

# Supplementary Material

## Part

### Table of Contents

<b>I. Matching Requirements.</b> . . . . .	<b>62</b>
<b>II. Time Complexity</b> . . . . .	<b>63</b>
II.1. Pipeline Components . . . . .	63
II.2. Operators' algorithmic implementation . . . . .	64
<b>III. Maximal Common Subset Problem</b> . . . . .	<b>66</b>
III.1. Use Case: Running $\text{Atom}_{A/T}^{\mathcal{L}, \tau}$ operators . . . . .	70
<b>IV. Data And Models</b> . . . . .	<b>73</b>

#### Supplement I Matching Requirements.

Depending on the binary  $\text{xtLTL}_f$  operator of choice, we might set a successful match if at least two events in a trace, one per operator, match anytime in the future ( $E_{\Theta}^i$ ), or if correlation conditions hold across activated and targeted events within the temporal window of interest (e.g., *Until*;  $A_{\Theta}^i$ ). Given some maps  $M_1$  and  $M_2$  associating events of interest from the same trace  $\sigma^i$  to a list  $L$  providing the activation and target conditions to be tested, we define such two distinct matching requirements as follows:

$$E_{\Theta}^i(M_1, M_2) = \left\{ M[h, k] \mid \begin{array}{l} \exists j \in \text{dom}(M_1) \cap \text{dom}(M_2), L_1 \in M_1(j), L_2 \in M_2(j). \\ (A(h) \in L_1 \vee \exists h'. M[h, h'] \in L_1), \\ (T(k) \in L_2 \vee \exists k'. M[k, k'] \in L_2), \Theta(\sigma_h^i, \sigma_k^i) \end{array} \right\} \quad (\text{S1})$$

$$A_{\Theta}^i(M_1, M_2) = \begin{cases} \emptyset & \exists j \in \text{dom}(M_1) \cap \text{dom}(M_2). E_{\Theta}^i([j \mapsto M_1(j)], [j \mapsto M_2(j)]) = \emptyset \\ E_{\Theta}^i(M_1, M_2) & \text{oth.} \end{cases} \quad (\text{S2})$$

where  $i$  indicates the trace id (e.g., of  $\sigma^i$ ) from which the events in  $L$  should be considered while reconstructing their associated payload while determining the binary predicate  $\Theta$ .

After expressing the different matching semantics as such, we skip testing correlations when (e.g.) only one of the two operands will contain activation (or target) events. We distinguish the case where the activation and target conditions of interest are missing because some basic operators did not return those from the scenario where the correlation condition was tested but failed. While in the former we return either an empty set or a set of just activations or targets, in the latter we want to return a placeholder **False** to remember that the correlation condition was falsified. We can then define a *testing functor* as follows:

$$\mathcal{T}_{\Theta}^{F,i}(M_1, M_2) = \begin{cases} \text{cod}(M_1) \cup \text{cod}(M_2) & \text{dom}(M_1) = \emptyset \vee \text{dom}(M_2) = \emptyset \\ F_{\Theta}^i(M_1, M_2) & F_{\Theta}^i(M_1, M_2) \neq \emptyset \\ \text{False} & \text{oth.} \end{cases} \quad (\text{S3})$$

where  $F$  subsumes one of the two matching semantics of choice, either  $E_{\Theta}^i$  or  $A_{\Theta}^i$ .

## Supplement II Time Complexity

### Supplement II.1 Pipeline Components

We now discuss the first two procedures from Algorithm 2.

**Lemma S1.** *BULKINSERTION can be computed in quasi-linear time with respect to the full log size (thus including the trace size, as well as the size of their associated payload).*

**Proof.** As most of the data structures are vectors (Lines 7 and 8) and the only two maps are a hashmap (Line 6) and a hashmap of ordered maps (Line 11), this phase has a quasi-linear time computational complexity with respect to the size of the log, thus showing that no additional overhead for pre-calculating the semantics associated to the declarative clauses is required in [60]. In fact, bulk inserting the data requires scanning all the traces and their associated events, as well as its possible associated payload, while the insertion in the  $\text{AttBulkMap}_k$  performs a sorted insertion for each value  $p(k)$  (Line 11); this pays a log-linear cost over the number of distinct values per key  $k$ .  $\square$

**Lemma S2.** *LOADINGANDINDEXING has linear time complexity with respect to the log size (trace and data payload inclusive).*

**Proof.** We show that the first outer for loop has a time complexity in  $O(|\mathcal{L}|\epsilon)$ . The initialization of the primary index of  $\text{ActivityTable}_{\mathcal{L}}$  requires a linear scan of every unique activity label of the log, with index assignments being performed in constant time (Line 18). The loading of the aforementioned table (minus assignments of  $\text{Prev}(\text{Next})$ ) requires a scan of every element of the log, thus providing a linear complexity with respect to the log size (inclusive of trace lengths, Line 23). The loading of the  $\text{CountingTable}_{\mathcal{L}}$  has complexity  $O(|\Sigma||\mathcal{L}|)$  (Line 20), which is always dominated by loading of the  $\text{ActivityTable}_{\mathcal{L}}$ , as it is guaranteed that  $|\Sigma| \leq |\mathcal{L}|\epsilon$ . In fact, in the worst-case scenario, every single event has a unique activity label, and therefore  $|\Sigma| = |\mathcal{L}|\epsilon = |\text{ActToEventBulkVector}|$ . In the best-case scenario, there is only one activity label ( $|\Sigma| = 1$ ) providing a linear complexity with log size.

In the best case scenario, no trace comes with a data payload, and therefore the entire outermost second loop is skipped as the set of keys is completely empty, and therefore no  $\text{AttributeTable}_{\mathcal{L}}^k$  will be initialised. Otherwise, loading an  $\text{AttributeTable}_{\mathcal{L}}^k$  associated to a key  $\kappa \in K$  requires a linear scan over each activity label  $a$  from  $\Sigma$ . While doing so, we access the  $\text{AttBulkMap}$  for every unique activity label and payload key, from which we also scan all the events from  $lst$  associated to a same value  $v$ . The initialisation of each  $\text{AttributeTable}_{\mathcal{L}}^k$  primary index is linear with respect to the number of unique activity labels (Line 35); the loading and secondary indexing of  $\text{AttributeTable}_{\mathcal{L}}^k$  is bounded by the number of events associating a value  $v$  to each key  $\kappa$  of interest. In the worst-case scenario, every event contains the same key labels. Overall, this provides a complexity dominated by the number of key-value associations in each event and therefore by  $|K|$ , as each event can contain an unlimited number of payload keys. This boils down to a worst-case time complexity of  $O(|\mathcal{L}|\epsilon|K|)$ .

We show that the last outer for loop also has a time complexity in  $O(|\mathcal{L}|\epsilon)$ . This includes the assignment of the  $\text{Prev}(\text{Next})$  pointers and the initialisation of the  $\text{ActivityTable}_{\mathcal{L}}$ 's secondary index, requiring only a linear scan of the log (trace inclusive). No other data structures are iterated, therefore we always perform computations in  $O(|\mathcal{L}|\epsilon)$  time.

Overall, in the best case scenario no event is associated to a payload, and therefore the time complexity is dominated by the first two loops, which provide a time complexity linear to the log size, trace size included:  $O(|\mathcal{L}|\epsilon)$ . Otherwise, we have to scan each value associated to each event which, in the worst case scenario, requires a full scan of  $|K|$  for each event. In this other scenario, the time complexity is heavily dominated by the data payload and therefore is  $O(|\mathcal{L}|\epsilon|K|)$ .  $\square$

Last, we discuss the time complexity for the query processing algorithms.

