

Worst-Case Execution Time Calculation for Query-Based Monitors by Witness Generation

MÁRTON BÚR, McGill University, Canada

KRISTÓF MARUSSY, Budapest University of Technology and Economics, Hungary

BRETT H. MEYER, McGill University, Canada

DÁNIEL VARRÓ, McGill University, Canada

Runtime monitoring plays a key role in the assurance of modern intelligent cyber-physical systems, which are frequently data-intensive and safety-critical. While graph queries can serve as an expressive yet formally precise specification language to capture the safety properties of interest, there are no timeliness guarantees for such auto-generated runtime monitoring programs, which prevents their use in a real-time setting. While worst-case execution time (WCET) bounds derived by existing static WCET estimation techniques are safe, they may not be tight as they are unable to exploit domain-specific (semantic) information about the input models. This paper presents a semantic-aware WCET analysis method for data-driven monitoring programs derived from graph queries. The method incorporates results obtained from low-level timing analysis into the objective function of a modern graph solver. This allows the systematic generation of input graph models up to a specified size (referred to as *witness models*) for which the monitor is expected to take the most time to complete. Hence the estimated execution time of the monitors on these graphs can be considered as safe and tight WCET. Additionally, we perform a set of experiments with query-based programs running on a real-time platform over a set of generated models to investigate the relationship between execution times and their estimates, and compare WCET estimates produced by our approach with results from two well-known timing analyzers, aiT and OTAWA.

CCS Concepts: • **Computer systems organization** → **Real-time system specification**; • **Software and its engineering** → *Automated static analysis*; Model-driven software engineering.

Additional Key Words and Phrases: real-time systems, worst-case execution time analysis, graph queries, model generation

ACM Reference Format:

Márton Búr, Kristóf Marussy, Brett H. Meyer, and Dániel Varró. 2018. Worst-Case Execution Time Calculation for Query-Based Monitors by Witness Generation. *J. ACM* 37, 4, Article 111 (August 2018), 36 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Runtime monitoring has become a key technique in the assurance of safety-critical and intelligent cyber-physical systems (CPS) such as autonomous vehicles [48] (e.g., self-driving cars, drones) where traditional upfront design time verification is problematic due to the dynamically changing

Authors' addresses: Márton Búr, marton.bur@mail.mcgill.ca, McGill University, 3480 Rue University, Montreal, Quebec, H3A 2K6, Canada; Kristóf Marussy, marussy@mit.bme.hu, Budapest University of Technology and Economics, Magyar tudósok körútja 2, Budapest, 1117, Hungary; Brett H. Meyer, brett.meyer@mcgill.ca, McGill University, Canada; Dániel Varró, daniel.varro@mcgill.ca, McGill University, Canada.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

0004-5411/2018/8-ART111 \$15.00

<https://doi.org/10.1145/1122445.1122456>

environment and the data-intensive nature of the system. However, it is an open challenge how to align real-time requirements with the ability to capture the highly dynamic operation context of the system. A promising line of research aims to leverage high-level, model-based techniques to manage system complexity and overcome current limitations [34, 58].

In traditional embedded systems, runtime monitoring programs are integral components of the system that analyze events and execution traces [5] in order to detect potentially critical situations that violate a requirement. Since this requires formal precision to capture safety requirements, logic-based formalisms (e.g., propositional logic, temporal logic) are frequently used to specify execution traces. Furthermore, monitoring programs can be automatically synthesized from such specifications that are ready to be used in traditional hard real-time systems without compromising task schedulability and real-time properties of the existing program [30, 49].

However, existing runtime monitoring approaches used in safety-critical applications have certain limitations, which are increasingly problematic for the new generation of data-intensive, intelligent, and self-adaptive, yet safety-critical CPSs. First, the *moderate expressiveness of the specification language* [29] makes it difficult for engineers to capture and understand complex rules. Moreover, while *safety-critical programs typically use statically allocated data with bounded input sizes* and they conservatively avoid many programming language constructs, dynamically evolving data and advanced language constructs are inherent parts of data-intensive programs.

Recent advances in runtime monitoring aim to overcome these limitations by (1) offering high-level and expressive query-based [8, 17] or rule-based [29] formalisms to capture the properties to be monitored, and (2) using runtime graph models as an in-memory knowledge base which capture dynamic changes in the system or its environment at a high-level of abstraction [8, 28]. Such *data-driven safety monitors* derived from high-level specifications can analyze aggregated changes triggered by complex sequences of atomic events by evaluating queries over a continuously evolving data model. For example, in the railway domain, queries can check if a path exists between two points along a railway track, or identify cargo waiting at stations for more than a specified duration. As such, query-based monitoring programs use a network of linked objects as data structures and *exhibit heavily input-dependent and semantic-aware complex control and data flow*.

To enable the use of data-driven safety monitors in hard real-time systems, the computation of safe worst-case execution time (WCET) estimates is required. While recent research has investigated data-driven runtime monitors for intelligent and critical CPSs in a distributed environment [8, 17, 28], and various testing approaches have been proposed [1, 55], the timeliness aspect of the problem has been neglected. In fact, only very few initial ideas are available [60], which suggest limiting the maximum graph size and employing optimized query plans. A wide range of existing timing analysis techniques and tools (e.g., aiT [21], Chronos [39], OTAWA [11], and SWEET [43]) can provide safe and tight WCET bounds for traditional critical embedded software. However, there is a high degree of inherent design-time uncertainty present in data-driven monitoring programs. In particular, the unknown contents of the runtime knowledge graph capturing the system and its operating environment constitutes an enormously large input space which can compromise the accuracy of existing techniques. Therefore, novel techniques are needed to complement existing WCET analysis techniques to efficiently incorporate domain-specific restrictions for program inputs on a high-level of abstraction and automatically incorporate this domain knowledge as flow facts.

