalisms have been developed, based

on timed temporal logics [23], and they can be used as basis for an alternative syntax for temporal constraints.

Alternatively, as discussed in the Introduction, SQL can be used to express temporal constraints. A set of matchings over the set of edge variables  $Y$  is naturally represented by a relation  $M$  over the relation schema  $Y$ , using temporal variables as attribute names. Together with a table  $Active(y, t)$  indicating when each edge is active, we obtain a relational representation of the temporal words that represent the matchings in  $M$ . It is an open question whether every timed automaton  $\Gamma$  over the alphabet  $\Sigma = 2^Y$  is equivalently expressible in SQL, meaning that the SQL query returns exactly all matchings from  $M$  that satisfy  $\Gamma$ . The converse question (whether timed automata can express all SQL-expressible constraints) is open as well. Expressiveness questions about SQL that include arithmetic operations (notably, arithmetic over timepoints, in our case), are notoriously hard to settle [39].

As a final alternative, also mentioned in the Introduction, rather than separating the BGP and the temporal constraints, one can attempt to use SQL holistically, representing both kinds of constraints in a SQL query. A natural way to proceed is to represent the temporal graph by three relations, *Node*, *Edge*, and *Active*, without separating the BGP and the temporal constraint. However, while this may be practical for a restricted class of temporal constraints, the SQL expressions tend to become very complex, and thus hard both to implement and to optimize. Furthermore, by treating temporal constraints as part of the query, we largely lose the opportunity to exploit the actual temporal nature of the data in query processing. We will support these claims experimentally in Sect. 6.4.

## 4 Algorithms for timed-automaton temporal graph pattern matching

In this section, we will discuss several algorithms for processing temporal BGPs with temporal constraints given as timed automata.

### 4.1 Temporal graph representation

Assume that we are given a temporal graph  $H = (G, \xi)$  where  $G = (N, E, \rho_G, \lambda_G)$ . We are also given a temporal BGP  $Q = (P, \Gamma)$  where  $P = (X, Y, \rho_P, \lambda_P)$  as defined in Sect. 2. Our goal is to compute all matchings  $\mu$  of  $Q$  in  $H$ ; recall that this means that  $\mu$  must be a matching of  $P$  in  $G$ , and, moreover, the temporal assignment  $\alpha_\mu$  must satisfy  $\Gamma$  (Definition 5). We assume that  $\Gamma$  is specified as a timed automaton over the alphabet  $\Sigma = 2^Y$ .

<sup>1</sup> In 2022 alone, eighty new publications appeared on the applications of timed automata.

We can naturally represent  $H$  by the following relations:  $Node(vid, label)$ ,  $Edge(eid, vid_1, vid_2, label)$ , and  $Active(eid, time)$ .

We recall the natural notion of snapshot from temporal databases:

**Definition 7 (Snapshot)** For any timepoint  $t_i$  in the temporal domain of  $H$ , the *snapshot* of  $H$  at time  $t_i$  is the subgraph of  $G$  induced by all edges that are active at time  $t_i$ . We denote this subgraph by  $H_{t_i}$ .

The snapshot can be represented by the relevant slices of the tables  $Node$ ,  $Edge$  and  $Active$ .

We will present three algorithms for finding the matchings of a temporal BGP in a temporal graph, when the temporal constraint is given by a timed automaton. We start by presenting our baseline algorithm that operates in two stages. First, it generates all matchings of the BGP; next, it filters out those matchings that violate the temporal constraint. Our second on-demand algorithm works incrementally. It considers the graph in temporal order, snapshot by snapshot. As time progresses, more edges of the graph are seen so that more and more BGP matchings are found. Also, at each timestep, the possible transitions of the automaton are evaluated to keep track of the possible states for each matching. For newly found matchings, however, the automaton has to catch up from the beginning. This catching up is avoided in our third partial-match algorithm, which incrementally maintains all *partial* matches of the BGP, refining them as time progresses.

### 4.2 Baseline algorithm

Assume a temporal BGP  $Q = (P, \Gamma)$ , where  $\Gamma$  is specified as a timed automaton. Given a temporal graph  $H = (G, \xi)$ , we want to find all the matchings of  $Q$  in  $H$ . We do this in two stages:

Find the BGP matchings: Find all matchings of  $P$  in  $G$  using any of the available algorithms for this task (e.g., a worst-case optimal join algorithm [3, 46]).

Run the automaton: For each obtained matching  $\mu$ :

- (a) Convert the assignment  $\alpha_\mu$  into a timed word over the temporal domain of  $H$ , as described in Sect. 3. We denote this timed word by  $w_\mu$ .
- (b) Check if  $w_\mu$  is accepted by automaton  $\Gamma$ .

We next describe how the automaton stages (a) and (b) can be done synchronously, for all matchings  $\mu$  together. Effectively, we can obtain a bulk-processing variant of the acceptance algorithm described at the end of Sect. 3, as follows.

We use a table  $States$  that holds triples  $(\mu, s, v)$ , where  $\mu$  is a matching;  $s$  is a state of the automaton; and  $v$  is an assignment of timepoints to the clocks of  $\Gamma$ . Since  $\Gamma$  is non-deterministic, the same  $\mu$  may be paired with different  $s$  and  $v$ . Naturally, in the initial content of  $States$ , each  $\mu$  is paired with state  $s_0$  and  $v_0$  that maps every clock to 0.

Let  $T = \{t_1, \dots, t_n\}$  with  $t_1 < \dots < t_n$  be the temporal domain of  $H$  as described in Sect. 3, and let  $t_0 := 0$ . Recall that the active timepoints for each edge are stored in the table  $Active(eid, time)$ . We obtain  $T$  by first sorting  $Active$  on time and then scanning through it. Now during this scan, for  $i = 1, \dots, n$ , we do the following:

1. Update each  $(\mu, s, v)$  in  $States$  by increasing every clock value by  $t_i - t_{i-1}$ .
2. Let  $Y$ , the set of edge variables of  $P$ , be  $\{y_1, \dots, y_k\}$ . Extend each  $(\mu, s, v)$  in  $States$  with Boolean values  $b_1, \dots, b_k$  defined as follows:  $b_j$  is true if edge  $\mu(y_j)$  is active at the current time  $t_i$ , and false otherwise. Observe that the bit vector  $b_1 \dots b_k$  represents the  $i$ -th letter of the timed word  $w_\mu$ .
3. Join all records  $(\mu, s, v, b_1 \dots b_k)$  from  $States$  with all transitions  $(s_1, \theta, \delta, R, s_2)$  from  $\Gamma$ , where the following conditions are satisfied:  $s = s_1$ ;  $b_1 \dots b_k$  satisfies  $\theta$ ; and  $v$  satisfies  $\delta$ .
4. Project every joined tuple on  $(\mu, s_2, v'')$ , where  $v''$  is  $v$  but with every clock from  $R$  reset to 0. The resulting projection is the new content of  $States$ .

*Complexity* Each of the above steps can be accomplished by relational-algebra-like dataflow operations over the  $States$  table. In particular, step 4.2 is done by successive left outer joins. For  $j = 1, \dots, k$ , let  $A_j$  be the  $Active$  table, filtered on  $time = t_i$ , and renaming  $eid$  to  $b_j$ . We left-outer join  $States$  with  $A_j$  on condition  $y_j = b_j$ . If, in the result,  $b_j$  is null, it is replaced by false; otherwise it is replaced by true. The entire second stage, for a fixed timed automaton, can be implemented in time  $O(A + nM)$ , where  $A$  is the size of the  $Active$  table,  $n$  is the size of the temporal domain, and  $M$  is the number of matchings returned from the first stage.

The first stage of the baseline algorithm, finding BGP matchings in a graph, is a well-researched and still active topic [7]. The naive complexity upper-bound is  $O(N^k)$ , where  $N$  is the size of  $G$  and  $k$  is the size of the BGP. However, the more advanced algorithms cited above are worst-case optimal [3, 46], meaning that their running time is proportional to  $N + g(N)$ , where  $g$  is the so-called AGM-bound of the BGP. Using such an algorithm, and noting that  $N \leq A$ , the total complexity for stages one and two combined is  $O(g(N) + A + nM)$ .

*Early acceptance or rejection* After the iteration for  $i = n$ , the matchings that are accepted by the automaton  $\Gamma$  are those that are paired in  $States$  with an accepting state. We may also

**Fig. 4** Content of the *States* relation at  $t = 0, \dots, 3$ , illustrating the execution of the timed automaton in Fig. 2c by the baseline algorithm of Sect. 4.2

t	$\mu$	s	v	$b_1 b_2$
0	$e_5, e_6$	$s_0$	$\square$	00
	$e_8, e_9$	$s_0$	$\square$	00
1	$e_5, e_6$	$s_1$	$[x=0]$	10
	$e_8, e_9$	$s_0$	$\square$	00
1.1	$e_5, e_6$	$s_1$	$[x=.1]$	00
	$e_8, e_9$	$s_0$	$\square$	00
1.9	$e_5, e_6$	$s_1$	$[x=.9]$	00
	$e_8, e_9$	$s_0$	$\square$	00
2	$e_5, e_6$	$s_0$	$[x=1]$	01
	$e_8, e_9$	$s_1$	$[x=0]$	10
3	$e_5, e_6$	$s_1$	$[x=0]$	10
	$e_8, e_9$			10

be able to accept results early: when, during the iteration, a matching  $\mu$  is paired with an accepting state  $s$ , and all states reachable from  $s$  in the automaton are also accepting, then  $\mu$  can already be output. On the other hand, when all states reachable from  $s$  are *not* accepting, we can reject  $\mu$  early.

**Example 5** Consider again the temporal graph in Fig. 1, and suppose that we want to find all cycles of length 2 shown in Fig. 2a, under the temporal constraint  $TA_2$  shown in Fig. 2c. The first stage of the baseline algorithm identified two matchings,  $\mu = (e_5, e_6)$  and  $\mu = (e_8, e_9)$ . These are considered by the timed automaton in the second stage.

Figure 4 shows the *States* relation with the timed words  $w_\mu$  at times between 0 to 3. At  $t = 0$ , both matchings are at  $s = s_0$ , no clocks have been set, and, since neither of the matchings has any edges,  $b_1 = 0$  and  $b_2 = 0$ . At time  $t = 1$ ,  $e_5$  is active, hence the bit  $b_1$  for the matching  $(e_5, e_6)$  is set to 1, and, since there is a rule in the timed automaton, state is updated to  $s_1$  and clock  $x$  is set to 1. Matching  $(e_8, e_9)$  does not exist at time 1, and so no change is made in that row of the *State* table. At times  $t = 1.1$  and  $t = 1.9$ , neither of the matchings' edges are active, hence the only change is that the clock is updated for  $(e_5, e_6)$ . At time  $t = 2$ , for matching  $(e_5, e_6)$ , the edge  $e_6$  is active and the clock  $x$  is less than 2, hence we move back to state  $s_0$ . For the matching  $(e_8, e_9)$ ,  $e_8$  is active so we move to  $s_1$  and set the clock. At time  $t = 3$ ,  $(e_5, e_6)$  continues to alternate, but for  $(e_8, e_9)$  we see that  $e_8$  is active, hence, we set  $b_1 = 1$  and  $b_2 = 0$ , and, seeing that the timed automaton does not have a transition, we drop this matching (shown as grayed out in Fig. 4). Between times 3–6, the matching  $(e_5, e_6)$  continues alternating between  $s_0$  and  $s_1$ . From time 7–9, we observe neither of  $e_5$  or  $e_6$ , and hence no change happens to the state of this matching. The final output of this algorithm is that matching  $(e_5, e_6)$  is accepted at state  $s_0$ .