**Lemma S3.** *ATOMISATIONPIPELINE (Algorithm 3) has polynomial time complexity with respect to the model, key-set, and elementary intervals' maximum size.*

**Proof.** In the best-case scenario, all the clauses have no payload conditions, for which the activation condition will only contain the activation (or target) activity label. For this, we only pay the linear cost of scanning the model,  $O(|\mathcal{M}|)$ . In the worst-case scenario, each clause has a non-trivial payload condition for activation and target (when available): this implies always computing an intersection over unordered sets, which has the computational cost of scanning the smaller set and checking which elements are contained in the larger, where each operation takes  $O(1)$  time. We can assume that each  $p$  is described in disjunctive normal form and that its size is negligible if compared to the size of both the collected elementary intervals and the model. Given that usually  $ak(B)$  contains all of the atoms referring to the data predicates to  $B$  and therefore  $Atom_{\mu,ad}(B, \kappa) \subseteq ak(B)$  for any  $\kappa \in K$ , we have a computational complexity of  $O(m|K||\mathcal{M}|)$  under the assumption that each clause's data payload condition contains predicates for any payload key in  $O(|K|)$  and that each interval has a size comparable to the maximum size of  $m = \mu(a, \kappa)$  for  $a \in \Sigma$  and  $\kappa \in K$ .  $\square$

**Lemma S4.** *QUERYSCHEUDLER (Algorithm 5) has a linear time complexity with respect to the number of operators appearing in the compiled query plan.*

**Proof.** If we exclude the root node, the branching factor of our graph is at most 2, as each operator (represented as a node in the graph) might contain at most two sub-expressions; therefore, we can estimate that the number of edges of  $\mathcal{G}$  is linear with respect to the graph nodes. As the scheduler is represented through a vector of vectors and the distance computation requires visiting each edge associated with each non-leaf node, while filling the scheduler layer-wise just requires visiting each operator, the time complexity of this step is linear with respect to the operators appearing in the formula.  $\square$

#### Supplement II.2 Operators' algorithmic implementation

We now discuss the computational complexity of the  $xtLTL_f$  operators. Let us assume that  $\ell$  is the size of the collected events in  $L$  that need to be matched. Let us also remember that the intermediate results  $\rho$  and  $\rho'$  are represented as vectors of triplets sorted by trace and event id. Let us also assume all traces have a maximum length of  $\epsilon$ . In the worst case scenario, we have  $|\rho| = |\rho'| = |\mathcal{L}|\epsilon$ .

**Lemma S5.** *Future( $\rho$ ) has linear time complexity with respect to the operand size while  $Future^\tau(\rho)$  has quadratic time complexity with respect to the maximum trace length size.*

**Proof.** The computational complexity associated with the untimed Future operator boils down to performing an aggregation for each trace represented in the operand as well as its represented event, where its collected activation/target conditions are all associated with the first event of the trace. Therefore, this has a linear time complexity with respect to the size of each operand, i.e.  $|\rho| \in O(|\mathcal{L}|\epsilon)$ .

On the other hand, its timed counterpart ( $Future^\tau$ ) needs both to list all of the events in the operand and, for each of them, to associate to it also the activation and target conditions happening in the future, thus requiring a quadratic cost over the size of each trace. This boils down to  $O(|\mathcal{L}|\epsilon^2)$ .  $\square$

**Lemma S6.** *TIMEDIINTERSECTION has computational complexity in  $O(\ell^2|\mathcal{L}|\epsilon)$ .*

**Proof.** The time complexity of the TIMEDIINTERSECTION is in the worst-case scenario in  $O(\ell^2|\mathcal{L}|\epsilon)$  if we assume that each matched event needs to be tested and that  $\rho$  and  $\rho'$  have non-empty  $L$ s. If either  $\Theta$  is always true or the  $L$ s of the two operands are empty, the time complexity boils down to  $O(|\mathcal{L}|\epsilon)$ .  $\square$

**Corollary S1.** For  $\Theta = \text{True}$ ,  $\text{FASTUNTIMEDAND}(\rho, \rho')$  is faster than  $\text{SLOWUNTIMEDAND}(\rho, \rho')$ . For  $\ell \approx 1$ , the theoretical speed-up is proportional to the increase of both trace length and number of traces.

**Proof.** The time complexity of the  $\text{SLOWUNTIMEDAND}$  is in the worst case scenario in  $O(\ell^2 |\mathcal{L}| \epsilon^2)$ : as the former is a minor adaptation of  $\text{TIMEDIINTERSECTION}$  where, nevertheless, we have to pay a quadratic cost for scanning the whole traces as in timed future ( $\text{Future}^\tau$ ). This has also an additional cost of  $O(|\mathcal{L}| \log |\mathcal{L}| \epsilon)$  gained while grouping the operands' results by trace id.

On the other hand, its  $\text{FASTUNTIMEDAND}$  avoids a quadratic cost if  $\Theta = \text{True}$  by only collecting the activation and target conditions through a linear scan of both operands, thus boiling down its computational complexity to  $O(2|\mathcal{L}| \epsilon \ell)$ .

At this point, we can compute the speed-up of the allegedly slower version over the faster one by checking when the following condition holds:

$$\frac{\ell^2 |\mathcal{L}| \epsilon^2 + 2|\mathcal{L}| \log |\mathcal{L}| \epsilon}{2|\mathcal{L}| \epsilon \ell} = \frac{\ell \epsilon}{2} + \frac{\log |\mathcal{L}|}{\ell} > 1$$

We can close the goal after observing that the aforementioned condition always holds, and that is heavily dominated by the average trace length.  $\square$

**Corollary S2.** Computing  $\text{FASTUNTIMEDAND}(\rho, \rho')$  is more efficient than its equivalent  $\text{FASTUNTIMEDAND}(\text{Future}(\rho), \text{Future}(\rho'))$  by a theoretical constant speed up.

**Proof.** By expanding and composing the definitions of the former computational complexities, we get that the speed-up is almost constant for small  $\ell$  values:

$$\frac{2|\mathcal{L}| \epsilon + 2\mathcal{L} \epsilon \ell}{2\mathcal{L} \epsilon \ell} = \frac{1}{\ell} + 1 > 1$$

$\square$

**Corollary S3.** Computing  $\text{FASTUNTIMEDOR}(\rho, \rho')$  is more efficient than *Choice* as per its original *Declare semantics*, i.e.  $\text{FASTUNTIMEDOR}(\text{Future}(\rho), \text{Future}(\rho'))$  by a theoretical constant speed up.

**Proof.** As when no traces match we have no additional matching costs as they are just collected in the final results with a linear scan, the worst-case scenario boils down to performing trace matches as per the previous corollary.  $\square$

**Lemma S7.** For  $\Theta = \text{True}$ ,  $\text{UNTIMEDUNTIL}^2(\rho, \rho')$  is faster than  $\text{UNTIMEDUNTIL}^1(\rho, \rho')$  for greater log sizes.