In order to obtain safe and tight WCET bounds for data-driven runtime monitors, major challenges in timing analysis need to be tackled. (i) First, domain-specific flow constraints pose several complex restrictions on the program flow, but the respective flow facts need to be manually formulated and the program needs to be annotated by experts. Such additional program flow information largely helps to enhance the precision of safe WCET bounds in existing timing analyzers. However, there are no generally applicable methods to automatically obtain such flow facts by exploiting high-level,

domain-specific information and constraints on program flow during timing analysis. Specifying the flow facts manually is highly error-prone and the resulting annotations need to be updated after subsequent modifications to the program [2]. (ii) Furthermore, runtime graph models have varying underlying structure and memory demands which makes WCET analysis problematic since the worst-case graph structure needs to be considered regardless of domain-specific constraints. While memory can be preallocated to allow the timing analyzer to produce WCET bounds [31], but a WCET estimate from the size of the preallocated memory for graph data without appropriate flow facts would still be overly conservative. Moreover, the contents of the runtime model are regularly updated at runtime. Therefore, value analysis has no upfront access (at design-time) to the data that the reserved memory space will store. (iii) Finally, traditional WCET analysis challenges also need to be tackled: detailed information is required about the executable binary and execution platform, including precise memory, pipeline, and cache descriptions [65].

Contributions. This paper aims to address WCET estimation in the challenging setting of query-based runtime monitoring programs. In particular, we present the following contributions.

- (1) We adapt query-based runtime monitoring programs derived from high-level graph query specifications [8] to real-time platforms (Section 3.3).
- (2) We provide a novel high-level static analysis technique for query-based monitors to estimate execution time on a given runtime model. The approach provides precise flow facts by counting the number of basic block executions during query evaluation w.r.t. the given model even if exact memory allocation information is unavailable. Moreover, it combines such flow facts with constraints obtained from existing low-level analysis tools (Section 5.3).
- (3) We estimate the WCET of query-based programs by estimating execution time over designated *witness models*. Such witness models have the highest estimated execution times for graph query programs executing over any input models up to a predefined model size and domain-specific constraints, and they are derived by a state-of-the-art *graph solver* [55] (Section 5.4).
- (4) We perform an extensive experimental assessment of query evaluation times over a variety of graph models executed on an industry-grade real-time platform, and we compare our WCET estimates with those provided by two popular timing analyzers (OTAWA and aiT) (Section 6).

Novelty. Our technique complements existing WCET estimation methods by extracting program flow information from domain-specific constraints of the abstract input space on top of existing constraints derived by traditional timing analysis. To our best knowledge, our approach is the first to provide safe and tight WCET bounds of real-time graph query programs by abstraction refinement using state-of-the-art model generation techniques. This enables the use of query-based runtime monitoring programs in a real-time context by providing safe WCET estimates for a desired set of input models satisfying domain-specific constraints.

2 RELATED WORK

Numerous static and probabilistic WCET analysis methods have been discussed in surveys [37, 65]. Abella et al. [2] compares the most common WCET estimation approaches for programs in real-time systems and highlights their strengths and limitations. Based on the categorization of approaches of this latter work, our approach is a *high-level, static deterministic timing analysis (SDTA)* which provides *safe* execution bounds for embedded programs executing complex graph queries. Furthermore, measurement-based WCET estimations [38, 64] and probabilistic methods [16, 27] are out of scope for our work. Nevertheless, we focus on *semantic-aware WCET estimation* [14], which aims at providing safe and tight estimates for programs where there are some semantic limitations on the input data, which cannot be automatically explored and exploited by current analysis

techniques, and often times manual annotations of the code are necessary. In our case, control flow is often constrained by complex rules which are based on various properties of graph models. We provide an overview of existing work related to *graph-based programs*, *runtime monitoring* and *program flow analysis*.

Graph query programs: *Existing platforms.* Graph models and queries have been often used in design models and tools of real-time systems [9, 25, 35]. Furthermore, graph-based techniques are used in various IoT and edge computing applications [40, 66]. However, due to the soft real-time requirements of such applications, the WCET analysis aspect is often neglected. The focus of our work is to provide safe and tight WCET of such programs, and thus extend their application area.

Graph query programs: *WCET estimation.* One of the few related works that investigates real-time properties of graph-based techniques is [10]. Motivated by the expressiveness of *story diagrams* [22], the authors evaluate the applicability of this high-level modeling formalism to recognize hazardous situations in real-time systems. Their work investigates worst-case execution times of imperative programs generated from such story diagrams by executing measurements of manually created worst-case inputs. In contrast, our work aims to automatically synthesize worst-case well-formed input models as part of static analysis.

Runtime monitoring: *Hard real-time embedded systems.* One of the earlier works in the field is The Temporal Rover [18]. This framework can generate monitoring code from temporal logic formulae with low overhead, but the verification of properties is done in a large part on a powerful remote host, while our method does not rely on any external component. The concept of *predictable monitoring* was introduced in [67] where static scheduling techniques were used to show that a monitor fits its allocated time frame, but the analysis of monitoring tasks is out of its scope which is the topic of this current paper. Finally, synchronous component execution and observable program states are the main assumptions made in [49] to support sampling-based monitoring of input streams in real-time systems, whereas our work targets monitors executing complex queries over a graph model capturing contextual information on a high-level of abstraction.