### 4.3 On-demand algorithm

A clear disadvantage of the baseline algorithm is that we must first complete the first stage (BGP matching on the whole underlying graph  $G$ ) before we can move to the automaton stage. This delay may be undesirable and prohibits return-

ing results early in situations where the temporal graph is streamed over time. We next describe our second algorithm, which works incrementally by processing *snapshots* in chronological order.

Recall Definition 7 of snapshots. We also define:

**Definition 8 (History)** The *history* of  $H$  until time  $t_i$ , denoted  $H_{\leq t_i}$ , is the union of all snapshots  $H_{t_j}$  for  $j = 1, \dots, i$ . For  $t_0 := 0$ , we define  $H_{\leq t_0}$  to be the empty graph.

The baseline algorithm is now modified as follows. We no longer have a first stage. Snapshots arrive chronologically at timepoints  $t_1, \dots, t_n$ ; it is not necessary for the algorithm to know the entire temporal domain  $\{t_1, \dots, t_n\}$  in advance. For  $i = 1, \dots, n$ :

1. We receive as input the next snapshot  $H_{t_i}$ . In previous iterations we have already computed all matchings of  $P$  in the preceding history  $H_{\leq t_{i-1}}$ . Using this information and the next snapshot, we compute the new matchings, i.e., the matchings of  $P$  in the current history  $H_{\leq t_i}$  that were not yet matchings of  $P$  in the preceding history. Incremental BGP matching is a well-researched topic, and any of the available algorithms can be used here [22, 25, 31, 69].
2. We use the table *States* as in the baseline algorithm. For each newly discovered matching  $\mu$ , we must catch up and run the automaton from the initial state on the prefix of  $w_\mu$  of length  $i - 1$ . We add to *States* all triples  $(\mu, s, v)$ , such that the configuration  $(s, v)$  is a possible configuration of the automaton after reading the prefix.
3. All matchings we already had remain valid; indeed, if  $\mu$  is a matching of  $P$  in  $H_{\leq t_{i-1}}$  then  $\mu$  is also a matching of  $P$  in  $H_{\leq t_i}$ . *States* is now updated for the  $i$ -th letter of the timed words of all matchings, new and old, as in the baseline algorithm.

We call this the “on-demand” algorithm because the automaton is run from the beginning, on demand, each time new matchings are found, in order to catch up with table *States* holding the possible automaton configurations.

**Example 6** Figure 5 shows the *State* relation for the on-demand algorithm, for the same BGP and temporal constraint



	$\mu$	$s$	$v$	$b_1 b_2$
$t = 2$	$(e_5, e_6)$	$s_0$	$[x=1]$	01
$t = 3$	$(e_5, e_6)$	$s_1$	$[x=0]$	10
$\dots$				
$t = 9$	$(e_5, e_6)$	$s_0$	$[x=0]$	00
	$(e_8, e_9)$			01

**Fig. 5** Content of the *States* relation at  $t = 0, \dots, 3$ , illustrating the execution of the timed automaton in Fig. 2c by the on-demand algorithm of Sect. 4.3

as in Example 5. The first time a cycle of length 2 exists in the graph in Fig. 1 is  $t = 2$ , hence there will be no matching in any iteration before  $t = 2$  and no temporal automaton would run. At time  $t = 2$  the incremental matching algorithm finds the matching  $(e_5, e_6)$  and passes it to the timed automaton that runs it for  $t = 0, t = 1, t = 1.1$  and  $t = 1.9$  as was described in Example 5. At time  $t = 9$ , edge  $e_9$  is received and gives rise to a new matching  $(e_8, e_9)$ . At that point, the timed automaton is invoked for all  $t < 9$ . The process is similar to what we described in Example 5, and the on-demand automaton will eliminate this matching at  $t = 3$  because no rule in the automaton can be satisfied. The final output is the same as for the baseline algorithm: matching  $(e_5, e_6)$  accepted at state  $s_0$ . Note that using on-demand algorithm, we can process the graph that arrives as a stream.

**Correctness and Complexity** The correctness of the on-demand algorithm follows because it is an incremental variant of the baseline algorithm. Every record  $(\mu, s, v)$  that occurs in table *States* in the baseline algorithm will also occur in the on-demand algorithm. Indeed,  $\mu$ , which is a matching of  $P$  in  $G$ , will be found in some iteration, at the very least in the last iteration, since  $H_{\leq t_n}$  equals  $G$ . Conversely, if  $(\mu, s, v)$  occurs in *States* in the on-demand algorithm, it will also occur in the baseline algorithm. Indeed, if  $\mu$  is a matching of  $P$  in some history  $H_{\leq t_i}$ , then  $\mu$  is also a matching of  $P$  in  $G$ , since  $H_{\leq t_i} \subseteq G$ .

The complexity of the on-demand algorithm is, theoretically, not better than that of the baseline algorithm. Let  $M_i$  be the number of new matchings found in iteration  $i$ , so that  $M_1 + \dots + M_n = M$ . We assume new matchings are found with a worst-case optimal join algorithm applied to the delta query [46], and, by slight abuse of notation, we, again, denote the AGM bound for that query by  $g$ . Still, we need to catch up the automaton configurations for the  $M_i$  new matchings, and also do the state updates for all  $M_1 + \dots + M_i$  matchings found so far. This leads to a complexity in the order of

$$\sum_{i=1}^n \left( N + g(N) + \sum_{j=1}^i M_j \right)$$

which appears worse than  $A + g(N) + nM$  for the baseline algorithm.

However, note that the  $O(N + g(N))$  estimate for each incremental query is a worst-case upper bound only. Moreover, the advantage remains that the on-demand algorithm can be applied in a streaming context. So, realistic running time comparisons must be made experimentally. Indeed, in our experiments, we will see that on-demand is typically better than the baseline.

#### 4.4 Partial-match algorithm

A disadvantage of the on-demand algorithm is the catching-up of the automaton on newly found matchings. Interestingly, we can avoid any catching-up and obtain a fully incremental algorithm, provided we keep not only the total matchings of  $P$  in the current history, but also all *partial* matchings.

Specifically, we will work with *maximal* partial matchings: these are partial matchings that cannot be extended to a strictly larger partial matching on the same graph. Now, for any partial matching  $\mu$  of  $P$  in  $G$ , we can define a timed word  $w_\mu$ , in the same way as for total matchings. Formally,  $w_\mu = (t_1, a_1) \dots (t_n, a_n)$ , where now  $a_i = \{y \in Y \mid \mu \text{ is defined on } y \text{ and } t_i \in \xi(\mu(y))\}$ . The following property now formalizes how a fully incremental approach is possible:

**Proposition 1** *Let  $\mu$  be a maximal partial matching of  $P$  in  $H_{\leq t_{i-1}}$ , and let  $\mu'$  be a partial matching of  $P$  in  $H_{\leq t_i}$ , such that  $\mu \subseteq \mu'$ . Then the timed words  $w_\mu$  and  $w_{\mu'}$  have the same prefix of length  $i - 1$ .*

**Proof** Let  $w_\mu = (t_1, a_1), (t_2, a_2) \dots (t_n, a_n)$  and  $w_{\mu'} = (t_1, b_1), (t_2, b_2) \dots (t_n, b_n)$ . We must show that  $a_j = b_j$  for  $j = 1, \dots, i - 1$ . The containment from left to right is straightforwardly verified. Indeed, take  $y \in a_j$ . Then  $\mu(y)$  is defined and  $t_j \in \xi(\mu(y))$ . Since  $\mu \subseteq \mu'$ , also  $\mu'(y) = \mu(y)$  is defined and we see that  $y \in b_j$  as desired.

For the containment from right to left, take  $y \in b_j$ . Then  $\mu'(y)$  is defined and  $t_j \in \xi(\mu'(y))$ . For the sake of contradiction, suppose  $\mu(y)$  would not be defined. Let  $\rho_P(y) = (x_1, x_2)$ , and strictly extend  $\mu$  to  $\mu''$  by mapping  $y$  to  $\mu'(y)$ ;  $x_1$  to  $\mu'(x_1)$ ; and  $x_2$  to  $\mu'(x_2)$ . Since  $\mu'$  is a partial matching of  $P$  in  $G$ , we know that  $\mu'(y)$  is an edge in  $G$  from node  $\mu'(x_1)$  to node  $\mu'(x_2)$ . Moreover, since  $t_j \in \xi(\mu'(y))$ , the edge  $\mu'(y)$  is present in  $H_{\leq t_j}$ , so certainly also in  $H_{\leq t_{i-1}}$  since  $j \leq i - 1$ . Thus,  $\mu''$  is a partial matching of  $P$  in  $H_{\leq t_{i-1}}$ , contradicting the maximality of  $\mu$ . We conclude that  $\mu(y)$  is defined, and  $y \in a_j$ .

Concretely, the *partial-match* algorithm incrementally maintains, for  $i = 1, \dots, n$ , all maximal partial matchings  $\mu$  of  $H_{\leq t_i}$ , along with the possible configurations  $(s, v)$  of the automaton after reading the  $i$ -th prefix of the timed word  $w_\mu$ .

The triples  $(\mu, s, v)$  are kept in the table *States* as before. Initially, *States* contains just the *single* triple  $(\emptyset, s_0, v_0)$ , where  $\emptyset$  is the empty partial matching, and  $s_0$  (initial automaton state) and  $v_0$  (every clock set to 0) are as in the initialization of the baseline algorithm. For  $i = 1, \dots, n$ , we receive the next snapshot  $H_{t_i}$  and do the following:

1. From previous iterations, *States* contains all tuples  $(\mu, s, v)$ , where  $\mu$  is a maximal partial matching of  $P$  in  $H_{\leq t_{i-1}}$  and  $(s, v)$  is a possible configuration of the automaton on the  $i - 1$ -th prefix of  $w_\mu$ . Now, using an incremental query processing algorithm, compute *Extend*: the set of all pairs  $(\mu, \mu')$  such that  $\mu$  appears in *States*,  $\mu'$  extends  $\mu$ , and  $\mu'$  is a maximal partial matching of  $P$  in  $H_{\leq t_i}$ .
2. With *Extend* computed in the previous iteration, update

$$\begin{aligned} \text{States} &:= \{(\mu', s, v) \mid (\mu, s, v) \in \text{States} \\ &\quad (\mu, \mu') \in \text{Extend}\} \end{aligned}$$

by a project-join query. By Proposition 1, *States* now contains all tuples  $(\mu, s, v)$ , where  $\mu$  is a maximal partial matching of  $P$  in  $H_{\leq t_i}$  (as opposed to  $H_{\leq t_{i-1}}$ ) and  $(s, v)$  is a possible configuration of the automaton on the  $i - 1$ -th prefix of  $w_\mu$ .

3. Exactly as in the baseline and on-demand algorithms, *States* is now updated for the  $i$ -th letter of the timed words of all partial matchings.

Note that in step 1 above, it is possible that  $\mu' = \mu$ , which happens when the new snapshot does not contain any edges useful for extending  $\mu$ , or when  $\mu$  is already a total matching. On the other hand, when  $\mu$  can be extended, there may be many different possible extensions  $\mu'$ , and table *States* will grow in size.