**Proof.** In both implementations of the untimed Until operator, for each log trace, we perform a logarithmic scan for reaching the end of each trace within the operand. Trace by trace, the computation time decreases with the number of the operand's traces being visited. This can be expressed as  $\sum_{i=0}^{|\mathcal{L}|-1} \log(|\mathcal{L}| - i \epsilon) = \sum_{i=1}^{|\mathcal{L}|} \log(i \epsilon)$ . Given  $\sum_{i=0}^n \log i = \log(i!)$ , we can rewrite the computational cost of such operation as follows:

$$\sum_{i=1}^{|\mathcal{L}|} \log(i \epsilon) = \sum_{i=1}^{|\mathcal{L}|} (\log i + \log \epsilon) = \log(|\mathcal{L}|!) + |\mathcal{L}| \log \epsilon \leq |\mathcal{L}| (\log |\mathcal{L}| \epsilon)$$

In its worst-case scenario,  $\text{UNTIMEDUNTIL}^1$  always run the same scan over each trace's targeted events for each activated event in the second operand, while its Fast- counterpart performs constant access for each of the  $\epsilon$  trace events. So, while  $\text{UNTIMEDUNTIL}^2$  has an overall time complexity of  $\epsilon |\mathcal{L}|$ , the former has the upper bound of:

$$= \sum_{i=1}^{|\mathcal{L}|} (\log i + \sum_{j=1}^{\epsilon} \log j) \leq |\mathcal{L}| \log |\mathcal{L}| + |\mathcal{L}| \epsilon \log \epsilon = |\mathcal{L}| \log(|\mathcal{L}| \epsilon^{\epsilon})$$

Furthermore, for both implementations we pay at most a quadratic cost for scanning each targeted event that needs to hold until the activation is met, thus including the cost of inserting the matched conditions in the intermediate result set. This implies a cost of  $O(\epsilon^2 |\mathcal{L}|)$  despite  $\Theta = \mathbf{True}$ . As for the speed-up we can neglect the sub-operations having the same computational complexity<sup>22</sup>, we can rewrite the asymptotic speed-up as:

$$\frac{|\mathcal{L}| \log(|\mathcal{L}| \epsilon^{\epsilon})}{\epsilon |\mathcal{L}|} = \frac{\log |\mathcal{L}|}{\epsilon} + \log \epsilon > 1$$

which is always true only for adequately large traces,  $\epsilon \gg 0$ . We can see that if  $\epsilon$  is fixed, the speed-up grows proportionally with the log size. As the equation is always true for non-empty datasets, the fast implementation is always faster than the slower counterpart.  $\square$

**Corollary S4.** *For  $\Theta = \mathbf{True}$ , the second variant of the (timed) ANDGLOBALLY (Algorithm S2) is theoretically more performant than the first one (Algorithm S1) when the number of traces is exponentially upper bounded by the maximum trace length size.*

**Proof.** In the second variant of the operator, for each trace, we always perform a logarithmic scan for both operators for determining the end of the trace. As seen in the previous lemma, this corresponds to  $O(|\mathcal{L}| \log(|\mathcal{L}| \epsilon^{\epsilon}))$ . While scanning the trace events backwards we have that, in the worst-case scenario, each event on the left operand corresponds an event on the right one. This operation, as well as the gradual backward creation of the result, boils down to an additional time complexity of  $2\epsilon |\mathcal{L}|$ .

In the first algorithm for the operator (referred to as “variant”), we scan all of the events per single trace, this reducing to  $2|\mathcal{L}| \epsilon$ , but, for each event in the left operand, we have to scan all of the events in the right one until the end of the trace, thus adding up to a quadratic computational complexity with respect to the trace size:  $\epsilon |\mathcal{L}|$ .

The asymptotic and theoretical speed-up can be then computed as follows:

$$\begin{aligned} \frac{\epsilon^2 |\mathcal{L}|}{2|\mathcal{L}| \log(|\mathcal{L}| \epsilon^{\epsilon})} &= \frac{\epsilon^2}{2(\log |\mathcal{L}| + \epsilon \log \epsilon)} > 1 \Leftrightarrow \frac{\epsilon^2}{2} - \epsilon \log \epsilon > \log |\mathcal{L}| \\ &\Leftrightarrow |\mathcal{L}| < 2^{\epsilon^2/2 - \epsilon \log \epsilon} \leq \sqrt{2^{\epsilon^2}} \end{aligned}$$

We can therefore conclude that the first variant will be faster when traces are shorter, while the other one will be always faster for “reasonably” long traces. We can also observe that, when the log size is fixed, the speed-up grows most quadratically with respect to the maximum trace length.  $\square$

### Supplement III Maximal Common Subset Problem

In order to minimise union and intersection operations over a sequence of atoms, we had to solve multiple times the MAXIMAL COMMON SUBSET PROBLEM<sup>23</sup>. This problem, different from the maximal common subsequence for sequences (e.g., strings) [61], can be stated as follows:

<sup>22</sup>  $\frac{A+B}{A+C} > 1 \Leftrightarrow B > C \Leftrightarrow \frac{B}{C} > 1$  when  $A+B > 0$ ,  $A+C > 0$ , and  $A > 0$  are always true.

<sup>23</sup> This problem was coded as the `partition_sets` function in [https://github.com/datagram-db/knobab/blob/main/include/yauc/structures/set\\_operations.h](https://github.com/datagram-db/knobab/blob/main/include/yauc/structures/set_operations.h).

---

**Algorithm S1** xtLTL<sub>f</sub> pseudocode implementation for derived operators, First Variant
 

---

```

1: function TIMEDANDFUTUREΘ1(ρ, ρ')
2:   it ← Iterator(ρ), it' ← Iterator(ρ')
3:   while it ≠ ⊥ and it' ≠ ⊥ do
4:     ⟨t, e, L⟩ ← current(it), ⟨t', e'', λ⟩ ← current(it')
5:     if t = t' and e ≥ e'' then
6:       L'' ← ∅; hasMatch ← Θ = True
7:       it'_* ← it'
8:       while it'_* ≠ ⊥ do
9:         ⟨t', e', L'⟩ ← current(it'_*);
10:        if t' ≠ t then break;
11:        tmp ← TΘE,i(L, L')
12:        if tmp ≠ False then
13:          hasMatch ← true; L'' ← L'' ∪ tmp
14:        end if
15:        next(it'_*)
16:      end while
17:      if hasMatch then yield ⟨t, e, L''⟩;
18:      next(it); next(it');
19:    else if t < t' or (t = t' and e < e') then
20:      next(it)
21:    else
22:      next(it')
23:    end if
24:  end while

25: function TIMEDANDGLOBALLYΘ1(ρ, ρ')
26:   it ← Iterator(ρ), it' ← Iterator(ρ')
27:   while it ≠ ⊥ and it' ≠ ⊥ do
28:     ⟨t, e, L⟩ ← current(it), ⟨t', e'', λ⟩ ← current(it')
29:     if t = t' and e = e'' then
30:       L'' ← ∅; hasMatch ← Θ = True; count ← 0
31:       it'_* ← it'
32:       while it'_* ≠ ⊥ do
33:         ⟨t', e', L'⟩ ← current(it'_*)
34:         if t' ≠ t then break;
35:         tmp ← TΘE,i(L, L')
36:         if tmp ≠ False then
37:           hasMatch ← true; L'' ← L'' ∪ tmp; count ← count + 1
38:         end if
39:         next(it'_*)
40:       end while
41:       if hasMatch and count = |σt| - e + 1 then yield ⟨t, e, L''⟩;
42:       next(it); next(it');
43:     else if t < t' or (t = t' and e < e') then
44:       next(it)
45:     else
46:       next(it')
47:     end if
48:  end while

```

▷ Algorithm 7

▷ Algorithm 7

**Problem S1.** Given a set of sets  $\mathcal{S} = \{S_1, \dots, S_n\}$ , we want to compute the maximal common subsets  $\Delta = \{\delta_1, \dots, \delta_k\}$  so that each set  $S_i \in \mathcal{S}$  can be defined as the union of some pairwise disjoint sets  $\delta_{i_1}, \dots, \delta_{i_m}$  in  $\Delta$ . This decomposition has to guarantee that there exists no other decomposition containing non-overlapping supersets of the sets given in  $\Delta$  providing such decomposition.