Runtime monitoring: *Real-time database queries.* In real-time databases [47], access to data has strict time constraints. The work in [32] presents a data sampling-based statistical method to evaluate aggregate queries in a database. There is a trade-off between time available for query execution and the precision of the estimate. Such estimations would not be acceptable in a monitoring setting where precise query results are expected. The real-time object-oriented database RODAIN [57], which targets telecommunication applications, does not support *hard real-time* transaction (i.e., query) types, because it is considered too costly for the target domain. However, our objective is exactly to provide such guarantees over graph models to support hard real-time applications.

Program flow analysis. Timing analyzers for program flow analysis often employ some version of the implicit path enumeration technique (IPET) [41]. The general idea behind this method is to use the control flow graph (CFG) of the program to create an integer linear program (ILP) where each variable encodes the number of executions of a corresponding basic blocks, and the objective function is to maximize their total execution time. Besides the IPET method, several *tree-based* methods exist which use a tree representation of the program (obtained from the source code or compiled binary) and apply some traversal to find the longest path in a program [4, 15, 42]. In any case, the effectiveness of these methods rely on precise program *flow facts* (e.g., loop bounds, infeasible paths) to be able to determine a safe and tight WCET estimate. Although there are several advanced (semi-)automated techniques available today to derive additional constraints on the program flow and thus improve the precision of the WCET estimate [13, 20, 26, 36, 43], there is

still a significant manual effort needed to specify flow facts [2]. Most closely related to our current work is [36], which uses abstraction refinement to reduce WCET estimates by *squeezing*. However, this approach cannot exclude longest execution paths from the program which are infeasible due to complex domain-specific constraints on the inputs.

3 QUERY-BASED RUNTIME MONITORS

In this section, we provide an informal overview of query-based runtime monitors, while their formal treatment is deferred to Section 4.

3.1 Running Example: the MoDeS3 CPS Demonstrator

Our key concepts are illustrated in the context of the open source *Model-Based Demonstrator for Smart and Safe Cyber-Physical Systems* (MoDeS3) [63] platform, which showcases various challenges of modern intelligent yet safety-critical CPS applications. The demonstrator (see Figure 1) is a model railway system with an added layer of safety to prevent train collision and derailment using runtime monitors. The railway track is equipped with several sensors and actuators, which are represented by black triangles in the lower part of Figure 1. Train shunt detectors can sense when trains pass by on a particular segment of the track, while direction of turnouts can be read and set.

The system is managed by a (distributed) monitoring service running on a network of heterogeneous computing units, such as Arduinos, Raspberry Pis, BeagleBone Blacks, etc. Relevant runtime information gained from sensor reads (e.g., the occupancy of a segment, or the status of a turnout) is uniformly captured in an in-memory *runtime graph model*, which is also deployed on the platform.

Safety monitors are formally captured as *graph queries*. Alerts from the monitoring services may trigger control commands of actuators (e.g., to change turnout direction) to guarantee safe operation. The monitoring and control programs are running in a real-time setting on the computing units.

While the MoDeS3 platform can demonstrate various challenges of CPSs, this paper exclusively focuses on the real-time aspect of query-based runtime monitoring programs deployed to some embedded devices with limited resources (memory, CPU, etc).

3.2 Graph Models at Runtime

3.2.1 Graph Models. The **models@run.time** paradigm [6] facilitates the capture of runtime knowledge about the system and its environment as a (typed and directed) *graph model* continuously maintained at runtime for the system. Such graphs are dynamically changing in-memory data structures which encode domain-specific instance models typed over a *domain metamodel*, which captures core concepts (classes) in a domain and the relations (references) between those concepts.

Example 3.1. The domain concepts of the MoDeS3 runtime model are captured in a metamodel excerpt shown in Figure 2(a) using the Eclipse Modeling Framework (EMF) notation [59]. One

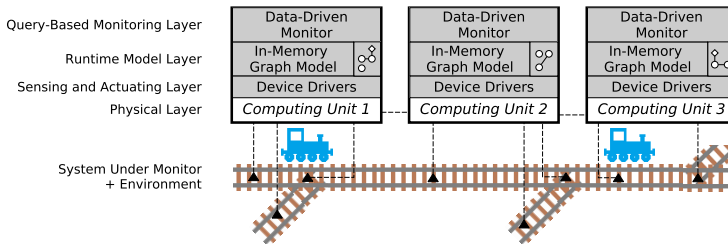


Fig. 1. Runtime monitoring by graph queries

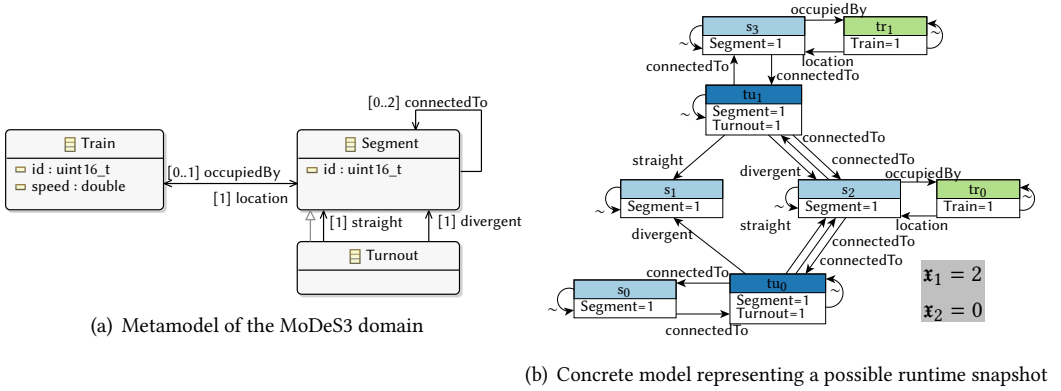


Fig. 2. An instance (concrete) model and a partial model in the MoDeS3 domain

domain concept is Train. Class Segment represents a section of the railway track with the `connectedTo` reference which describes what other segments it is linked to (up to two). Moreover, each train maintains a location reference to a segment to describe its current position. Instances of class Segment record if they are occupied by a train with the `occupiedBy` reference. Moreover, Turnout is a special Segment that can change its connections between straight and divergent segments.