**Example 7** We now illustrate the algorithm for the same BGP and temporal constraint as in Examples 5 and 6. Figure 6 shows the *States* relation with the timed words at  $t = 0, \dots, 3$ . (To streamline presentation, we omit edges that are not part of any cycle of length 2, but note that there are 16 such partial matchings in this relation). At  $t = 0$ , we only have one partial matching, denoted by  $\emptyset$ . At time  $t = 1$ ,  $e_5$  is active for the first time, and we create two partial matchings  $(e_5, -)$  and  $(-, e_5)$ . For  $(e_5, -)$ ,  $b_1 = 1$  and, since the second edge is not set yet,  $b_2 = 0$ . Based on this, the automaton will update this matching state to  $s_1$  and set the clock to 0. For  $(-, e_5)$ , we have  $b_1 = 0$  and  $b_2 = 1$ , and, as there is no transition in the automaton for this situation, this partial matching is dropped early. At  $t = 2$ , two new edges  $e_6$  and  $e_8$  are observed, and  $(e_6, -)$ ,  $(-, e_6)$ ,  $(e_8, -)$ ,  $(-, e_8)$  partial matching are added to *Extend*. Additionally,  $e_6$  can extend  $(e_5, -)$ , creating the full matching  $(e_5, e_6)$ . In this

	$\mu$	$s$	$v$	$b_1 b_2$
t=0	$\emptyset$	$s_0$	$\square$	00
t=1	$\emptyset$	$s_0$	$\square$	00
	$(e_5, -)$	$s_1$	$[x=1]$	10
	$(-, e_5)$		$\square$	01
...				
t=2	$\emptyset$	$s_0$	$\square$	00
	$(e_5, -)$	$s_1$	$[x=1]$	00
	$(e_6, -)$	$s_1$	$[x=0]$	10
	$(-, e_6)$		$\square$	01
	$(e_8, -)$	$s_1$	$[x=0]$	10
	$(-, e_8)$		$\square$	01
	$(e_5, e_6)$	$s_0$	$[x=0]$	01
t=3	$\emptyset$	$s_0$	$\square$	00
	$(e_5, -)$		$[x=2]$	10
	$(e_6, -)$	$s_1$	$[x=1]$	00
	$(e_8, -)$		$[x=1]$	10
	$(e_5, e_6)$	$s_1$	$[x=1]$	10

**Fig. 6** Content of the *States* relation at  $t = 0, \dots, 3$ , illustrating the execution of the timed automaton in Fig. 2c by the partial matching algorithm of Sect. 4.4

timepoint, as there is no transition for  $(-, e_6)$  and  $(-, e_8)$ , they are rejected. At  $t = 3$ ,  $e_8$  is active and the partial matching  $(e_8, -)$  is rejected. Another observation is that, at  $t = 3$  we see  $e_5$  again, and so we have  $b_1 = 1, b_2 = 0$ . We thus drop the partial matching  $(e_5, -)$ , since no edge can extend this matching. An early rejection such as this can reduce the computation time for the partial matching algorithm. For matching  $(e_5, e_6)$ , the algorithm continues as in Example 5, producing the same result.

**Correctness and Complexity** Correctness of the partial-match algorithm was already established by Proposition 1. As to the complexity, the same analysis applies as for the on-demand algorithm, except that we must replace  $M_i$  by  $M'_i$ , the number of new *partial* matches. Since there may be more partial matches than total matches, we see that theoretically the complexity of partial-match is not better than that of on-demand.

Of course, partial-match does avoid the catching up on the automaton configurations of newly found matches, which was needed in on-demand. But this win is a constant factor and hard to reflect accurately in the asymptotic complexity. Nevertheless, again, in our experiments, we will see that partial-match does often perform better than on-demand, depending on characteristics of datasets and BGPs, which influence whether a large number of partial matches will be generated that never can be completed to a total match.

### 4.5 Avoiding quadratic blowup

A well-known problem with partial BGP matching, in the non-temporal setting, is that the number of partial matchings may be prohibitively large.

For a simple example, consider matching a path of length 3,  $x_1 \xrightarrow{y_1} x_2 \xrightarrow{y_2} x_3 \xrightarrow{y_3} x_4$ , in some graph  $G$ . Note that any edge in  $G$  gives rise to a partial matching for  $y_1$ ,  $y_2$ , and  $y_3$ . What is worse, however, is that any pair of edges gives rise to a partial matching for  $y_1$  and  $y_3$  together. We thus immediately get a quadratic number of partial matchings, irrespective of the actual topology of the graph  $G$ . For example,  $G$  may have no 3-paths at all, or even no 2-paths. Such a quadratic blowup may not occur for  $y_1$  and  $y_2$  together. Indeed, since  $y_1$  and  $y_2$  form a connected subpattern, only pairs of adjacent edges give rise to a partial matching.

Of course, in the above example,  $G$  may still have many 2-paths but very few 3-paths, so connectivity is not a panacea. Still, we may expect connected subpatterns to have a number of partial matchings that is more in line with the topology of the graph. At the very least, working only with connected subpatterns avoids generating the Cartesian product of sets of partial matchings of two or more disconnected subpatterns.

Interestingly, in the temporal setting, the very presence of a temporal constraint (timed automaton) may avoid disconnected partial matchings. This happens when the temporal constraint enforces that only partial matchings of connected subpatterns can ever satisfy the constraint, allowing early rejection of when partial matchings of disconnected subpatterns. We can formalize the above hypothesis as follows.

Consider a temporal BGP  $Q = (P, \Gamma)$ . As usual, let  $Y$  be the set of edge variables of  $P$ . Consider a total ordering  $<$  on  $Y$ . We say that:

- $<$  is connected with respect to  $P$  if, for every  $y \in Y$ , the subgraph of  $P$  induced by all edge variables  $z \leq y$  is connected.
- $<$  is compatible with  $\Gamma$  if, for any  $y_1 < y_2$  in  $Y$ , and any timed word  $w$  satisfying  $\Gamma$  in which both  $y_1$  and  $y_2$  appear, the first position in  $w$  where  $y_1$  appears does not come after the first position where  $y_2$  appears.

Now, when a connected, compatible ordering is available, we can modify the partial-match algorithm in the obvious manner so as to focus only on partial matchings based on the subsets of variables  $\{y_1, \dots, y_j\}$ , for  $1 \leq j \leq n$ . By the connectedness property, we avoid Cartesian products. Moreover, by the compatibility property, we do not lose any outputs.

As a simple example, consider the path of length 3 BGP and the timed automaton  $\Gamma$  from Fig. 8. The ordering  $y_1 < y_2 < y_3$  is connected with respect to  $P$ , and is compatible with  $\Gamma$ . So our theory would predict the partial matching algorithm to work well for this temporal BGP  $(P, \Gamma)$ . We will show effectiveness of the partial matching algorithm in Sect. 6.5.

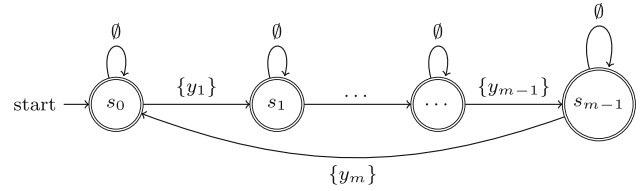


Fig. 7  $TA_0$  generalizes  $TA_1$ , specifying that edges should appear repeatedly in the given temporal order

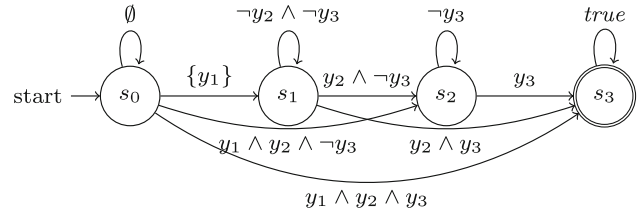


Fig. 8  $TA_4$ : Each of  $y_1, y_2$  and  $y_3$  is active at some point, with first time  $y_1 \leq$  first time  $y_2 \leq$  first time  $y_3$

Whether or not an ordering of the edge variables is connected with respect to  $P$  is straightforward to check, by a number of graph connectivity tests. Moreover, when  $\Gamma$  is given by a timed automaton, it also possible to effectively check whether an ordering is compatible with  $\Gamma$ .

*Verifying compatibility* We offer the following algorithm for verifying that an ordering  $y_1 < \dots < y_m$  is compatible with a temporal constraint  $\Gamma$ , specified by a timed automaton.

1. Compute an automaton defining the intersection of  $\Gamma$  with all regular languages of the form

$$(\neg y_j \wedge \neg y_i)^* \cdot (y_j \wedge \neg y_i) \cdot true^* \cdot y_i \cdot true^*,$$

for  $1 \leq i < j \leq n$ . These languages contain the words where both  $y_i$  and  $y_j$  appear, but  $y_j$  appears first, which we do not want when  $i < j$ .

2. The resulting timed automaton should represent the empty language, i.e., should not accept any timed word.

Effective algorithms for computing the intersection of timed automaton and for emptiness checking are known [2]. Note that it actually suffices here to intersect a timed automaton ( $\Gamma$ ) with an NFA (the union of the regular languages from step 1). An interesting question for further research is to determine the precise complexity of the following problem:

**Problem:** Compatible and connected ordering  
**Input:** A temporal BGP  $(P, \Gamma)$   
**Output:** An ordering of the edge variables that is connected w.r.t.  $P$  and compatible with  $\Gamma$ , or ‘NO’ if no such ordering exists.

## 5 Implementation

The algorithms described in Sect. 4 have been implemented using Rust and the `Itertools` library [52] as a single-threaded application. Our algorithms are easy to implement using any system supporting the dataflow model such as Apache Spark [71], Apache Flink [11], Timely Dataflow [43], or Differential Dataflow [40].

A temporal graph is stored on disk as relational data in CSV files *Node*, *Edge*, and *Active*. In the initial stage of the program, we load all data into memory, loading edges into two hash-tables with *vid1* and *vid2* as keys. We added a “first” meta-property field to the *Edge* relation and use it for lazy evaluation of matchings in the baseline and on-demand algorithms.

**BGP matching** We implement BGP matching as a select-project-join query. For cyclic BGPs such as triangles and rectangles, we use worst-case optimal join [3], meaning that instead of the traditional pairwise join over the edges, we use a vertex-growing plan. We use a state-of-the-art method in our implementation but note that (non-temporal) BGP matching in itself is not the focus of this paper, and so any other BGP matching algorithm can be used in conjunction with the timed automata-based methods we describe.

On-demand and partial matching algorithms are both designed to work in online mode, computing new matchings at each iteration. To implement online matching for the on-demand algorithm, we build on join processing in streams [67]. We can use information from the temporal constraint to avoid useless joins in the incremental computation of matchings. For example, consider two edge variables *y* and *z* coming from the BGP. With *E* the current history of active edges and  $\Delta E$  the edges from the new snapshot, we must in principle update the join of  $\rho_y(E)$  with  $\rho_z(E)$  by three additional joins  $\rho_y(E) \bowtie \rho_z(\Delta E)$ ;  $\rho_y(\Delta E) \bowtie \rho_z(E)$ ; and  $\rho_y(\Delta E) \bowtie \rho_z(\Delta E)$ . When the temporal constraint implies, for example, that *y* is never active before *z*, the first of these three additional joins can be omitted. Such order information can be inferred from a timed automaton using similar techniques already described in the paragraph on verifying compatibility in Sect. 4.5.