Please observe that, if all of the sets in  $\mathcal{S}$  are pairwise disjoint, then  $\mathcal{S} = \Delta$ .

**Example S1.** Given  $\mathcal{S} = \{\{1, 2, 3\}, \{2, 3\}, \{4, 5\}, \{1, 2, 3, 4, 5\}, \{2, 3, 4, 5\}\}$ , we obtain the maximal common subsets  $\Delta = \{\{2, 3\}, \{1\}, \{4, 5\}\}$ . Given this, we might characterize the sets in  $\mathcal{S}$  from the constituents in  $\Delta$ , and therefore decomposition returned by the MAXIMALCOMMON-SUBSET function allows the characterization of  $\mathcal{S}$  as  $\{\delta_2 \cup \delta_1, \delta_1, \delta_3, \delta_2 \cup \delta_1 \cup \delta_3, \delta_1 \cup \delta_3\}$  via the returned decomposition  $\mathcal{I}$ .

Algorithm S3 sketches a solution to the maximal common subsets problem: by remembering the sets  $S_i$  in which each set item is contained (Line 4) and by subsequently remembering the items that are shared among the sets in  $\mathcal{S}$  (Line 7), the latter map identifies

---

**Algorithm S2**  $\text{xtLTL}_f$  pseudocode implementation for derived operators, Second Variant
 

---

```

1: function TIMEDANDFUTURE $^2_{\Theta}(\rho, \rho')$ 
2:   if  $\rho' = \emptyset$  then return  $\emptyset$  ▷ Yielding no records
3:    $it \leftarrow \text{Iterator}(\rho), it' \leftarrow \text{Iterator}(\rho')$ 
4:   while  $it \neq \uparrow$  do
5:      $toBeReversed \leftarrow \emptyset; \langle i_c, j_c, L \rangle \leftarrow \text{current}(it); \langle i', j', L' \rangle \leftarrow \text{current}(it')$ 
6:     if  $i_c > i'$  or  $(i = i' \text{ and } j_c > j')$  then
7:        $it' \leftarrow \text{LOWERBOUND}(\rho, it', \uparrow, \langle i_c, j_c, \top_{\Omega} \rangle);$  if  $it' = \uparrow$  break
8:     else if  $i < i'$  then
9:        $it \leftarrow \text{LOWERBOUND}(\rho, it, \uparrow, \langle i', 1, \top_{\Omega} \rangle);$ 
10:    else
11:      if  $it' = \uparrow$  break;  $tmp \leftarrow \emptyset$ 
12:       $bend \leftarrow \text{UPPERBOUND}(\rho', it', \uparrow, \langle i', |\sigma^{it'}| + 1, \top_{\Omega} \rangle); aend \leftarrow \text{UPPERBOUND}(\rho, it, \uparrow, \langle i_c, |\sigma^{it_c}| + 1, \top_{\Omega} \rangle)$ 
13:      if  $\text{current}(it' - 1).i \neq i$  then
14:         $it \leftarrow aend;$ 
15:        while  $(it \neq \uparrow) \text{ and } \text{current}(it).i = i$  do  $it++$ 
16:        continue
17:      end if
18:       $it \leftarrow --bend; it' \leftarrow --aend; \langle i, j, L \rangle \leftarrow \text{current}(it); \langle i', j', L' \rangle \leftarrow \text{current}(it')$ 
19:      for  $J \leftarrow j'$  downto 1 by 1 do
20:        if  $bend \geq it'$  and  $J = j'$  then
21:           $tmp \leftarrow tmp \cup L'; it'--; \langle i', j', L' \rangle \leftarrow \text{current}(it')$ 
22:        end if
23:        while  $i = i'$  and  $j > j'$  do
24:           $aend--; \langle i, j, L \rangle \leftarrow \text{current}(it).$ 
25:        end while
26:        if  $i \neq i'$  then break
27:        if  $j < j_c$  or  $j < J$  then continue
28:         $tmp \leftarrow \mathcal{T}_{\Theta}^{E,i}(L, L')$  ▷ Algorithm 7
29:        if  $tmp \neq \text{False}$  then  $toBeReversed.add(\langle i, j, tmp \rangle)$ 
30:      end for
31:      forall  $\langle i, j, L'' \rangle \in \text{REVERSED}(toBeReversed)$  do yield  $\langle i, j, L'' \rangle$ 
32:       $it \leftarrow aend;$  if  $it' = \uparrow$  then break;
33:    end if
34:  end while

35: function TIMEDANDGLOBALLY $^2_{\Theta}(\rho, \rho')$ 
36:   if  $\rho' = \emptyset$  then return  $\emptyset$  ▷ Yielding no records
37:    $it \leftarrow \text{Iterator}(\rho), it' \leftarrow \text{Iterator}(\rho')$ 
38:   while  $it \neq \uparrow$  do
39:      $toBeReversed \leftarrow \emptyset; \langle i, j, L \rangle \leftarrow \text{current}(it); \langle i', j', L' \rangle \leftarrow \text{current}(it')$ 
40:     if  $i > i'$  or  $(i = i' \text{ and } j > j')$  then
41:        $it++;$  if  $it' = \uparrow$  break
42:     else
43:       if  $it' = \uparrow$  break;  $tmp \leftarrow \emptyset$ 
44:        $bend \leftarrow \text{UPPERBOUND}(\rho', it', \uparrow, \langle i', |\sigma^{it'}| + 1, \top_{\Omega} \rangle); aend \leftarrow \text{UPPERBOUND}(\rho, it, \uparrow, \langle i, |\sigma^i| + 1, \top_{\Omega} \rangle)$ 
45:        $it \leftarrow --bend; it' \leftarrow --aend; \langle i, j, L \rangle \leftarrow \text{current}(it); \langle i', j', L' \rangle \leftarrow \text{current}(it')$ 
46:        $J \leftarrow j'$ 
47:       for  $J \leftarrow j'$  downto 1 by 1 do
48:         if  $bend \geq it'$  and  $\text{DISTANCE}(it', bend) = |\sigma^{it'}| - J + 1$  then
49:            $tmp \leftarrow tmp \cup L'; it'--; \langle i', j', L' \rangle \leftarrow \text{current}(it')$ 
50:         else break
51:       end if
52:       while  $i = i'$  and  $j > j'$  do
53:          $it--; \langle i, j, L \rangle \leftarrow \text{current}(it).$ 
54:       end while
55:       if  $i \neq i'$  then break
56:       if  $j < J$  then continue
57:        $tmp \leftarrow \mathcal{T}_{\Theta}^{E,i}(L, L')$  ▷ Algorithm 7
58:       if  $tmp \neq \text{False}$  then  $toBeReversed.add(\langle i, j, tmp \rangle)$ 
59:     end for
60:     forall  $\langle i, j, L'' \rangle \in \text{REVERSED}(toBeReversed)$  do yield  $\langle i, j, L'' \rangle$ 
61:      $it \leftarrow aend;$  if  $it' = \uparrow$  then break
62:   end if
63: end while

```

---

the desired maximal common subsets. Therefore, each set  $S_i$  can be defined through the union of the maximal common subsets  $\delta_k \in \Delta$  (Line 10).