A runtime (instance) model captures a snapshot of the underlying system in operation [6, 56]. Relevant changes in the system are reflected in the runtime model and operations executed on the runtime model (e.g., setting values of controllable attributes of objects or updating links between objects) are reflected in the system itself (e.g., by executing scripts or calling services). In this work, we use *concrete (graph) models* to formally capture runtime model snapshots.

Example 3.2. Figure 2(b) shows a concrete model in a graphical syntax. The graph has six Segment objects (including two Turnouts) with their respective `connectedTo` links. Turnouts tu_0 and tu_1 can switch between segments s_1 and s_2 (see their straight and divergent edges). In the depicted state, both turnouts are switched to s_2 and the trains tr_0 and tr_1 are located on s_2 and s_3 , respectively.

3.2.2 Graph Data Structures in Embedded Systems. Runtime monitors captured by graph queries are continuously evaluated over the runtime models. This section informally summarizes our assumptions and requirements about such programs while the theoretical background is introduced in Section 4.

For data-driven monitors, the structure of the underlying graph model directly impacts the performance of query evaluation. Since an embedded device may have limited available CPU and memory resources, a lightweight data structure is needed to efficiently capture runtime graph models. While the in-depth discussion of such a graph data structure is out of scope for this paper, we make the following assumptions about the supported operations of the underlying graph:

- **Dynamic element creation and deletion.** The runtime model serves as the knowledge base about the underlying system and its environment. For this reason, it needs to accommodate graph models without a theoretical a priori upper bound for model size. Based on [31], one way to support this is to allocate the maximum amount of memory that is physically possible to be used for storing the graph. However, only the allocated memory is determined at compile time, the type (and distribution) of objects stored in the graph is runtime information.

```

1 typedef struct {
2   uint16_t segment_id;
3   Train *train;
4   Segment *connected_to[2];
5   uint8_t connected_to_count;
6 } Segment;

```

Listing 1. Segment class

```

1 typedef struct {
2   uint16_t train_id;
3   double speed;
4   Segment *location;
5 } Train;

```

Listing 2. Train class

```

1 struct Modes3ModelRoot {
2   Segment segments[SEGMENTS];
3   uint16_t segment_count;
4   Train trains[TRAINS];
5   uint16_t train_count;
6 } runtime_model;

```

Listing 3. Graph model root

- **Indexing of objects by type using unique identifiers.** As query evaluation typically starts by iterating over all elements of a given type or accessing specific objects, it necessitates efficient object access, e.g. by maintaining a real-time index for memory resident data [12].
- **Navigability along edges.** Query evaluation often navigates along the edges of selected objects to find further appropriate variable substitutions for unbound query variables. This feature can be supported by, e.g., maintaining direct pointers to reachable objects.

It is also important to note that the same graph model can be represented in memory in many ways, because different placements of the same data can cause different run times. For example, two memory images of the same graph may differ in the order the objects are stored in the array. For this reason, two different in-memory representations of the same graph may not necessarily yield identical run times, which must be considered when computing WCET of graph query programs.

Example 3.3. Listing 1 and Listing 2 show a possible C implementation of data structures for Segment and Train classes in the metamodel of Figure 2(a). Line 2 in Listing 1 and lines 2 and 3 in Listing 2 are fields created from respective attributes present in the metamodel, e.g., the speed attribute of class Train is represented by line 3. For each type, an id attribute encodes the type of the object for indexing and model manipulation purposes. Uniqueness of this attribute needs to be guaranteed at runtime to distinguish objects. Furthermore, in this example, we implement graph edges as pointers (line 4 in Listing 2) or pointer arrays with sizes (lines 4 and 5 in Listing 1). Representing links between objects with pointers is highly efficient from a performance viewpoint.

Listing 3 shows how a simple graph model container Modes3ModelRoot can allocate static memory for graph objects in C. The maximum memory used by the graph is preallocated in lines 2 and 4 by the segments and trains arrays which have a length of the maximum expected number of trains (denoted by the constant TRAINS) and the maximum expected number of segments (SEGMENTS). The id attribute of a given object used for indexing these arrays, i.e., encodes their positions in the arrays.

3.3 Graph Query Programs in Real-Time Systems

Data-driven runtime monitors defined by graph queries can check structural properties of the runtime model representing a snapshot of the system. In other words, they focus on the most up-to-date data (maintained either by periodic updates with a certain frequency, or by certain event-driven triggers [6, 56]) available on the underlying system's state at a given point of time.

Classical *event-based runtime monitors* rely on some temporal logic formalism to detect sequences of events occurring in the system at different points in time, while the underlying data model is restricted to atomic propositions. As such, data-driven and event-based monitors are complementary techniques. While graph queries can be extended to express temporal behavior, our current work is restricted to (structural) safety properties where safety violations are expressible by graph queries.