**Timed automata** We represent a timed automaton as a relation *Automaton*( $s_c, \theta, \delta, R, s_n$ ), in which each tuple corresponds to a transition from the current state  $s_c$  to the next state  $s_n$ . For example, the timed automaton of Fig. 2c is represented as follows:

$s_c$	$\theta$	$\delta$	<i>R</i>	$s_n$
0	00	<i>true</i>	[]	0
0	10	$c.0 < 3$	[0]	1
1	01	$c.0 < 3$	[0]	0
1	00	<i>true</i>	[]	1

The specification of the timed automaton is loaded into memory as a hash table, with  $(s_c, \theta)$  as the key. The timed word  $\theta$  (see Sect. 3), is encoded as a bitset. For example, in the timed automaton in Fig. 2c, we encode  $y_1 \wedge \neg y_2$  as 10, where the first bit corresponds to  $y_1$  and the second to  $y_2$ . If the transition condition is *true* then, for a matching with two variables, we add 4 rows to *Automaton*, one for each 00, 01, 10, and 11. Using bitsets makes automaton transitions efficient, as we will show in Sect. 6.6. Table *Automaton* also stores the clock acceptance condition  $\delta$ , and a nested field *R* with an array of clocks to be reset during the transition to the next state. To update the state of a matching, we execute a hash-join followed by a projection between *Automaton* and *States*.

Updating the clock for each matching will be computationally expensive. Instead, during the automaton transition, for each matching, we store the current time (of last snapshot visited) value for that clock instead of setting it to zero. This way, instead of updating all clocks in every iteration, we can just get the correct value of the clock when needed and compute the current value of the clock by subtracting the value of the clock from the current time.

In many temporal graphs, due to the nature of their evolution, most edges appear for the first time during the last few snapshots. To optimize performance we implemented a simple but effective optimization for our baseline and on-demand algorithms: when the initial state of the timed automaton self-loops on the empty letter, we will not run on a matching until at least one of its edges is seen. This can be determined using the “first” meta-property of the *Edge* relation. This optimization is not necessary for the partial matching algorithm, where it is essentially already built-in.

We also implement the early acceptance and early rejection optimizations.

## 6 Experiments

We now describe an extensive experimental evaluation with several real datasets and temporal BGPs, and demonstrate that using timed automata is practical. We investigate the relative performance of our methods, and compare them against two state-of-the-art in-memory relational systems, DuckDB [49] and HyPer [44, 45].

*In summary*, we observe that the on-demand and partial-match algorithms are effective at reducing the running time compared to the baseline. Interestingly, while no single algorithm performs best in all cases, the trade-off in performance can be explained by the properties of the dataset, of the BGP, and of the temporal constraint. Our results indicate that partial-match is most efficient for acyclic BGPs such as paths of bounded length, while on-demand performs best for cyclic BGPs such as triangles, particularly when evaluated over

sparse graphs. Interestingly, the performance gap between on-demand and partial-match is reduced with increasing graph density or BGP size, and partial-match outperforms on-demand in some cases. We also show that algorithm performance is independent of timed automaton size and of the number of clocks.

We show that our methods substantially outperform state-of-the-art relational implementations in most cases. We also demonstrate that temporal BGPs are more concise than the corresponding relational queries, pointing to better usability of our approach.

*Experimental setup* Our algorithms were executed as single-thread Rust programs, running on a single cluster node with an Intel Xeon Platinum 8268 CPU, using the Slurm scheduler [70]. We used DuckDB v.0.3.1 and the HyPer API 0.0.14109 provided by Tableau.<sup>2</sup> All systems were run with 32GB of memory on a single CPU. Execution time of DuckDB and HyPer parsing, optimizing and executing the SQL query, and does not include the time to create database tables and load them into memory. Similarly, execution time of our algorithms includes loading the timed automaton and executing the corresponding algorithm. All execution times are averages of 3 runs; the coefficient of variation of the running times was under 10% in most cases, and at most 12%.

**Remark 2** Since we are dealing with graph data and BGPs, one may ask why we implemented our algorithms in a dataflow environment, and compare to relational systems. Why not work on top of a graph database system, and compare to graph databases? The reason is that a BGP is, in essence, a multiway join query, for which the best performance is realized with the help of worst-case optimal join algorithms, or relational query processors with very good optimization. It is exactly with respect to these two environments that we conduct our experiments. On the other hand, the main advantage of graph database systems is their support of reachability queries or regular path queries, which are not part of our basic notion of BGP. Rather, our contributions lie in expressive temporal filtering of the matchings of a BGP, for which our experimental set-up provides a suitable empirical evaluation.

*Datasets.* Experiments were conducted on 4 real datasets, summarized in Table 1, where we list the number of distinct nodes and edges, temporal domain size (“snaps”), the number of active edges across snapshots (“active”), structural density (“struct”, number of edges in the graph, divided by number of edges that would be present in a clique over the same number of nodes), and temporal density (“temp”, number of timepoints during which an edge is active, divided by temporal domain size, on average).

**Table 1** Description of the experimental datasets

	Nodes	Edges	Active	Snaps	Density	
					Struct	Temp
EPL	50	1500	35K	25	0.6	0.93
Contact	541	3349	21K	48	0.16	0.13
Bitcoin	1704	4845	268K	1036	0.026	0.049
Email	776	65K	1.9M	800	0.1	0.03
<b>FB-Wall</b>	46K	264K	856K	850	0.0001	0.003
Superuser	194k	854K	1.2M	2774	0.00002	0.0005

*EPL*, based on the English soccer dataset [19], represents 34,800 matches between 50 teams over a 25-year period. We represent this data as a temporal graph with 1-year temporal resolution, where each node corresponds to a team and a directed edge connects a pair of teams that played at least one match during that year. The direction of the edge is from a team with the higher number of goals to the team with the lower number of goals in the matches they played against each other that year; edges are added in both directions in the case of a tied result. This is the smallest dataset in our evaluation, but it is very dense both structurally and temporally.

*Contact* is based on trajectory data of individuals at the University of Calgary over a timespan of 4 h [47]. We created a bipartite graph with 500 person nodes and 41 location nodes, where the existence of an edge from a person to a location indicates that the person has visited the location. The original dataset records time up to a second. To make this data more realistic for a contact tracing application, we made the temporal resolution coarser, mapping timestamps to 5-min windows, and associated individuals with locations where they spent at least 2.5 min.

*Bitcoin* is a temporal graph based on the decentralized digital currency and payment system [33, 48]. In this graph, nodes represent bitcoin addresses, and directed edges from  $u$  to  $v$  at time  $t$  signify that bitcoin was transferred from address  $u$  to address  $v$  at time  $t$ . In-line with the work by Semertzidis and Pitoura [55], we used data up to December 31, 2013. The graph used in our experiments has 1,704 nodes, 4,845 unique edges, and 268K active edges over a 1,036-day period.

*Email*, based on a dataset of email communications within a large European research institution [37], represents about 1.9 M email messages exchanged by 776 users over an 800-day period, with about 65K distinct pairs of users exchanging messages. This dataset has high structural density (10% of all possible pairs of users are connected at some point during the graph’s history), and intermediate temporal density (3%).

*FB-Wall*, derived from the Facebook New Orleans user network dataset [62], represents wall posts of about 46K

<sup>2</sup> [https://help.tableau.com/current/api/hyper\\_api](https://help.tableau.com/current/api/hyper_api)

users over a 850-day period, with 264K unique pairs of users (author / recipient of post).

*Superuser* is based on StackOverflow, a stack exchange website where users post questions and receive answers from other users. In this temporal graph, derived by Paranjape et al. [48], nodes represent users, and a directed edge from  $u$  to  $v$  at time  $t$  signifies that user  $u$  answered a question from user  $v$ , or that  $u$  commented on  $v$ 's question or on  $v$ 's answer. This graph has 194K nodes, 854K edges, and 1.2 M temporal edges that span over 2774 days, making this the largest dataset in our experiments, both in the number of edges and in duration.

We also use synthetic datasets to study the impact of data characteristics on performance, and describe them in the relevant sections.

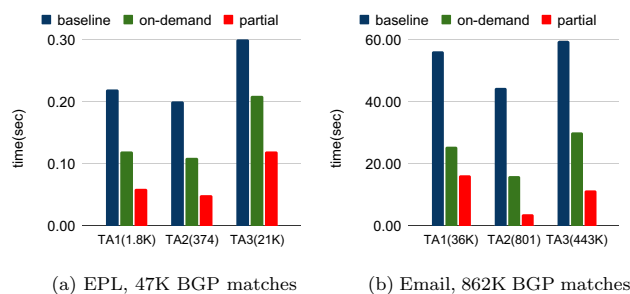
**Example 8** Our chosen test datasets illustrate the practical applicability of temporal BGPs with expressive temporal constraints on real-life data. Indeed, recall Examples 2 and 3 from the Introduction. Example 2 looks for patterns of “ping-pong email messages,” which is naturally applicable to datasets Email, FB-Wall, and Superuser. Then Example 3, looking for patterns of location usage without overlap (using automaton  $TA_6$  from the Introduction), is naturally applicable to the Contact dataset. Furthermore, on the Bitcoin dataset, we could envisage looking for “money carousels” where one party transfers money to a second party, the second party subsequently transfers money to a third party, upon which the third party transfers money back to the first party, and this cycle can repeat. This uses as BGP a cycle of length three, and the automaton  $TA_0$  from Fig. 7. Our experiments, interestingly, found out that no such carousel exists in the Bitcoin dataset.

## 6.1 Relative performance of the algorithms

In our first set of experiments, we evaluate the relative performance of the baseline (Sect. 4.2), on-demand (Sect. 4.3), and partial-match (Sect. 4.4) algorithms. Note that the baseline algorithm can only be used when a graph's evolution history is fully available (rather than arriving as a stream), and that partial-match is only used when the matching is guaranteed to be connected (Sect. 4.5).

We use the BGP that looks for paths of length 2, with timed automata  $TA_1$ ,  $TA_2$  and  $TA_3$  from Fig. 2. Automaton  $TA_3$  is interesting for showing the impact of early acceptance and rejection on performance.

Figure 9 shows the execution time of the baseline, on-demand and partial-match algorithms for EPL and Email, also noting the number of temporal matches. The BGP, which is in common for all executions in this experiment, returns 47K matches on EPL and 862K matches on Email. When the temporal constraint is applied, the number of matches is



**Fig. 9** Running time for path of length 2 with timed automata  $TA_1$ ,  $TA_2$ ,  $TA_3$

reduced, and is presented on the  $x$ -axis. For example,  $TA_1$  returns 1.8K matches on EPL and 36K on Email.

We observe that partial-match is the most efficient algorithm for all queries and both datasets, returning in under 0.12 sec for EPL, and in under 17 sec for Email in all cases. The on-demand algorithm outperforms the baseline in all cases, but is slower than partial-match. The performance difference between the baseline and on-demand is due to a join between two large relations in baseline, compared to multiple joins over smaller relations in on-demand.

We also observe that the relative performance of the algorithms depends on the number of matches, and explore this relationship further in the next experiment. To compare algorithm performance across BGPs and datasets, we use the timed automaton  $TA_0$  of Fig. 7, which generalizes  $TA_1$  from 2 to  $m$  edges.  $TA_0$  specifies that edges in a matching should appear repeatedly in a strict temporal order. We use  $TA_0$  (with  $m = 2, 3$ , or 4 as appropriate) as the temporal constraint for paths of length 2 and 3, and for cycles of length 2, 3, 4, for three of the datasets. Because the Contact dataset is a bipartite graph, we used it for in-star ( $x_2 \xrightarrow{y_1} x_1 \xleftarrow{y_2} x_3$ ) and out-star ( $x_2 \xleftarrow{y_1} x_1 \xrightarrow{y_2} x_3$ ) BGPs of size 2 and 3. For Superuser, the number of matches for path2 and path3 exceeded our available memory, and we used the o-star2 pattern instead in our experiments. (We will discuss memory usage in Sect. 6.2 below.)

Table 2 summarizes the results. It shows number of BGP matchings (“BGP”), number of matchings accepted by  $TA_0$  (“BGP+ $TA_0$ ”), and running times of computing the BGP match only (“match”), and of computing both BGP and temporal matches using to the baseline, on-demand, and partial-match algorithms.