The primary aim of this algorithm in KnoBAB is to minimize the computations of Or or And operators where the computation of further intermediate subsets is appropriate. As we might indeed see from the former example, the reconstruction of the original sets in  $\mathcal{S}$  from their constituents in  $\Delta$  implies evaluating the same unions multiple times, e.g.  $\delta_2 \cup \delta_1$  and  $\delta_1 \cup \delta_3$ . This overhead can be limited as follows: first, we might sort the



---

**Algorithm S3** Maximal Common Subsets
 

---

```

1: function MAXIMALCOMMONSUBSET( $\mathcal{S}$ )
2:   setsToSharedValues  $\leftarrow \{\}$ ; itemInSets  $\leftarrow \{\}$ ;  $\mathcal{I} \leftarrow \{\}$ ;  $k \leftarrow 1$ ;  $\Delta \leftarrow \emptyset$ 
3:   for all  $S_i \in \mathcal{S}$  and item  $\in S_i$  do
4:     itemInSets[item].put( $i$ )
5:   end for
6:   for all  $\langle \text{item}, \{i, j, \dots\} \rangle \in \text{itemInSets}$  do
7:     setsToSharedValues[ $\{i, j, \dots\}$ ].put(item)  $\triangleright \text{item} \in S_i \cap S_j \cap \dots$ 
8:   end for
9:   for all  $\langle \{i, j, \dots\}, \{\text{item}_\alpha, \text{item}_\beta, \dots\} \rangle \in \text{setsToSharedValues}$  and  $\iota \in \{i, j, \dots\}$  do
10:     $\mathcal{I}[\iota].\text{put}(k)$   $\triangleright \delta_k \subseteq S_\iota, S_\iota \in \mathcal{S}$ 
11:     $\Delta.\text{put}(\{\text{item}_\alpha, \text{item}_\beta, \dots\})$   $\triangleright \delta_k = \{\text{item}_\alpha, \text{item}_\beta, \dots\}$ 
12:     $k \leftarrow k + 1$ 
13:   end for
14:   return  $\langle \Delta, \mathcal{I} \rangle$ 

15: function MCSREFINEMENT( $\Delta, \mathcal{I}$ )
16:    $\Delta' \leftarrow \emptyset$ ;  $k \leftarrow |\Delta| + 1$ ; elemMap  $\leftarrow \{\}$ 
17:   SORT( $\mathcal{I}, (\langle s, I \rangle, \langle s', I' \rangle) \rightarrow I \subseteq I'$ )
18:   for all  $\langle s, I \rangle_i \in \mathcal{I}$  from  $i = 1$  to  $|\mathcal{I}|$  s.t.  $|I| \neq 1$  do
19:     hasElem  $\leftarrow \text{False}$ ;
20:     for all  $\langle s', I' \rangle_j \in \mathcal{I}$  from  $j = |\mathcal{I}| - 1$  to  $i + 1$  s.t.  $|I'| \neq 1$  do
21:       if  $|I| = |I'|$  then break
22:       else if  $I \subseteq I'$  then
23:         if not hasElem then
24:           hasElem  $\leftarrow \text{True}$ 
25:            $\Delta'.\text{put}(I)$ 
26:         end if
27:          $I' \leftarrow I' \setminus I \cup \{k\}$ 
28:       end if
29:     end for
30:     if hasElem then
31:       elemMap[ $i$ ]  $\leftarrow k$ 
32:        $k \leftarrow k + 1$ 
33:     end if
34:   end for
35:   return  $\langle \Delta, \Delta', \mathcal{I} \rangle$ 
36:   for all  $\langle i, k \rangle \in \text{elemMap}$  do
37:      $\mathcal{I}[i] \leftarrow \{k\}$ 
38:   end for

```

---

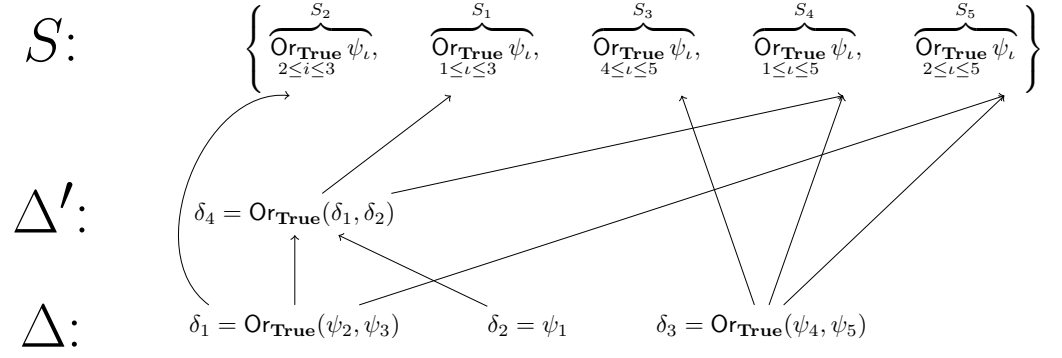
decomposed itemsets by their inclusion relationship (Line 17): this is possible as the set of all possible subsets is a partially ordered set (poset) where the partial order is the subset-equal relationship [40]. Last, we iterate over all the possible decomposed itemsets  $I$  which do not appear as a singleton, as their replacement is trivial (Line 18): if the decomposition  $I'$  for a set  $S_j$  contains the decomposition  $I$  for a set  $S_i$ , we can elect  $I$  as a new refined subset  $\delta_k$  (Line 25) which is now a component for  $I'$  (Line 27).

**Example S2.** Let us continue the former example: the refinement step computes a new refinement  $\delta_4 \in \Delta'$  which is defined as  $\delta_1 \cup \delta_2$ ; the characterization of  $\mathcal{S}$  has now become  $\{\delta_4, \delta_1, \delta_3, \delta_4 \cup \delta_3, \delta_1 \cup \delta_3\}$ .

**Example S3.** Let us now suppose to have an array of sets containing  $\mathbf{x}\mathbf{t}\mathbf{LTL}_f$  formulae  $\mathcal{S} = \{\{\psi_1, \psi_2, \psi_3\}, \{\psi_2, \psi_3\}, \{\psi_4, \psi_5\}, \{\psi_1, \psi_2, \psi_3, \psi_4, \psi_5\}, \{\psi_2, \psi_3, \psi_4, \psi_5\}\}$ , where we want to express each set in  $\mathcal{S}$  as a timed disjunction  $\text{Or}_{\text{True}}^{\tau}$  formulae. While doing so, we want also to minimise as possible the computation of unions that are shared among different sets. Figure S1 shows the desired result, where  $\Delta$  and  $\Delta'$  contain the intermediate untimed unions or just one single  $\mathbf{x}\mathbf{t}\mathbf{LTL}_f$  formula, while  $\mathcal{S}$  will contain the nodes that will perform the final union operations. FINITARYSETOPERATIONS( $\mathcal{S}, \text{Or}_{\text{True}}^{\tau}$ ) returns the desired  $\mathbf{x}\mathbf{t}\mathbf{LTL}_f$  DAG in Figure S1, rooted in each distinct formula associated to a set in  $\mathcal{S}$ .

*Time Complexity.* For MAXIMALCOMMONSUBSET, if we assume that each subset is at most of size  $m$ , then the time complexity of the first loop is in  $O(|\mathcal{S}|m)$ ; the second loop iterates over the whole possible items contained in each set in  $\mathcal{S}$ , and therefore has a time complexity of  $O(|\bigcup \mathcal{S}|)$ ; the last iteration performs a number of iterations which is





**Figure S1.** Decomposing untimed unions of  $\text{xTLTL}_f$  formulae into maximal shared sub-expressions via an abstract syntax DAG, thus attempting at computing each untimed union shared among different expressions at most once.