3.3.1 Graph Queries. A *graph query* is a declarative description of a model fragment to be identified by a set of variables and a set of constraints (type, reference, and equality assertions) [62]. A *match* of the query is a binding of the query variables to objects in the model such that the constraints are

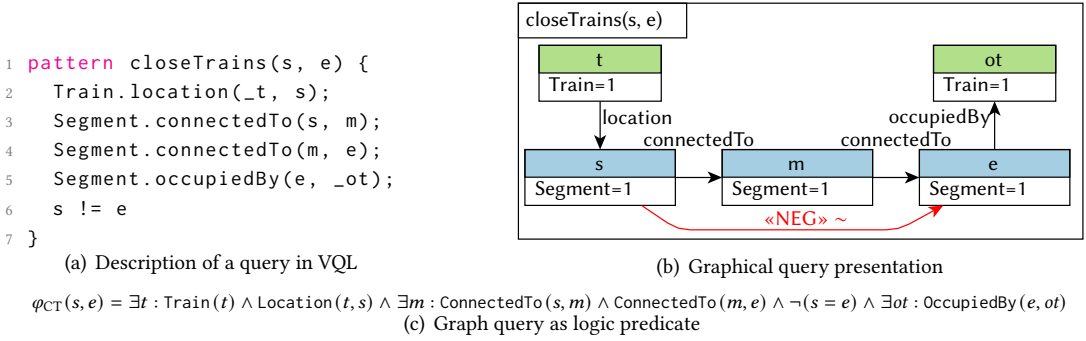


Fig. 3. Monitoring goal formulated as a graph query φ_{CT} for `closeTrains`

satisfied. In data-driven runtime monitors, such high-level descriptions allow us to automatically generate and optimize the monitor program by adaptive *query planning*.

Example 3.4. Trains are required to keep long headway distances to ensure that they can safely decelerate without collision [19]. The safety case captured by the *closeTrains* (CT) graph query represents a situation when the headway distance between trains located on connecting segments is reduced below the safety limit. Any match of this query highlights segments where immediate action (e.g., stopping the trains) is required. The declarative query specification is presented in Listing 3(a) in a textual syntax to identify the violating situation as a hazardous case. The query returns pairs of segments s, e where there is a train located on segment s that is one segment away (i.e., there is a middle segment m) from a different segment e , which is also occupied by a train. Any variables not appearing in the parameter list of the query are existentially quantified.

Figure 3(b) shows the same query in a graphical presentation often employed by modeling tools, and Figure 3(c) presents it as a first-order logic (FOL) formula φ_{CT} (discussed later in Section 4.2).

In the runtime snapshot Figure 2(b), the variable bindings $\{s \mapsto s_2, e \mapsto s_3\}$ and $\{s \mapsto s_3, e \mapsto s_2\}$ are matches of the *closeTrains* query.

3.3.2 Local search-based graph query evaluation. Among the many possible query evaluation strategies [24], our runtime monitoring framework uses *local search-based query evaluation* [62] to find matches of monitoring over the entire runtime model. This strategy at its core uses a tailored depth-first search graph traversal. This keeps the memory footprint of the query evaluation algorithm constant. To obtain efficient performance at runtime, query evaluation is guided by a *search plan* [62], which maps each constraint in the query to a single pair of $\langle \text{Step index}, \text{Operation type} \rangle$. In this tuple, *Step index* specifies the order in which query evaluation should attempt to satisfy the respective constraint. *Operation type* can be one of the followings:

- An **extend operation** evaluates a constraint with at least one free variable. Execution of such operations requires iterating over all potential variable substitutions and selecting the ones for which the constraint evaluates to 1.
- A **check operation** evaluates a constraint with only bound variables. Execution of such operations determines if the constraint evaluates to 1 over the actual variable binding.

Example 3.5. Table 1 shows a possible search plan for the φ_{CT} query. Each row represents a search operation. The first column shows which constraint is enforced by the given step where free parameters at the start of the execution of the operation are underlined. The second column shows the ordering of steps, i.e., the step index, and the third column shows the search operation

Algorithm 1: Code generation from search plans

```

1 Function CompileSearchPlan(sp, idx) is
2   if idx > sp.size() then return code for storing a match;
3   step = sp[idx]
4   matcherCode = ""
5   if step is extend then
6     for uv ∈ step.getFreeVariables() do
7       matcherCode +=
8         AddAssignmentFor(uv, step.getConstraintFor(uv))
9   else if step is check then
10    matcherCode +=
11      AddIfFor(step.getAllVariables(), step.getConstraint())
12  return matcherCode + CompileSearchPlan(sp, idx + 1)

```

Table 1. A possible search plan for query closeTrains where free variables are underlined

Constraint	Step#	Op. type
Train(<u>t</u>)	1	extend
Location(<u>t</u> , <u>s</u>)	2	extend
ConnectedTo(<u>s</u> , <u>m</u>)	3	extend
ConnectedTo(<u>m</u> , <u>e</u>)	4	extend
$\neg(s = e)$	5	check
OccupiedBy(<u>e</u> , <u>ot</u>)	6	extend

type (check or extend) which is based on the variable bindings prior to the execution of the search operation: if the constraint parameters are all bound, then it is a check, otherwise, it is an extend.

3.3.3 Implementations of Query Programs. Although constructing effective search plans for graph queries is a complex challenge, it is outside of the scope of the current paper and has been formerly extensively studied (see, e.g., [62] for a possible solution). However, we present pseudo-code that generates embedded query code from a search plan in Algorithm 1. The function *CompileSearchPlan* takes a search plan and a search step index as parameters. Line 2 returns a code snippet to register a match if the provided index is beyond the index of the final search step. Otherwise, the search step is extracted (line 3) and the variable *matcherCode* to hold the generated code is initialized to an empty string (line 4). Then, the different operation types of the query search plan are translated to structured imperative code:

- Each **extend operation** binds all free variables of the respective constraint (lines 5–6). For each variable, this translates to either a single assignment or a for loop iterating over a set of candidate variable bindings, depending on the multiplicity of the respective navigation edge (reference constraint) (lines 7–8).
- Each **check operation** (line 9) is mapped to an if statement to check if the current variable binding satisfies a given condition created from the query constraint (lines 10–11).