We observe that, for acyclic patterns (e.g., paths, i-star, o-star), partial-match is significantly faster than on-demand and baseline. For such patterns, partial matchings are shared by many total matchings and by larger partial matchings, benefiting the running time. Interestingly, for cycles of size 2 and 3, on-demand is fastest, followed by baseline. However, for cycles of size 4 partial-match is once again the fastest algorithm. The reason for this is that there are far fewer cycles than

**Table 2** Relative performance of baseline, on-demand and partial-match, for different BGPs and automaton  $TA_0$ 

Pattern	# of matches		Time (s)				Memory (GB)		
	BGP	BGP+TA	Match	Baseline	On-demand	Partial	Baseline	On-demand	Partial
EPL									
path2	47K	1.8K	0.01	0.22	0.12	<b>0.06</b>	0.005	0.005	0.005
path3	1.5 M	4.6K	0.21	5.11	3.66	<b>0.32</b>	0.005	0.005	0.005
cycle2	1.1K	35	0.01	<b>0.03</b>	<b>0.01</b>	0.02	0.005	0.005	0.005
cycle3	35K	66	0.02	0.12	<b>0.09</b>	0.36	0.005	0.005	0.005
cycle4	1.1 M	106	0.22	3.31	2.92	<b>0.53</b>	0.005	0.005	0.005
Email									
pattern	# of matches		Time (sec)				Memory (GB)		
	BGP	BGP+TA	match	baseline	on-demand	partial	baseline	on-demand	partial
path2	862K	36K	0.25	56.27	25.56	<b>16.35</b>	0.17	0.01	0.01
path3	42 M	126K	41.00	1475.20	766.33	<b>74.77</b>	10.00	4.00	0.30
cycle2	18K	843	0.12	0.92	<b>0.54</b>	2.41	0.02	0.02	0.02
cycle3	205K	309	0.19	5.13	<b>3.49</b>	14.34	0.02	0.02	0.02
cycle4	8.4 M	1352	4.38	196.50	125.85	<b>92.40</b>	1.50	1.00	1.20
Bitcoin									
Pattern	# of matches		Time (sec)				Memory (GB)		
	BGP	BGP+TA	Match	Baseline	On-demand	Partial	Baseline	On-demand	Partial
path2	89K	1387	0.19	3.93	2.28	0.04	0.005	0.005	0.005
path3	1.1 M	820	0.19	21.87	20.96	<b>1.70</b>	0.005	0.005	0.005
cycle2	839	6	0.08	<b>0.15</b>	0.17	0.21	0.005	0.005	0.005
cycle3	3621	0	0.09	<b>0.16</b>	0.34	0.24	0.005	0.005	0.005
Contact									
pattern	# of matches		Time (sec)				Memory (GB)		
	BGP	BGP+TA	Match	Baseline	On-demand	Partial	Baseline	On-demand	Partial
i-star2	10.1K	0	0.01	0.03	0.03	<b>0.01</b>	0.01	0.01	0.01
o-star2	207K	0	0.03	0.84	0.79	<b>0.67</b>	0.01	0.01	0.01
i-star3	21.4K	0	0.01	0.07	0.06	<b>0.03</b>	0.01	0.01	0.01
o-star3	10.7 M	0	2.86	36.24	35.26	<b>1.73</b>	15.00	10.00	0.50
FB-Wall									
pattern	# of matches		Time (sec)				Memory (GB)		
	BGP	BGP+TA	Match	Baseline	On-demand	Partial	Baseline	On-demand	Partial
path2	4.4 M	681K	1.08	839.30	387.31	<b>318.14</b>	1.20	0.62	0.37
path3	91 M	235K	13.44	9477.47	5093.05	<b>998.54</b>	27.20	21.60	2.62
cycle2	160K	16K	0.94	12.34	<b>8.32</b>	66.76	0.01	0.01	0.22
cycle3	272K	4.1K	0.84	17.69	<b>10.94</b>	341.39	0.01	0.01	0.57
Superuser									
pattern	# of matches		Time (sec)				Memory (GB)		
	BGP	BGP+TA	Match	Baseline	On-demand	Partial	Baseline	On-demand	Partial
cycle2	280K	36K	10.00	79.49	<b>59.31</b>	383.56	0.18	0.17	0.35
cycle3	2 M	100K	113.52	561.63	<b>383.56</b>	1372.80	0.35	0.30	3.20
o-star2	43 M	22 M	113.52	> 12h	28,855.00	<b>14,189</b>	17.4	10.1	3.20

Best-performing results are highlighted in bold

possible partial matchings, and in smaller cycles this causes partial-match to run slower. As cycle size increases, performance of partial-match becomes comparable to, or better, than of the other two algorithms. Another graph characteristic that can affect partial-match performance is graph density, which we discuss in the next section. (Our machine's RAM could not fit cycle4 for FB-Wall, so we did not conduct that experiment.)

Finally, we note that—not surprisingly—intermediate result size and output size substantially impact the running time of all algorithms. This is particularly clear in our two largest datasets, FB-Wall and Superuser. Note the Superuser is the most challenging because it is both the largest dataset in terms of the number of edges, and also has the largest temporal domain ( $3\times$  compared to FB-Wall). For Superuser, the o-star2 pattern produced 22 M temporal matches. This com-

**Table 3** Relative performance of baseline, on-demand and partial-match, for different BGPs, and for automata  $TA_3$  and  $TA_5$ , over FB-Wall and Superuser

Dataset	Pattern	TA	# matches	Time (s)			Memory (GB)		
				Baseline	On-demand	Partial	Baseline	On-demand	Partial
FB-Wall	Cycle3	ta3	8910	11.87	19.66	306.00	0.10	0.10	0.57
FB-Wall	Cycle3	ta5	57	3.69	5.87	4.49	0.10	0.10	0.10
FB-Wall	Path3	ta3	4,339,730	9758	4888	1685	22	17	2.6
FB-Wall	Path3	ta5	52,626	3048	1344	20.25	22	15	0.10
Superuser	Cycle3	ta3	151,883	767	425	751	0.47	0.32	3.61
Superuser	Cycle3	ta5	147,212	175	191	989	0.35	0.30	3.2
Superuser	Path2	ta3	0	<b>OM</b>	26,086	118	<b>OM</b>	30.2	0.2
Superuser	Path2	ta5	705,005	<b>OM</b>	<b>OM</b>	1322	<b>OM</b>	<b>OM</b>	0.70

Best-performing results, and cases where an out-of-memory (OM) error was encountered, are highlighted in bold

putation did not complete under the baseline algorithm after 12h, took about 4h for the partial match, and about 8h for the on-demand algorithm.

## 6.2 Memory usage

Table 2 also reports the maximum memory usage of our proposed algorithms for different BGPs and datasets, for the timed automaton  $TA_0$ . We start with some general observations. The memory used by our baseline algorithm is a function of the number of matchings, as it produces all the matching at time zero. For acyclic patterns using partial-match, we can see that in most cases maximum memory use is a function of the number of accepted matching. The partial-match algorithm produces matchings incrementally and, while it keeps partial matching, because of early rejection it will remove the non-accepted and partial matchings as early as possible. The on-demand algorithm follows the same outline as partial-match, however, because it also needs to keep track of active edges from the previous time points to run the automata on demand, it will have higher memory overhead. For cyclic patterns, especially on larger graphs where most matching will ultimately be accepted, the partial-match algorithm needs a larger amount of memory because the number of partial matchings is higher (and often substantially so) as compared to the number of total matchings.

According to Table 2, the relatively small datasets EPL and Bitcoin run with less than 5MB of memory for all BGPs in our experiments, and with no significant difference in terms of memory overhead among the different algorithms. For the Email dataset, patterns cycle2 and cycle3 return relatively few matchings, and can be computed using less than 20MB of memory. For patterns that return more results, such as path3, the baseline algorithm needs 10GB of memory, the on-demand needs 4GB, and partial-match can run with less than 300MB of memory. For the cycle4 pattern, the baseline algo-

rithm needs 1.5GB, on-demand needs 1GB and partial-match needs 1.2GB, due to a higher number of partial matchings.

An interesting observation in Table 2 is for the Contact dataset, with o-star3. For this pattern, baseline and on-demand need a significant amount of memory due to the large number of matchings, but partial-match can process it with  $20\times$  less memory. This is because partial-match benefits from there being to accepted o-star2 partial patterns in this case. For the cyclic patterns over FB-Wall and Superuser datasets, our largest datasets, baseline is more memory-efficient for cyclic patterns, while partial-match is more memory-efficient for acyclic patterns (e.g., o-star2, path2, and path3).

## 6.3 Additional scalability experiments

FB-Wall and Superuser are the two largest datasets in our experiments. Here, we further investigate the relative performance of our algorithms over these datasets with several patterns and two timed automata,  $TA_3$  and  $TA_5$ . We used cycle3 and path3 for FB-Wall, and cycle3 and path2 for Superuser. The path3 pattern on Superuser returns more than 1.8 billion matchings, and we could not process it within 32GB of memory that we had available.

Path3 returns 90 million matchings on FB-Wall, and path2 returns 180 million matchings on Superuser which, considering the large temporal domain of this graphs, is challenging to handle. Table 3 shows the result of this experiment. Similarly to previous experiments, baseline and on-demand show better performance than partial-match for cyclic patterns, in terms of both processing time and memory overhead. Acyclic patterns are where partial-match works best: for path3 on FB-Wall, partial-match computed the matching more than  $67\times$  faster than the other algorithms, with  $100\times$  less memory. For Superuser under  $TA_3$ , partial-match computed the result  $800\times$  faster using 200MB of memory, while on-demand needed 26GB of memory and baseline ran out of memory.



**Table 4** Best-performing timed-automaton algorithm (TAA) compared to DuckDB v.0.3.1 and HyPer API v.0.0.14109, for SQL queries in Fig. 1

	TA	# matches	Time (s)		
			TAA	DuckDB	HyPer
EPL Path2	$TA_e$	35,866	<b>1.09</b>	1.11	1.44
	$TA_1$	1801	<b>0.06</b>	60.14	17.84
	$TA_2$	374	<b>0.05</b>	69.89	11.36
	$TA_3$	21,035	<b>0.04</b>	0.07	0.13
	$TA_4$	29,726	<b>0.06</b>	0.07	0.13
	$TA_5$	1714	<b>0.07</b>	18.57	0.39
	$TA_6$	19,578	0.54	0.12	<b>0.09</b>
	$TA_7$	257	<b>1.3</b>	8.22	5.56
EPL Cycle2	$TA_e$	933	0.04	0.597	<b>0.01</b>
	$TA_1$	35	<b>0.01</b>	2.67	0.41
	$TA_2$	22	<b>0.01</b>	3.33	0.26
	$TA_3$	418	<b>0.01</b>	0.1	<b>0.01</b>
	$TA_4$	740	<b>0.01</b>	0.1	<b>0.01</b>
	$TA_5$	90	<b>0.02</b>	0.53	1.19
	$TA_6$	312	0.07	0.07	<b>0.01</b>
	$TA_7$	0	<b>0.05</b>	1.18	5.15
Bitcoin Path2	$TA_e$	80,008	118.53	470.81	<b>5.13</b>
	$TA_1$	1387	<b>0.04</b>	<b>OM</b>	<b>OM</b>
	$TA_2$	34	<b>0.08</b>	<b>OM</b>	<b>OM</b>
	$TA_3$	56,943	<b>3.31</b>	24.45	6.89
	$TA_4$	57,782	<b>3.32</b>	24.96	6.86
	$TA_5$	452	<b>0.94</b>	<b>OM</b>	20.19
	$TA_6$	53,524	24.25	7.20	<b>1.17</b>
	$TA_7$	2083	<b>41.28</b>	<b>OM</b>	<b>OM</b>
Bitcoin Cycle2	$TA_e$	801	1.36	286.99	<b>0.39</b>
	$TA_1$	6	<b>0.15</b>	<b>OM</b>	<b>OM</b>
	$TA_2$	0	<b>0.15</b>	<b>OM</b>	4471
	$TA_3$	184	<b>0.11</b>	0.80	0.38
	$TA_4$	654	<b>0.12</b>	0.79	0.39
	$TA_5$	438	0.30	15.83	<b>0.10</b>
	$TA_6$	124	<b>0.17</b>	7.64	0.38
	$TA_7$	0	<b>0.73</b>	<b>OM</b>	<b>OM</b>
Contact O-Star2	$TA_e$	169,410	<b>8.87</b>	22.43	19.29
	$TA_1$	0	<b>0.67</b>	<b>OM</b>	316.5
	$TA_2$	0	<b>0.72</b>	<b>OM</b>	205.95
	$TA_3$	101,268	<b>0.59</b>	0.95	0.69
	$TA_4$	107,650	<b>0.68</b>	0.91	<b>0.68</b>
	$TA_5$	4154	<b>0.72</b>	<b>OM</b>	2.035
	$TA_6$	10,832	4.68	0.99	<b>0.54</b>
	$TA_7$	76	<b>15.12</b>	<b>OM</b>	102.89
Email Path2	$TA_e$	719,609	501.07	649.69	<b>239.82</b>
	$TA_1$	35,594	<b>22.35</b>	<b>OM</b>	<b>OM</b>