---

**Algorithm S4** Finitary set operations for commutative and associative set operations  $\oplus$

---

```

1: function FINITARYSETOPERATIONS( $\mathcal{S}, \oplus$ ) ▷  $\oplus \in \{\cup, \cap\}$ 
2:    $\Delta, \mathcal{I} \leftarrow \text{MAXIMALCOMMONSUBSETS}(\mathcal{S})$ 
3:    $\Delta, \Delta', \mathcal{I} \leftarrow \text{MCSREFINEMENT}(\Delta, \mathcal{I})$ 
4:   for all  $s' \in \Delta'$  (parallel) do
5:      $\delta_{s'} \leftarrow \oplus_{\delta_i \in \Delta} \delta_i$ 
6:   end for
7:   for all  $\langle i, D \rangle \in \mathcal{I}$  do ▷ Decomposition  $D$  of  $S_i$  through intermediate sets  $\delta_s$  with  $s \in D$ 
8:     yield  $\langle i, \oplus_{s \in D} \delta_s \rangle$  ▷ Result associated to  $S_i \in \mathcal{S}$ 
9:   end for

10: function FINITARYSETOPERATIONS( $\mathcal{S}, \text{Results}, \oplus$ ) ▷  $\oplus \in \{\cup, \cap\}$ 
11:    $\Delta, \mathcal{I} \leftarrow \text{MAXIMALCOMMONSUBSETS}(\mathcal{S})$ 
12:    $\Delta, \Delta', \mathcal{I} \leftarrow \text{MCSREFINEMENT}(\Delta, \mathcal{I})$ 
13:   for all  $s \in \Delta$  (parallel) do
14:      $\delta_s \leftarrow \oplus_{i \in \mathcal{S}} \text{Results}[i]$ 
15:   end for
16:   for all  $s' \in \Delta'$  (parallel) do
17:      $\delta_{s'} \leftarrow \oplus_{i \in s'} \delta_i$ 
18:   end for
19:   for all  $\langle i, D \rangle \in \mathcal{I}$  do ▷ Decomposition  $D$  of  $S_i$  through intermediate sets  $\delta_s \in \Delta \cup \Delta'$ 
20:     yield  $\langle i, \oplus_{s \in D} \delta_s \rangle$ 
21:   end for

```

---

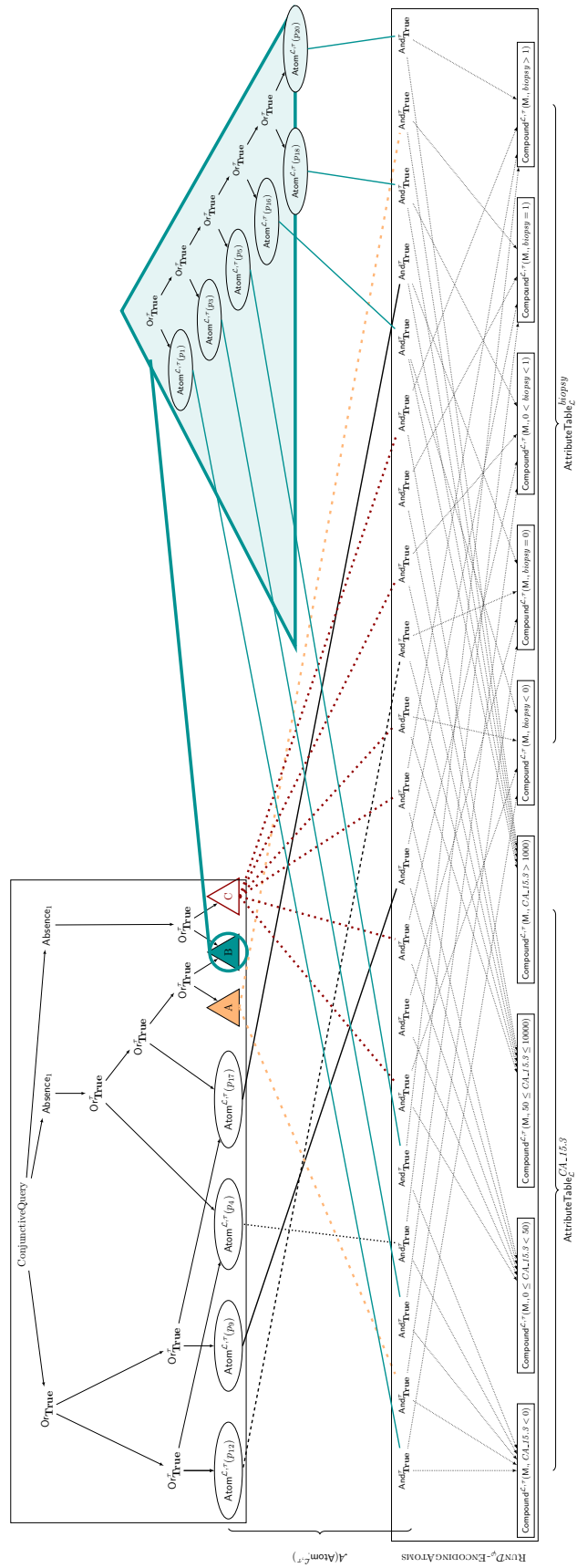
comparable to the size of  $\Delta$  which, in the worst-case scenario, is composed of singletons of size  $k = 1$  comprising elements from  $\cup \mathcal{S}$ , and therefore has a time complexity in  $O(|\cup \mathcal{S}|)$ . We also assume that all the insertion operations can be performed in  $O(1)$  time, as all the data structures that might be exploited are either hash maps or vectors. Overall, this makes a worst case scenario time complexity of  $O(|\mathcal{S}|m + |\cup \mathcal{S}|)$ .

In MCSREFINEMENT, as checking the subset relationship among itemsets has at most a cost proportional to the size  $k$  of each itemset, this makes such sorting  $O(|\cup \mathcal{S}| \log(|\cup \mathcal{S}|))$  when  $k = 1$  and  $O(|\cup \mathcal{S}| \log(c))$  when  $k = |\cup \mathcal{S}|/c$  for  $c$  constant; so, even in this scenario, we have the worst case scenario computational complexity for  $k = 1$ . Last, in the worst case scenario, the double for loop has a computational complexity comparable to  $O(|\mathcal{S}|^2)$ , as we perform an iteration over all the possible characterizations of the sets in  $\mathcal{S}$  while trying to aggregate the remaining shared union operations further.

The overall time complexity of both phases becomes  $O(|\mathcal{S}|(m + |\mathcal{S}|) + |\cup \mathcal{S}| \log(|\cup \mathcal{S}|))$ .

#### Supplement III.1 Use Case: Running $\text{Atom}_{A/T}^{\mathcal{L}, \tau}$ operators

This subsection describes how efficiently computing the results associated to the  $\text{Atom}_{A/T}^{\mathcal{L}, \tau}$  operator also requires the computation of the MAXIMALCOMMONSUBSET problem. This is achieved not only by performing all the Compound $_{A/T}^{\mathcal{L}, \tau}$  queries to be run on the same AttributeTable $_{\mathcal{L}}^k$  in one single scan of the former but by also running the final result associated to the  $\text{Atom}_{A/T}^{\mathcal{L}, \tau}$  by computing the shared intersections among different atoms' definition at most once. This detailed pseudocode is shown in Algorithm S5.