Finally, in line 12, the generation continues recursively appending the code generated from the subsequent steps to the result. The query code for the entire search plan *sp* can be generated by calling *CompileSearchPlan*(*sp*, 1). As a result, the source code contains a deep hierarchy of embedded for-loops and if-statements based on the ordering of constraints prescribed by the search plan.

Besides obtaining a WCET, we also need to estimate the number of matches of a query to allocate appropriate space in memory in advance. In the case of runtime monitors of safety properties, we can assume that only a few violating matches will be detected [61], thus the query result set is expected to be small and memory required for storing matches can be reserved at compile time.

Example 3.6. Listing 4 shows the C code generated from the query specification of closeTrains. Assuming that a global variable *model* points to the root of the entire graph model including its up-to-date model statistics, calling the function *close_trains_matcher* with a pointer to the result set structure *results* will compute and store all matches over the model in *results*.

In the example, initially all variables are assumed to be free, as indicated in line 2 with NULL values, because we aim to find all matches in the entire model. In line 3, the size of the result set is initialized to 0. The for loop in line 6 represents step 1 from the search plan (see Table 1) and iterates over all trains in the model, binding the variable *vars.t* to all possible objects in line 7. Lines 8–10 together represent search step 2. In line 9, *vars.s* is assigned a segment referred by

```

1 void close_trains_matcher(CloseTrainsMatchSet *results) {
2   CTVars vars = {t=NULL, ot=NULL, s=NULL, m=NULL, e=NULL};
3   int match_cntr = 0;
4   // Constraint:  $\exists t: \text{Train}(t)$ 
5   int loop_bound0 = model->train_count;
6   for (int i0 = 0; i0 < loop_bound0; i0++) {
7     vars->t = model->trains[i0];
8     // Constraint:  $\text{Location}(t,s)$ 
9     vars->s = vars->t->location;
10    if (vars->s != NULL) {
11      // Constraint:  $\exists m: \text{ConnectedTo}(s,m)$ 
12      int loop_bound1 = vars->s->connected_to_count;
13      for (int i1 = 0; i1 < loop_bound1; i1++) {
14        vars->m = vars->s->connected_to[i1];
15        // Constraint:  $\text{ConnectedTo}(m,e)$ 
16        int loop_bound2 = vars->m->connected_to_count;
17        for (int i2 = 0; i2 < loop_bound2; i2++) {
18          vars->e = vars->m->connected_to[i2];
19          // Constraint:  $\neg(s=t)$ 
20          if (vars->s != vars->e) {
21            // Constraint:  $\exists ot: \text{OccupiedBy}(e,ot)$ 
22            vars->ot = vars->e->train;
23            if (vars->ot != NULL) {
24              // Register match
25              results->matches[match_cntr].s = vars->s;
26              results->matches[match_cntr++].e = vars->e;
27            } } } } } }
28    results->size = match_cntr; }

```

Listing 4. Source code generated for query closeTrains

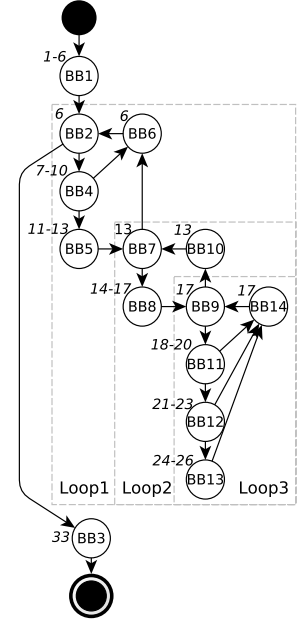


Fig. 4. CFG of example query program

`vars.t` via a single location link. If such a segment exists in line 10, execution continues with the third search operation that is mapped to lines 11–14, which iterates over segments connected to `vars.s` and assigns them to `vars.m`, one at a time. The next step in lines 15–18 does the same but with the connecting segments of `vars.m` and assigns them to `vars.e`. Search step 5 is a check, which is mapped to lines 19–20 to ensure that the segments referred by `vars.s` and `vars.e` are not the same. The final step of the search plan is mapped to lines 21–23. Here the train occupying the segment stored in `vars.e` is assigned to `vars.ot`. If such a train exists, a match is found and registered by assigning the corresponding variable values to parameter variables in a new match (lines 24–26) and incrementing the matches found counter `match_cntr`. The execution concludes with saving the number of matches (line 28).

Static analysis of the query code itself in Listing 4 would not impose any restrictions on line 20 despite the fact that the domain-specific constraints prescribe `connectedTo` links to be symmetrical. Therefore, at least every other execution of line 20 will jump back to line 17 instead of proceeding to line 22, yielding a flow constraint that is not discoverable by analysis of the code only.

Cyclomatic complexity (CC) is frequently used as a metric in safety-critical software to estimate code complexity [51]. As a general recommendation, code with high CC is traditionally avoided in a safety-critical system as it requires extra efforts to test and maintain. However, the derived imperative source code of data-driven monitoring programs is inherently complex even for small queries, which is largely attributed to the declarative nature of query specifications. For example, the CC of Listing 4 is 7, which already indicates substantial complexity.

While modern WCET analyzers can analyze complex code fragments, they heavily rely on manual annotation of the code (loop bounds, in particular) and design time information about variable values to be able to come up with an estimate that is both safe and tight. A key contribution of the

current paper is to complement the existing WCET analysis by providing means to automatically exploit domain-specific restrictions of input data and tighten the resulting WCET estimate. For data-driven monitors, this is a key step in order to enable their use in a safety-critical context.

4 FORMAL BACKGROUND

This section provides the formal background for the static analysis of data-driven runtime monitors, introduces definitions from traditional IPET-based approaches for WCET estimation, and revisits the state of the art of domain-specific graph modeling and graph model generation.