**Table 4** continued

	TA	# matches	Time (s)		
			TAA	DuckDB	HyPer
	$TA_2$	801	<b>3.73</b>	<b>OM</b>	1748.85
	$TA_3$	443,431	<b>7.78</b>	17.84	8.59
	$TA_4$	455,977	<b>7.82</b>	17.78	8.26
	$TA_5$	2474	<b>1.27</b>	<b>OM</b>	48.87
	$TA_6$	693,956	320.04	8.98	<b>5.61</b>
	$TA_7$	390	<b>327.33</b>	<b>OM</b>	<b>OM</b>
	$TA_8$	12,919	19.32	<b>8.02</b>	21.58
Email Cycle2	$TA_e$	14,188	<b>4.71</b>	254.7	15.3
	$TA_1$	843	<b>0.83</b>	<b>OM</b>	1788.87
	$TA_2$	240	<b>0.64</b>	<b>OM</b>	1733.05
	$TA_3$	6333	0.46	<b>0.38</b>	0.68
	$TA_4$	11,396	0.44	<b>0.38</b>	0.68
	$TA_5$	1800	<b>0.97</b>	206.26	1.19
	$TA_6$	4230	24.09	8.63	<b>5.48</b>
	$TA_7$	134	<b>7.1</b>	<b>OM</b>	<b>OM</b>
	$TA_8$	3441	1.79	<b>0.219</b>	49.79

Best-performing results, and cases where an out-of-memory (OM) error was encountered, are highlighted in bold

For Superuser under  $TA_5$ , both baseline and on-demand algorithm ran out of memory, while partial-match computed the result in 22 min using less than 1GB of memory.

## 6.4 Comparison to in-memory databases

In this set of experiments, we compare the running time of temporal BGP matching with equivalent relational queries. We used four datasets (EPL, Bitcoin, Contact and Email) and queried them with cyclic and acyclic BGPs of size 2, with temporal constraints specified by 9 timed automata:  $TA_e$  specifies an existential constraint (Fig. 3), while  $TA_1 \dots TA_8$  express constraints that are not existential. We did not run this set of experiments with our largest datasets, FB-Wall and Superuser, because we were unable to load these datasets into DuckDB and HyPer for processing.

We showed SQL queries QTA1 and QTA7 in the introduction. In the Supplementary Materials we give a complete listing of the SQL queries, with Path2, Cycl2 and Star2 expressing the BGPs used in our experiments, and QTAE, QTA1–QTA8 implementing the temporal constraints. Most of these queries are quite complicated compared to their equivalent timed automata.

For DuckDB and HyPer, we loaded the relations Node, Edge and Active into memory. To improve performance of DuckDB, we defined indexes on Edge(src), Edge(dst), and Active(eid,time). To the best of our knowledge, HyPer does not support indexes.

Table 4 shows the execution time for each query, comparing the running time of the best method based on timed automata (on-demand or partial-match, column “TAA”) with DuckDB and HyPer. Observe that our algorithms are significantly faster for  $TA_1$ ,  $TA_2$ ,  $TA_5$  and  $TA_7$  for all 4 datasets, and have comparable performance to the best-performing relational system for  $TA_3$  and  $TA_4$ . Relational systems outperform our algorithms on  $TA_e$  and  $TA_6$ . For  $TA_e$ , relational databases compute all possible matchings at all the time points in one shot, and then filter out those that fail the temporal constraint, which can be faster than an iterative process. Similarly, for  $TA_6$ , the XOR operator can be implemented as a join-antijoin. Interestingly, for  $TA_8$  (set containment join, see Introduction) DuckDB is most efficient, followed by our algorithms, and then by HyPer. For the majority of other cases, DuckDB either ran out of memory (OM in Table 4) or was the slowest system. Our methods were able to handle all queries within the allocated memory, while DuckDB and HyPer both ran out of memory in some cases.

## 6.5 Impact of graph properties on performance

In this set of experiments, we explore the effect of structural density and temporal domain size. We synthetically generated a complete graph with 50 nodes and 2450 edges (the same size as the EPL dataset) and temporal density of 0.5. We then sampled edges to create a graph with different structural densities. We use 4 representative BGPs: paths of length 3, and cycles of length 3. As temporal constraint we use  $TA_4$

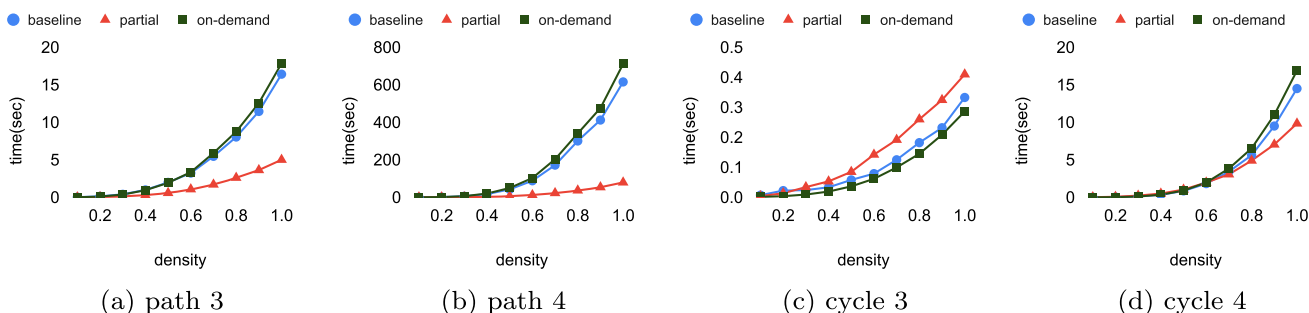


Fig. 10 Running time as a function of structural density for 4 common BGPs, with timed automaton  $TA_4$  in Fig. 8

from Fig. 8, as it has low early rejection rate, thus serving as a worst case.

Figure 10 shows the execution time of each algorithm as a function of graph density, varying from 0.1 to 1.0, where density 1.0 corresponds to a complete graph. Observe that partial-match outperforms on-demand for paths, particularly as graph density increases. For the cycle of size 3, baseline and on-demand have better performance than partial-match, but the performance gap decreases with increasing graph density. Our experiments for BGPs: path of length 4 and cycle of length 4 show similar trend. Notably, performance of on-demand is very close to the baseline.

Next, we consider the impact of temporal domain size on performance. In general, we expect execution times to increase with increasing temporal domain size. To measure this effect without changing the structure or the size of the graph, we synthetically changed the temporal resolution of the Email dataset, creating graphs with between 25 and 800 snapshots, and thus keeping the number of BGP matchings fixed.

Figure 11a shows the result of executing temporal BGP with path of length 2 and time automaton  $TA_1$  (Fig. 2b). Observe that the execution time of all algorithms increases linearly, with partial-match scaling best with increasing number of snapshots.

Finally, we study the relationship between result set size of a temporal BGP and algorithm performance. For this, we executed the path of length 2 BGP on the Email dataset, with timed automaton  $TA_2$  in Fig. 2c, and manipulated selectivity by varying the clock condition from  $c < 0$  to  $c < 1024$  on the logarithmic scale. With these settings, the temporal BGP accepts between 0 and 36K matchings. Figure 11b presents the result of this experiment, showing the running time (in sec) on the x-axis and the number of temporal BGP matchings (in thousands) on the y-axis. Observe that the running time increases linearly with increasing number of accepted matchings for all algorithms, and the slope of increase is small.

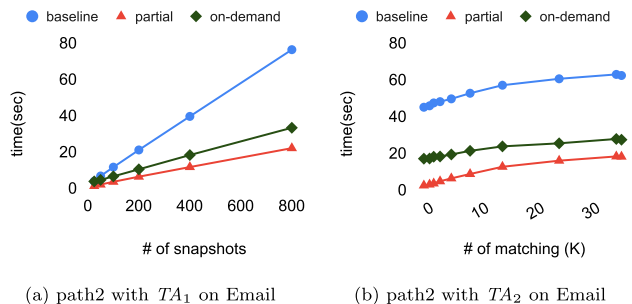


Fig. 11 a Running time vs. temporal domain size; b relationship between running time and result set size

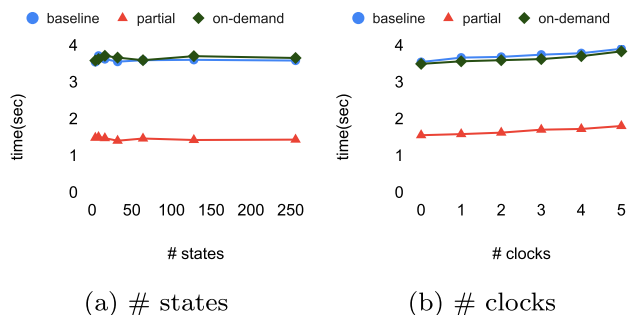


Fig. 12 Running time as a function of automaton size for cycle of size 4 with  $TA_0$  (Fig. 7) on EPL

### 6.6 Impact of the number of clocks and automaton size on performance

In our final set of experiments, we investigate the impact of automaton size and of the number of clocks on performance, while keeping all other parameters fixed to the extent possible. To do this, we fix the BGP and vary the size of the automaton, as follows. We fix the BGP to cycle4 and take  $TA_0$  (Fig. 7) with  $m = 4$  as a starting point. We can unfold the cycle of states, thus doubling the number of states but resulting in an equivalent automaton. We do this doubling seven times, until we obtain 256 states.

Figure 12a shows the execution times on the EPL dataset. Observing that the execution times remain constant, we

conclude that automaton size does not significantly impact performance.

Finally, to investigate the impact of the number of clocks on performance, we added multiple clocks to the automaton in Fig. 7 and we reset all clocks at every state transition. To ensure that any possible difference is not due to the change of output size, the condition of each clock is set to *true*. Figure 12b shows the result of this experiment, with the number of clocks on the  $x$ -axis and the execution time in seconds on the  $y$ -axis. Observe that the running time of all algorithms increases very slightly with increasing number of clocks. We thus conclude that the computational overhead of storing and updating clocks is low.