**Figure S2.** Detailed view of the query plan represented in Figure 5 for Example 15.

---

**Algorithm S5** Query Plan Initialisation and Execution for  $\mathcal{D}_\varphi$ -encoding  $p_i$  atoms.

---

```

1: global keyToLabelToSortedIntervals  $\leftarrow \{\}; S_\Sigma \leftarrow \{\}; Results \leftarrow \{\}; Cache \leftarrow \{\}$ 

2: procedure EXECUTE RANGEQUERIES( $\mathcal{L}$ )
3:   for all  $\langle \kappa, map \rangle \in \text{keyToLabelToSortedIntervals}$  do
4:     if table  $\text{AttributeTable}_{\mathcal{L}}^\kappa$  not exists then continue
5:     for all  $\langle a, L \rangle \in map$  (parallel) do
6:       if  $\exists \langle beg, end \rangle . \langle beg, end \rangle \in \text{AttributeTable}_{\mathcal{L}}^\kappa.\text{primary\_index}[\beta(a)]$  then
7:         SORT( $L$ )
8:         for all  $\langle low_\kappa \leq \kappa \leq up_\kappa, idx \rangle \in L$  do ▷ Computing Compound $_{\mathcal{L}, \tau}^{\mathcal{L}, \tau}(a, \kappa, [low_\kappa, up_\kappa])$ 
9:            $beg \leftarrow \text{LOWERBOUND}(\text{AttributeTable}_{\mathcal{L}}^\kappa, beg, low_\kappa, end)$ 
10:           $mid \leftarrow \text{UPPERBOUND}(\text{AttributeTable}_{\mathcal{L}}^\kappa, beg, up_\kappa, end)$ 
11:          while  $beg < mid$  and  $beg < end$  do
12:             $i, j \leftarrow \text{ActivityTable}_{\mathcal{L}}[\text{AttributeTable}_{\mathcal{L}}^\kappa[beg](\text{Offset})](\text{Trace}, \text{Event})$ 
13:             $Results[idx].\text{put}(\langle i, j \rangle)$ 
14:             $beg++$ 
15:          end while
16:           $beg \leftarrow mid$ ; SORT( $Results[idx]$ )
17:          if  $beg \geq end$  then break
18:        end for
19:      end if
20:    end for
21:  end for

22: procedure RUN  $\mathcal{D}_\varphi$ -ENCODING ATOMS( $\mathcal{L}$ )
23:   EXECUTE RANGEQUERIES( $\mathcal{L}$ )
24:   for all  $\langle p_i, Result \rangle \in \text{FINITARY SET OPERATIONS}(S_\Sigma, Results, \cap)$  do ▷ Algorithm S4
25:      $Cache[p_i] \leftarrow Result$ 
26:   end for
27:    $S_\Sigma \leftarrow \{\}; Results \leftarrow \{\}; \text{keyToLabelToSortedIntervals} \leftarrow \{\}$ 

28: function  $\mathcal{A}(\text{ATOM}_{\mathcal{L}, \tau}^{\mathcal{L}, \tau})(\psi)$  ▷ require  $\psi.\text{atom} = \{p_i\}$  with  $p_i$  from the  $\mathcal{D}_\varphi$ -encoding pipeline
29:   if  $\psi.\text{isActivation}$  then ▷  $\text{Atom}_A^{\mathcal{L}, \tau}(p_i)$ 
30:     return  $\{ \langle i, j, \{A(j)\} \rangle \mid \langle i, j \rangle \in Cache[p_i] \}$ 
31:   else if  $\psi.\text{isTarget}$  then ▷  $\text{Atom}_T^{\mathcal{L}, \tau}(p_i)$ 
32:     return  $\{ \langle i, j, \{T(j)\} \rangle \mid \langle i, j \rangle \in Cache[p_i] \}$ 
33:   else ▷  $\text{Atom}^{\mathcal{L}, \tau}(p_i)$ 
34:     return  $\{ \langle i, j, \emptyset \rangle \mid \langle i, j \rangle \in Cache[p_i] \}$ 
35:   end if

```

---

First, we evaluate each compound condition over the  $\text{AttributeTable}_{\mathcal{L}}^\kappa$ , as each atom is a conjunction of interval queries expressed as compound conditions. As the query compiler grouped such intervals by  $\text{AttributeTable}_{\mathcal{L}}^\kappa$  and associated activity labels by calling the `RETRIEVEINTERVALS` method in the `keyToLabelToSortedIntervals` map, we now need to iterate over such a map and access the  $\text{AttributeTable}_{\mathcal{L}}^\kappa$  for each key  $\kappa$  appearing in it (Line 4). Each  $\text{AttributeTable}_{\mathcal{L}}^\kappa$  is then accessed by activity label  $a$  through its primary index (Line 6): if the primary index points to a specific region in the table delimited by a beginning and end table offset, we can access it and query it via the sorted compound condition. For each of those (representing elementary intervals), we get the offset  $beg$  (and  $end$ ) to the first element within the block equal or greater to the lower bound  $low_\kappa$  (or strictly greater to the upper bound  $up_\kappa$ ). Such bounds are returned in Lines 9 and 10. For each record existing in such a range, we reconstruct the trace and event id containing the desired value (Line 12) and then store it into a preliminary result<sup>24</sup> (Line 13). As both the elementary intervals and the table are sorted in ascending order by value and given that each elementary interval is non-overlapping by construction, we can query the next interval in  $L$  from the remaining portion of the table, thus steadily reducing the table visiting cost by further reducing the data over which the scan is performed. As these results are not sorted by increasing trace and event id rather than by increasing value in  $V$ , we have to sort them accordingly (Line 16) to be compatible with the intermediate results  $\rho$  requirements. Please observe that despite each of these intervals could appear in the definition of multiple atoms, this construction guarantees that each of these queries is performed at most once per occurring interval.

---

<sup>24</sup> Please observe: **not** the intermediate representation  $\rho$ !

Second, we can compute the events associated with each  $p_i$  atom returned from the  $\mathcal{D}_\varphi$ -encoding pipeline by intersecting the result associated with each of the intervals. The result of such intersections is then stored in an intermediate cache (Line 25). We can perform these intersections efficiently by guaranteeing the computation of each sub-intersection shared across different atoms at most once through the calculation of maximal common subsets occurring in intersections' sub-expressions.

Third, while evaluating each atom from our scheduled operator in  $\psi$  (Line 7 from Algorithm 6), we can retrieve the result from the cache and, dependingly on whether this expression was associated to an activation (Line 29 from Algorithm S5), a target condition (Line 31), or neither of these (Line 34), we can now effectively represent the result as an intermediate result with the appropriate label markers in  $L$ .

**Example S4.** The lowermost boxes from Figure S2 represent the compound conditions extracted from the atoms appearing on the leaf nodes of the query plan, as well as their associated result stored in "Results" after running the EXECUTERANGEQUERIES. Their connection with the corresponding atoms could be tracked by navigating the edges backwards. The result of the intersection between all of the possible compound conditions as per atom definitions is then computed in the **for** loop in  $\text{RUN}\mathcal{D}_\varphi\text{-ENCODINGATOMS}$ . At this point, the execution of the leaf nodes on the query plan boils down to accessing the upper nodes in the box above and converting such results according to the activation or target condition associated with them. This mechanism guarantees that the results and the tables are accessed at most once, independently from the activation/target mark associated with the atom nodes.

#### Supplement IV Data And Models

We wanted to test our claim that exploiting data access minimisation techniques ~~to~~ can quicken conformance checking as well as model mining time. In addition, we want to prove that even in the worst-case scenario, we still achieve orders of magnitude of gains versus state-of-the-art solutions. Therefore, we provide two models, Table S2 and Table S1a (also available at <sup>17</sup>). The former was adapted from Burattin et al. [6], where we exploited the constraints they provided, allowing for a direct comparison. These tables are displayed under the assumption from [6], where the trace payloads are injected into every event. In our actual implementation, we inject a `__trace_payload` event at the beginning every trace. Therefore,  $M_{11}$  from Table S1a has a KnoBAB representation as Table S1b, where the resulting model size is five (seven) for Burattin (KnoBAB).