4.1 Implicit Path Enumeration Technique for Estimating WCET

Timing analysis frequently relies on the Implicit Path Enumeration Technique (IPET) [41] that uses the control flow graph (CFG) of a program to estimate the WCET. This section revisits definitions from [36, 50] to introduce this classic WCET estimation approach.

Definition 4.1. A program is a pair $\langle BB, Loops \rangle$, where BB is the set of *basic blocks* and $Loops$ is the set of (well-structured) loops. Each loop $\ell \in Loops$ has a *header* $bb_{\ell,h} \in \ell$, while the rest of its blocks $bb_i \in \ell$ constitute its *body*.

Definition 4.2. A *weighted control flow graph* corresponding to a program $\langle BB, Loops \rangle$ is a tuple $CFG = \langle V, E, s, t, w, tr \rangle$, where

- V is a finite set of *nodes*;
- $E \subseteq V \times V$ is the set of *edges*;
- s and $t \in V$ are the *program start* and *end* nodes, respectively;
- $w: E \rightarrow \mathbb{N}$ is the *weight function*, that assigns execution times (clock cycles) to the edges;
- $tr: V \rightarrow BB$ is the *traceability* function that maps nodes to their originating program blocks.

In the simplest case, $V = BB$ and tr is the identity function. However, even on simple embedded processors, basic block execution times may vary due to microarchitectural effects (e.g., pipeline or cache state), which necessitates representing basic blocks with multiple nodes to encode context-sensitive execution times. The extended CFG encodes information from a *low-level timing analysis*. For example, the VIVU approach [44] addresses this concern by *virtual loop unrolling* and creates additional nodes and edges in the CFG to explicitly model the first executions of loops.

Definition 4.3. Every execution of the program (i.e., program trace) can be represented by a path $\pi = \langle e_1, \dots, e_m \rangle$ in the CFG from s to t . Let $\pi \# e$ denote the frequency of the edge e appearing in π . The program execution time $\tau(\pi)$ can be estimated by the sum of the product of weights and frequencies of edges, i.e., $\tau(\pi) = \sum_{e \in E} w(e) \cdot \pi \# e$.

The goal of timing analysis is to find the path with the longest possible execution time of the program. Instead of explicitly enumerating all possible paths in the CFG, Puschner and Schedl [41] create a system of linear inequalities whose solutions over-approximate the set of possible paths and define an integer linear programming (ILP) problem for estimating WCET. First, we review some notations for systems of linear inequalities and ILP problems that will be used in this paper.

Definition 4.4. Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_{|\mathcal{X}|}\}$ be a large (but finite) *reserve* of linear equation variable symbols. Then $\mathcal{S} = \left\{ \sum_{\mathbf{x}_j \in \mathcal{X}} a_{i,j} \cdot \mathbf{x}_j \leq y_i \right\}_{i=1}^{|\mathcal{S}|}$ is a *system of linear equations*.

Definition 4.5. Let functions $k: \mathcal{X} \rightarrow \mathbb{Z}$ be *valuations*. A valuation k is a *solution* of \mathcal{S} (written as $k \models \mathcal{S}$) if $\sum_{\mathbf{x}_j \in \mathcal{X}} a_{i,j} \cdot k(\mathbf{x}_j) \leq y_i$ for all inequalities $\sum_{\mathbf{x}_j \in \mathcal{X}} a_{i,j} \cdot \mathbf{x}_j \leq y_i$ in \mathcal{S} . The system of linear equations \mathcal{S}_1 *entails* \mathcal{S}_2 (written as $\mathcal{S}_1 \models \mathcal{S}_2$) if, for any solution $k \models \mathcal{S}_1$, we also have $k \models \mathcal{S}_2$.

Definition 4.6. An *integer linear program* (ILP) is of the form $\max \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot \mathbf{x}_i$ *subject to* \mathcal{S} , where \mathcal{S} is a system of linear equations.

While ILP with both maximization (max) and minimization (min) objectives can be considered, we restrict our attention to max objectives w.l.o.g., since any min objective can be transformed into a max objective after multiplication by -1 .

Definition 4.7. Let $g(k) = \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot k(\mathbf{x}_i)$ denote the *cost* of valuation k . The valuation k^* is an *optimal solution* of the ILP if $k^* \models \mathcal{S}$ and $g(k^*) \geq g(k)$ for all $k \models \mathcal{S}$, i.e., its cost is maximal.

An ILP can be derived from the CFG $\langle V, E, s, t, tr \rangle$ as follows. Let $\mathfrak{f}: E \rightarrow \mathcal{X}$ be a function that associates variable symbols to CFG edges. In the system of linear equations \mathcal{S} , each feasible execution path P is associated with a solution k_π such that $k_\pi(\mathfrak{f}(e)) = \pi \# e$. Thus, linear constraints on variable obtained as $\mathfrak{f}(e)$ restrict frequency of e in all feasible paths. Moreover, $\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb_i} k_\pi(\mathfrak{f}(e))$ is the number of times the basic block $bb_i \in BB$ was executed in π .

- We add $\sum_{e=\langle s, n \rangle \in E} \mathfrak{f}(e) = 1$ and $\sum_{e=\langle n, t \rangle \in E} \mathfrak{f}(e) = 1$ to \mathcal{S} , because the program is entered and exited exactly once.
- Except for s and t , each node is entered and exited the same number of times, so we add $\sum_{e=\langle n_1, n_2 \rangle \in E} \mathfrak{f}(e) - \sum_{e=\langle n_2, n_3 \rangle \in E} \mathfrak{f}(e) = 0$ for each $n_2 \in V \setminus \{s, t\}$.
- Any additional *flow facts* regarding the execution frequency of program parts (e.g., loop execution counts) are added in the form $\sum_{e_{i,j} \in E} a_{i,j} \cdot \mathfrak{f}(e_{i,j}) \leq y_j$.
- For each edge $e \in E$, we also have $-\mathfrak{f}(e) \leq 0$, since the execution frequency is non-negative.