## 7 Related Work

During the past several decades, researchers considered different aspects of graph pattern matching, see Conte et al. [16] for a survey. The majority of temporal graph models use either time points [26, 34, 48] or intervals [42, 51, 68, 74] to enrich graphs with temporal information. We discuss some of these approaches below. In our work, we associate each edge in a graph with a set of time points, which is an appropriate representation when events—such as messages between users or citations—are instantaneous and so do not have a duration. It is an interesting topic for further research to investigate when and how an interval-based approach can be encoded by a point-based approach. This depends also on the considered graph model, and the considered class of queries and temporal constraints.

As an example of an interval-based approach, Xu et al. [68] consider temporal constraints that impose Allen relations [1] between the intervals of pairs of matched edges from the graph pattern. Their notion of subgraph matching is isomorphism-based, and their algorithmic approach is based on a form of signature pruning.

A prominent line of work where the point-based approach is adopted is that of mining frequent temporal subgraphs, called temporal motifs [26, 34, 48]. There, the focus is typically on existential temporal constraints, aiming to identify graph patterns with a specific temporal order among the edges, such as in our Example 1. Timed automata can easily specify such constraints. An important type of a motif is a  $\delta$  temporal motif [48], where all the edges occur inside the period of  $\delta$  time units. Timed automata can use one or multiple clocks to enforce such constraints.

An example of an approach that uses existential temporal constraints is the work of Semertzidis and Pitoura [56], who define a notion of “interaction graph pattern matching” with the help of inequalities among edge variables, and further support a “delta” constraint that imposes an upper-bound on the difference between the earliest and the latest time

points. Such constraints are expressible by timed automata. Semertzidis and Pitoura [56] propose efficient algorithms for handling their specific temporal constraints under injective semantics.

Züfle et al. [74] consider a particular class of temporal constraints where the time points within a query range are specified *exactly* up to the translation of the query range into the temporal range of the graph. Such constraints are more general than existential constraints, in that they can represent gaps. An interesting aspect of this work is that the history of each subgraph is represented as a string, and the temporal constraint is checked using substring search. While this method of expressing constraints can work over a set of time points, it is limited to ordered temporal constraints and does not support reoccurring edges. This work uses injective semantics and proposes optimization methods that exploit it.

There are various lines of research on querying temporal graphs that are complementary to our focus in this paper. For example, durable matchings [38, 54] count the number of snapshots in which a matching exists. Much attention has also been paid to tracing unbounded paths in temporal graphs, under various semantics, e.g., fastest, earliest arrival, latest departure, time-forward, time travel, or continuous [10, 20, 30, 64–66]. A focus on unbounded paths is complementary to our work on patterns without path variables, but with powerful temporal constraints. Extending our framework with path variables is an interesting direction for further research.

An important aspect of pattern matching in graphs is efficiently extracting the matchings. Early work started with the back-tracking algorithm by Ullmann [59], with later improvements [15, 60]. Pruning strategies for brute-force algorithms have been investigated as well [12, 17, 18]. Approaches suitable for large graphs typically build up the set of matchings in a relational table [36] by a series of natural joins over the edge relation; the aim is then to find an optimal join order. Until recently, the best-performing approaches were based on edge-growing pairwise join plans [35, 53, 57], but a new family of vertex-growing plans, known as worst-case optimal joins, have emerged [5, 46, 61], with better performance for cyclic patterns such as triangles. While we use the latter approach and implement our algorithms using relational operators, any method capable of finding matchings on a static graph can be combined with our timed automaton-based algorithms.

Another relevant direction is incremental graph pattern matching [3, 6, 13, 22, 27, 29, 31, 32], where the goal is to find and maintain patterns in an updating graph.

## 8 Conclusions and future work

In this paper, we proposed to use timed automata as a simple but powerful formalism for specifying temporal con-

straints in temporal graph pattern matching. We introduced algorithms that retrieve all temporal BGP matchings in a large graph, and presented results of an experimental evaluation, showing that this approach is practical, and identifying interesting performance trade-offs. Our code and data are available at <http://github.com/amirpouya/TABGP>.

An interesting open problem is how timed automata exactly compare to SQL in expressing temporal constraints (pinpointing the expressive power of SQL on ordered data is notoriously hard [39]). It is also interesting to investigate the decidability and complexity of the containment problem for temporal BGPs based on timed automata. Another natural direction for further research is to adapt our framework to a temporal graph setting where edges are active at durations (intervals), rather than at separate timepoints. Our hypothesis is that we can encode any set of non-overlapping intervals by the set of border-points. We conjecture that timed automata under such an encoding can express common constraints on intervals, such as Allen's relations [1].

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1007/s00778-023-00795-z>.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Allen, J.F.: Maintaining knowledge about temporal intervals. *Commun. ACM* **26**(11), 832–843 (1983)
- Alur, R., Dill, D.: A theory of timed automata. *Theoret. Comput. Sci.* **126**, 183–235 (1994)
- Ammar, K., McSherry, F., Salihoglu, S., Joglekar, M.: Distributed evaluation of subgraph queries using worst-case optimal low-memory dataflows. *PVLDB* **11**(6), 691–704 (2018)
- Angles, R., Arenas, M., Barceló, P., Hogan, A., Reutter, J.L., Vrgoc, D.: Foundations of modern query languages for graph databases. *ACM Comput. Surv.* **50**(5), 681–6840 (2017). <https://doi.org/10.1145/3104031>
- Arroyuelo, D., Hogan, A., Navarro, G., Reutter, J.L., Rojas-Ledesma, J., Soto, A.: Worst-case optimal graph joins in almost no space. In: *SIGMOD* (2021)
- Bindschaedler, L., Malicevic, J., Lepers, B., Goel, A., Zwaenepoel, W.: Tesseract: distributed, general graph pattern mining on evolving graphs. In: Barbalace, A., Bhatotia, P., Alvisi, L., Cadar, C. (eds.) *EuroSys '21: Sixteenth European Conference on Computer Systems*, Online Event, United Kingdom, April 26–28, 2021, pp. 458–473. ACM (2021). <https://doi.org/10.1145/3447786.3456253>
- Bonifati, A., Fletcher, G., Voigt, H., Yakovets, N.: *Querying Graphs. Synthesis Lectures on Data Management*. Morgan & Claypool (2018)
- Bouros, P., Mamoulis, N., et al.: Set containment join revisited. *Knowl. Inf. Syst.* **49**, 375–402 (2016)
- Bouyer, P., Fahrenberg, U., Larsen, K.G., Markey, N., Ouaknine, J., Worell, J.: Model checking real-time systems. In: Clarke, E., Henzinger, T., Veith, H., et al. (eds.) *Handbook of Model Checking*, pp. 1001–1046. Springer (2018)
- Byun, J., Woo, S., Kim, D.: Chronograph: enabling temporal graph traversals for efficient information diffusion analysis over time. *IEEE Trans. Knowl. Data Eng.* **32**(3), 424–437 (2020). <https://doi.org/10.1109/TKDE.2019.2891565>
- Carbone, P., Katsifodimos, A., Ewen, S., Markl, V., Haridi, S., Tzoumas, K.: Apache flink: Stream and batch processing in a single engine. *Bull. IEEE Comput. Soc. Tech. Committee Data Eng.* **36**(4) (2015)
- Carletti, V., Foggia, P., Saggese, A., Vento, M.: Challenging the time complexity of exact subgraph isomorphism for huge and dense graphs with VF3. *IEEE Trans. Pattern Anal. Mach. Intell.* **40**(4), 804–818 (2018). <https://doi.org/10.1109/TPAMI.2017.2696940>
- Chen, L., Wang, C.: Continuous subgraph pattern search over certain and uncertain graph streams. *IEEE Trans. Knowl. Data Eng.* **22**(8), 1093–1109 (2010). <https://doi.org/10.1109/TKDE.2010.67>
- Cheng, J., Yu, J.X., Ding, B., Philip, S.Y., Wang, H.: Fast graph pattern matching. In: 2008 IEEE 24th International Conference on Data Engineering, pp. 913–922. IEEE (2008)
- Čibej, U., Mihelič, J.: Improvements to Ullmann's algorithm for the subgraph isomorphism problem. *Int. J. Pattern Recognit. Artif. Intell.* **29**(07), 1550025 (2015)
- Conte, D., Foggia, P., Sansone, C., Vento, M.: Thirty years of graph matching in pattern recognition. *Int. J. Pattern Recognit. Artif. Intell.* **18**(03), 265–298 (2004)
- Cordella, L.P., Foggia, P., Sansone, C., Vento, M.: Performance evaluation of the vf graph matching algorithm. In: *Proceedings 10th International Conference on Image Analysis and Processing*, pp. 1172–1177. IEEE (1999)
- Cordella, L.P., Foggia, P., Sansone, C., Vento, M.: A (sub) graph isomorphism algorithm for matching large graphs. *IEEE Trans. Pattern Anal. Mach. Intell.* **26**(10), 1367–1372 (2004)
- Curley, J.: Engsoccerdata: English soccer data 1871–2106. *R Package Version 0.1* **5** (2016)
- Debrouvier, A., Parodi, E., Perazzo, M., Soliani, V., Vaisman, A.: A model and query language for temporal graph databases. *VLDB J.* **30**(5) (2021). <https://doi.org/10.1007/s00778-021-00675-4>
- Fan, W., Li, J., Ma, S., Tang, N., Wu, Y., Wu, Y.: Graph pattern matching: from intractable to polynomial time. *Proc. VLDB Endow.* **3**(1–2), 264–275 (2010)
- Fan, W., Wang, X., Wu, Y.: Incremental graph pattern matching. *ACM Trans. Database Syst.* **38**(3), 1–47 (2013)
- Ferrère, T., Maler, O., Nickovic, D., Pnueli, A.: From real-time logic to timed automata. *J. ACM* **66**(3), 191–1931 (2019)
- Grez, A., Riveros, C., Ugarte, M., Vansummeren, S.: A formal framework for complex event recognition. *TODS* **46**(4) (2021)
- Gupta, A., Mumick, I.S., Subrahmanian, V.S.: Maintaining views incrementally. *ACM SIGMOD Record* **22**(2), 157–166 (1993)
- Gurukar, S., Ranu, S., Ravindran, B.: COMMIT: a scalable approach to mining communication motifs from dynamic networks. In: Sellis, T.K., Davidson, S.B., Ives, Z.G. (eds.) *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, Melbourne, Victoria, Australia, May 31–June 4, 2015, pp. 475–489. ACM (2015). <https://doi.org/10.1145/2723372.2737791>