By performing conjunctions among these, we were able to provide 11 models, the largest consisting of five (seven, due to trace payload injection) clauses. The queries constructing these larger models are similar, for example, even the model  $q_1 \wedge q_2$  hold the clause `Response(A_SUBMITTED,true,A_ACCEPTED,true)` twice. These clauses need not be duplicated in the query plan, and therefore these models should better demonstrate the gains from data access minimisation. Conversely, Table S1a denotes an entirely novel model, that avoids similar queries. With each sub-model, entirely new events are considered, and even within those payload conditions vary. As such, our pipeline can gain much less from DAG optimization techniques. This model, therefore, provides a more appropriate case for analysing these proposed enhancements.

**Table S1.** Worst-case scenario (SCENARIO 1) model representation for the BPIC\_2012 dataset.

(a) Clause constituents of each model, where each sub-model is guaranteed to be a subset of the larger.

Model	Clauses
$M_1 =$	$\left\{ \begin{array}{l} q_1 := \text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true}) \\ q_2 := \text{Response}(\text{A\_SUBMITTED}, \text{AMOUNT\_REQ} \geq 10^3, \text{A\_ACCEPTED}, \text{true}) \\ q_3 := \text{Response}(\text{A\_SUBMITTED}, \text{AMOUNT\_REQ} < 10^3, \text{A\_ACCEPTED}, \text{true}) \\ q_4 := q_1 \text{ where } \text{A\_SUBMITTED.org:resource} = \text{A\_ACCEPTED.org:resource} \\ q_5 := q_1 \text{ where } \text{A\_SUBMITTED.org:resource} \neq \text{A\_ACCEPTED.org:resource} \end{array} \right.$
$M_2 = M_1 +$	$\left\{ \begin{array}{l} q_6 := \text{Response}(\text{W\_Completeren aanvraag}, \text{true}, \text{W\_Valideren aanvraag}, \text{true}) \\ q_7 := \text{Response}(\text{W\_Completeren aanvraag}, \text{true}, \text{O\_CANCELLED}, \text{true}) \\ q_8 := q_6 \text{ where } \text{W\_Valideren aanvraag.org:resource} \neq \text{W\_Valideren aanvraag.org:resource} \\ q_9 := \text{Response}(\text{W\_Valideren aanvraag}, \text{AMOUNT\_REQ} = 5 \cdot 10^3, \text{O\_CANCELLED}, \text{true}) \\ q_{10} := q_9 \text{ where } \text{W\_Valideren aanvraag.org:resource} = \text{O\_CANCELLED.org:resource} \end{array} \right.$
$M_3 = M_2 +$	$\left\{ \begin{array}{l} q_{11} := \text{Response}(\text{O\_SELECTED}, \text{true}, \text{O\_CANCELLED}, \text{true}) \\ q_{12} := q_{11} \text{ where } \text{O\_SELECTED.org:resource} = \text{O\_CANCELLED.org:resource} \\ q_{13} := \text{Response}(\text{O\_SELECTED}, \text{AMOUNT\_REQ} < 8 \cdot 10^3, \text{O\_CANCELLED}, \text{true}) \\ q_{14} := q_{13} \text{ where } \text{O\_SELECTED.org:resource} = \text{O\_CANCELLED.org:resource} \\ q_{15} := \text{Response}(\text{O\_SELECTED}, \text{AMOUNT\_REQ} > 10^3, \text{O\_CANCELLED}, \text{true}) \\ \text{where } \text{O\_SELECTED.org:resource} \neq \text{O\_CANCELLED.org:resource} \end{array} \right.$
$M_4 = M_3 +$	$\left\{ \begin{array}{l} q_{16} := \text{Response}(\text{A\_PARTLYSUBMITTED}, \text{true}, \text{A\_DECLINED}, \text{true}) \\ q_{17} := q_{16} \text{ where } \text{A\_PARTLYSUBMITTED.org:resource} = \text{A\_DECLINED.org:resource} \\ q_{18} := \text{Response}(\text{A\_PARTLYSUBMITTED}, \text{AMOUNT\_REQ} > 2 \cdot 10^4, \text{A\_DECLINED}, \text{true}) \\ q_{19} := \text{Response}(\text{A\_PARTLYSUBMITTED}, \text{AMOUNT\_REQ} > 2 \cdot 10^4, \text{A\_CANCELLED}, \text{true}) \\ q_{20} := q_{18} \text{ where } \text{A\_PARTLYSUBMITTED.org:resource} = \text{A\_DECLINED.org:resource} \end{array} \right.$

(b) Distinguishing consisting clauses of from Table S2, where [6] must inject the trace payload into each event.

Clause	Declare Analyzer Representation	KnoBAB representation
$q_1$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$
$q_2$	$\text{Response}(\text{A\_SUBMITTED}, \text{AMOUNT\_REQ} \geq 10^3, \text{A\_ACCEPTED}, \text{true})$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$ $\text{Exists}(\text{\_\_trace\_payload}, \text{AMOUNT\_REQ} \geq 10^3, \geq 1)$
$q_3$	$\text{Response}(\text{A\_SUBMITTED}, \text{AMOUNT\_REQ} < 10^3, \text{A\_ACCEPTED}, \text{true})$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$ $\text{Exists}(\text{\_\_trace\_payload}, \text{AMOUNT\_REQ} < 10^3, \geq 1)$
$q_4$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$ $\text{where } \text{A\_SUBMITTED.org:resource} = \text{A\_ACCEPTED.org:resource}$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$ $\text{where } \text{A\_SUBMITTED.org:resource} = \text{A\_ACCEPTED.org:resource}$
$q_5$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$ $\text{where } \text{A\_SUBMITTED.org:resource} \neq \text{A\_ACCEPTED.org:resource}$	$\text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$ $\text{where } \text{A\_SUBMITTED.org:resource} \neq \text{A\_ACCEPTED.org:resource}$

**Table S2.** Declare Models for a data access best-case scenario for the BPIC\_2012 dataset. These queries are performed under the assumption the trace payload has been injected into every event of the trace.

Model Name	Conjunctive Query
$M_1$	$q_1 := \text{Response}(\text{A\_SUBMITTED}, \text{true}, \text{A\_ACCEPTED}, \text{true})$
$M_2$	$q_2 := \text{Response}(\text{A\_SUBMITTED}, \text{AMOUNT\_REQ} \geq 10^3, \text{A\_ACCEPTED}, \text{true})$
$M_3$	$q_3 := \text{Response}(\text{A\_SUBMITTED}, \text{AMOUNT\_REQ} < 10^3, \text{A\_ACCEPTED}, \text{true})$
$M_4$	$q_4 := q_1 \text{ where } \text{A\_SUBMITTED.org:resource} = \text{A\_ACCEPTED.org:resource}$
$M_5$	$q_5 := q_1 \text{ where } \text{A\_SUBMITTED.org:resource} \neq \text{A\_ACCEPTED.org:resource}$
$M_6$	$q_1 \wedge q_2$
$M_7$	$q_1 \wedge q_2 \wedge q_4$
$M_8$	$q_1 \wedge q_3 \wedge q_4$
$M_9$	$q_1 \wedge q_2 \wedge q_5$
$M_{10}$	$q_1 \wedge q_3 \wedge q_5$
$M_{11}$	$q_1 \wedge q_2 \wedge q_3 \wedge q_4 \wedge q_5$