27. Han, W., Lee, J., Lee, J.: Turbo<sub>iso</sub>: towards ultrafast and robust subgraph isomorphism search in large graph databases. In: Ross, K.A., Srivastava, D., Papadias, D. (eds.) Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2013, New York, NY, USA, June 22–27, 2013, pp. 337–348. ACM (2013). <https://doi.org/10.1145/2463676.2465300>
28. Hopcroft, J., Ullman, J.: Introduction to Automata Theory, Languages, and Computation. Addison-Wesley (1979)
29. Imran, M., Gévy, G.E., Quiané-Ruiz, J.A., Markl, V.: Fast datalog evaluation for batch and stream graph processing. World Wide Web (2022)
30. Johnson, T., Kanza, Y., Lakshmanan, L.V.S., Shkapenyuk, V.: Nepal: a path query language for communication networks. In: Arora, A., Roy, S., Mehta, S. (eds.) Proceedings of the 1st ACM SIGMOD Workshop on Network Data Analytics, NDA@SIGMOD 2016, San Francisco, California, USA, July 1, 2016, pp. 6:1–6:8. ACM (2016). <https://doi.org/10.1145/2980523.2980530>
31. Kim, K., Seo, I., Han, W., Lee, J., Hong, S., Chafi, H., Shin, H., Jeong, G.: Turboflux: A fast continuous subgraph matching system for streaming graph data. In: Das, G., Jermaine, C.M., Bernstein, P.A. (eds.) Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10–15, 2018, pp. 411–426. ACM (2018). <https://doi.org/10.1145/3183713.3196917>
32. Ko, S., Lee, T., Hong, K., Lee, W., Seo, I., Seo, J., Han, W.: iturbo-graph: Scaling and automating incremental graph analytics. In: Li, G., Li, Z., Idreos, S., Srivastava, D. (eds.) SIGMOD '21: International Conference on Management of Data, Virtual Event, China, June 20–25, 2021, pp. 977–990. ACM (2021). <https://doi.org/10.1145/3448016.3457243>
33. Kondor, D., Pósfai, M., Csabai, I., Vattay, G.: Do the rich get richer? an empirical analysis of the bitcoin transaction network. CoRR abs/1308.3892 (2013). <http://arxiv.org/abs/1308.3892>
34. Kovanen, L., Karsai, M., Kaski, K., Kertész, J., Saramäki, J.: Temporal motifs in time-dependent networks. CoRR abs/1107.5646 (2011). <http://arxiv.org/abs/1107.5646>
35. Lai, L., Qin, L., Lin, X., Chang, L.: Scalable subgraph enumeration in mapreduce. Proc. VLDB Endow. **8**(10), 974–985 (2015). <https://doi.org/10.14778/2794367.2794368>. (<http://www.vldb.org/pvldb/vol8/p974-lai.pdf>)
36. Lai, L., Qing, Z., Yang, Z., Jin, X., Lai, Z., Wang, R., Hao, K., Lin, X., Qin, L., Zhang, W., et al.: Distributed subgraph matching on timely dataflow. Proc. VLDB Endow. **12**(10), 1099–1112 (2019)
37. Leskovec, J., Kleinberg, J.M., Faloutsos, C.: Graph evolution: densification and shrinking diameters. TKDD **1**(1), 2 (2007). <https://doi.org/10.1145/1217299.1217301>
38. Li, F., Zou, Z., Li, J.: Durable subgraph matching on temporal graphs. IEEE Trans. Knowl. Data Eng. (2022)
39. Libkin, L.: Expressive power of SQL. Theoret. Comput. Sci. **296**, 379–404 (2003)
40. McSherry, F., Murray, D.G., Isaacs, R., Isard, M.: Differential dataflow. In: CIDR (2013)
41. Moffitt, V.Z., Stoyanovich, J.: Temporal graph algebra. In: Proceedings of The 16th International Symposium on Database Programming Languages, DBPL 2017, Munich, Germany, September 1, 2017, pp. 10:1–10:12 (2017). <https://doi.org/10.1145/3122831.3122838>
42. Moffitt, V.Z., Stoyanovich, J.: Temporal graph algebra. In: Proceedings of The 16th International Symposium on Database Programming Languages, DBPL 2017, Munich, Germany, September 1, 2017, pp. 10:1–10:12 (2017). <https://doi.org/10.1145/3122831.3122838>
43. Murray, D.G., McSherry, F., Isaacs, R., Isard, M., Barham, P., Abadi, M.: Naiad: a timely dataflow system. In: Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles, pp. 439–455 (2013)
44. Neumann, T.: Efficiently compiling efficient query plans for modern hardware. Proc. VLDB Endow. **4**(9), 539–550 (2011). <https://doi.org/10.14778/2002938.2002940>. (<http://www.vldb.org/pvldb/vol4/p539-neumann.pdf>)
45. Neumann, T., Mühlbauer, T., Kemper, A.: Fast serializable multi-version concurrency control for main-memory database systems. In: Sellis, T.K., Davidson, S.B., Ives, Z.G. (eds.) Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31–June 4, 2015, pp. 677–689. ACM (2015). <https://doi.org/10.1145/2723372.2749436>
46. Ngo, H.Q., Ré, C., Rudra, A.: Skew strikes back: new developments in the theory of join algorithms. SIGMOD Rec. **42**(4), 5–16 (2013). <https://doi.org/10.1145/2590989.2590991>
47. Ojagh, S., Saeedi, S., Liang, S.H.: A person-to-person and person-to-place covid-19 contact tracing system based on ogc indoorgml. ISPRS Int. J. Geo Inf. **10**(1), 2 (2021)
48. Paranjape, A., Benson, A.R., Leskovec, J.: Motifs in temporal networks. In: de Rijke, M., Shokouhi, M., Tomkins, A., Zhang, M. (eds.) Proceedings of the Tenth ACM International Conference on Web Search and Data Mining, WSDM 2017, Cambridge, UK, February 6–10, 2017, pp. 601–610. ACM (2017). <https://doi.org/10.1145/3018661.3018731>
49. Raasveldt, M., Mühleisen, H.: Duckdb: an embeddable analytical database. In: Boncz, P.A., Manegold, S., Ailamaki, A., Deshpande, A., Kraska, T. (eds.) Proceedings of the 2019 International Conference on Management of Data, SIGMOD Conference 2019, Amsterdam, The Netherlands, June 30–July 5, 2019, pp. 1981–1984. ACM (2019). <https://doi.org/10.1145/3299869.3320212>
50. Reza, T., Ripeanu, M., Sanders, G., Pearce, R.: Approximate pattern matching in massive graphs with precision and recall guarantees. In: Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data, pp. 1115–1131 (2020)
51. Rost, C., Gomez, K., Täschner, M., Fritzsche, P., Schons, L., Christ, L., Adameit, T., Junghanns, M., Rahm, E.: Distributed temporal graph analytics with gradoop. VLDB J. **31**(2), 375–401 (2022)
52. Rust-Itertools: rust-iterools/iterools. <https://github.com/rust-iterools/iterools>
53. Selinger, P.G., Astrahan, M.M., Chamberlin, D.D., Lorie, R.A., Price, T.G.: Access path selection in a relational database management system. In: Readings in Artificial Intelligence and Databases, pp. 511–522. Elsevier (1989)
54. Semertzidis, K., Pitoura, E.: Durable graph pattern queries on historical graphs. In: 2016 IEEE 32nd International Conference on Data Engineering (ICDE), pp. 541–552. IEEE (2016)
55. Semertzidis, K., Pitoura, E.: A hybrid approach to temporal pattern matching. In: 2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pp. 384–388. IEEE (2020)
56. Semertzidis, K., Pitoura, E.: A hybrid approach to temporal pattern matching. In: Atzmüller, M., Coscia, M., Missaoui, R. (eds.) IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2020, The Hague, Netherlands, December 7–10, 2020, pp. 384–388. IEEE (2020). <https://doi.org/10.1109/ASONAM49781.2020.9381453>
57. Sun, Z., Wang, H., Wang, H., Shao, B., Li, J.: Efficient subgraph matching on billion node graphs. Proc. VLDB Endow. **5**(9), 788–799 (2012). <https://doi.org/10.14778/2311906.2311907>. ([http://vldb.org/pvldb/vol5/p788\\_zhaosun\\_vldb2012.pdf](http://vldb.org/pvldb/vol5/p788_zhaosun_vldb2012.pdf))
58. Ullman, J.: Principles of Database and Knowledge-Base Systems, vol. II. Computer Science Press (1989)
59. Ullmann, J.R.: An algorithm for subgraph isomorphism. J. ACM **23**(1), 31–42 (1976)

60. Ullmann, J.R.: Bit-vector algorithms for binary constraint satisfaction and subgraph isomorphism. *J. Exp. Algorithm.* **15**, 161–1664 (2011). <https://doi.org/10.1145/1671970.1921702>
61. Veldhuizen, T.L.: Leapfrog triejoin: a worst-case optimal join algorithm. In: *Proceedings 17th International Conference on Database Theory*, pp. 96–106 (2014)
62. Viswanath, B., Mislove, A., Cha, M., Gummadi, K.P.: On the evolution of user interaction in facebook. In: *Proceedings of the 2nd ACM SIGCOMM Workshop on Social Networks (WOSN'09)* (2009)
63. Wood, P.: Query languages for graph databases. *SIGMOD Rec.* **41**(1), 50–60 (2012)
64. Wu, H., Cheng, J., Huang, S., Ke, Y., Lu, Y., Xu, Y.: Path problems in temporal graphs. *Proc. VLDB Endow.* **7**(9), 721–732 (2014). <https://doi.org/10.14778/2732939.2732945>. (<http://www.vldb.org/pvldb/vol7/p721-wu.pdf>)
65. Wu, H., Cheng, J., Ke, Y., Huang, S., Huang, Y., Wu, H.: Efficient algorithms for temporal path computation. *IEEE Trans. Knowl. Data Eng.* **28**(11), 2927–2942 (2016). <https://doi.org/10.1109/TKDE.2016.2594065>
66. Wu, H., Huang, Y., Cheng, J., Li, J., Ke, Y.: Reachability and time-based path queries in temporal graphs. In: *32nd IEEE International Conference on Data Engineering, ICDE 2016, Helsinki, Finland, May 16–20, 2016*, pp. 145–156. IEEE Computer Society (2016). <https://doi.org/10.1109/ICDE.2016.7498236>
67. Xie, J., Yang, J.: A survey of join processing in data streams. In: *Data Streams*, pp. 209–236. Springer (2007)
68. Xu, Y., Huang, J., Liu, A., Li, Z., Yin, H., Zhao, L.: Time-constrained graph pattern matching in a large temporal graph. In: Chen, L., Jensen, C.S., Shahabi, C., Yang, X., Lian, X. (eds.) *Web and Big Data*, pp. 100–115. Springer, Cham (2017)
69. Yang, J., Widom, J.: Incremental computation and maintenance of temporal aggregates. *VLDB J.* **12**(3), 262–283 (2003). <https://doi.org/10.1007/s00778-003-0107-z>
70. Yoo, A.B., Jette, M.A., Grondona, M.: Slurm: Simple linux utility for resource management. In: *Workshop on job scheduling strategies for parallel processing*, pp. 44–60. Springer (2003)
71. Zaharia, M., Xin, R.S., Wendell, P., Das, T., Armbrust, M., Dave, A., Meng, X., Rosen, J., Venkataraman, S., Franklin, M.J., Ghodsi, A., Gonzalez, J., Shenker, S., Stoica, I.: Apache spark: a unified engine for big data processing. *Commun. ACM* **59**(11), 56–65 (2016). <https://doi.org/10.1145/2934664>
72. Zhao, Q., Tian, Y., He, Q., Oliver, N., Jin, R., Lee, W.C.: Communication motifs: A tool to characterize social communications. In: *Proceedings of the 19th ACM International Conference on Information and Knowledge Management, CIKM '10*, p. 1645–1648. Association for Computing Machinery, New York, NY, USA (2010). <https://doi.org/10.1145/1871437.1871694>
73. Zhu, K., Fletcher, G., Yakovets, N.: Leveraging temporal and topological selectivities in temporal-clique subgraph query processing. In: *ICDE* (2021)
74. Züfle, A., Renz, M., Emrich, T., Franzke, M.: Pattern search in temporal social networks. In: Böhlen, M.H., Pichler, R., May, N., Rahm, E., Wu, S., Hose, K. (eds.) *Proceedings of the 21st International Conference on Extending Database Technology, EDBT 2018, Vienna, Austria, March 26–29, 2018*, pp. 289–300. OpenProceedings.org (2018). <https://doi.org/10.5441/002/edbt.2018.26>

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.