We have an ILP $\max \sum_{e \in E} w(e) \cdot \mathfrak{f}(e)$ *subject to* \mathcal{S} . The objective function $g(k)$ overapproximates the execution time of P . Therefore the value $g(k^*)$ of any solution k^* is a WCET bound.

4.2 Metamodels and Partial Models

In order to extend static analysis of data-driven graph query programs with domain-specific flow information, we will formally capture metamodels by a logic signature and their instance models as logic structures following [45, 55].

Definition 4.8. A *metamodel* is formally represented as a first-order logic signature $\langle \Sigma, \alpha \rangle$, where

- $\Sigma = \{C_1, \dots, C_{m_C}, R_1, \dots, R_{m_R}, \varepsilon, \sim\}$ is a finite set of *symbols*, where $\{C_i\}_{i=1}^{m_C}$ are unary *class symbols*, $\{R_j\}_{j=1}^{m_R}$ are binary *relation symbols*, ε is the *object existence* symbol, and \sim is the *object equality*;
- $\alpha: \Sigma \rightarrow \mathbb{N}$ is the *arity* function with $\alpha(C_i) = 1$ for all $i = 1, \dots, m_C$, $\alpha(R_j) = 2$ for all $j = 1, \dots, m_R$, $\alpha(\varepsilon) = 1$, and $\alpha(\sim) = 2$.

The definition of a metamodel may also include binary *attribute symbols* such as in [8]. However, their handling is analogous to binary relation symbols, thus their discussion is excluded from here.

Partial models explicitly capture uncertainty in models as well as the design decisions yet to be made using 3-valued logic [46, 52]. In addition to the usual 1 and 0 truth values, the $1/2$ truth value corresponds to uncertainties in the model. We also add systems of linear equations as *scopes* [45] to partial models to impose numerical constraints on the sizes of the models by *polyhedron abstraction*. Later, variables $\mathbf{x}_j \in \mathcal{X}$ in the scopes will be connected to the number of model objects and graph pattern matches through *theories* of 3-valued logic expressions.

Definition 4.9. A *scoped partial model* over a signature $\langle \Sigma, \alpha \rangle$ is a triple $P = \langle \mathcal{O}_P, \mathcal{I}_P, \mathcal{S}_P \rangle$, where

- \mathcal{O}_P is a finite set of *objects*;
- \mathcal{I}_P is a 3-valued logical interpretation $\mathcal{I}_P(\sigma): \mathcal{O}_P^{\alpha(\sigma)} \rightarrow \{1, 0, 1/2\}$ for all $\sigma \in \Sigma$; and
- the *scope* $\mathcal{S}_P = \left\{ \sum_{\mathbf{x}_j \in \mathcal{X}} a_{i,j} \cdot \mathbf{x}_j \leq y_i \right\}_{i=1}^{|\mathcal{S}_P|}$ is a system of linear equations.

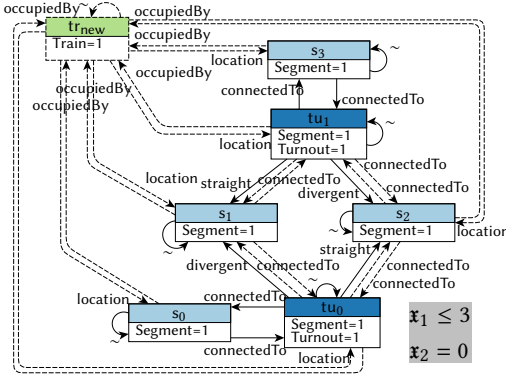


Fig. 5. Partial model with uncertain elements

$$\begin{aligned}
 \llbracket 1 \rrbracket_Z^P &:= 1 & \llbracket C_i(v) \rrbracket_Z^P &:= \mathcal{I}_P(C_i)(Z(v)) \\
 \llbracket 0 \rrbracket_Z^P &:= 0 & \llbracket R_j(v_1, v_2) \rrbracket_Z^P &:= \mathcal{I}_P(R_j)(Z(v_1), Z(v_2)) \\
 \llbracket \neg\varphi \rrbracket_Z^P &:= 1 - \llbracket \varphi \rrbracket_Z^P & \llbracket v_1 = v_2 \rrbracket_Z^P &:= \mathcal{I}_P(\sim)(Z(v_1), Z(v_2)) \\
 \llbracket \varphi_1 \vee \varphi_2 \rrbracket_Z^P &:= \max\{\llbracket \varphi_1 \rrbracket_Z^P, \llbracket \varphi_2 \rrbracket_Z^P\} \\
 \llbracket \varphi_1 \wedge \varphi_2 \rrbracket_Z^P &:= \min\{\llbracket \varphi_1 \rrbracket_Z^P, \llbracket \varphi_2 \rrbracket_Z^P\} \\
 \llbracket \exists v : \varphi \rrbracket_Z^P &:= \max_{o \in \mathcal{O}_P} \{\min\{\mathcal{I}_P(\varepsilon)(o), \llbracket \varphi \rrbracket_{Z, v \rightarrow o}^P\}\} \\
 \llbracket \forall v : \varphi \rrbracket_Z^P &:= \min_{o \in \mathcal{O}_P} \{\max\{1 - \mathcal{I}_P(\varepsilon)(o), \llbracket \varphi \rrbracket_{Z, v \rightarrow o}^P\}\}
 \end{aligned}$$

Fig. 6. 3-valued semantics of logic predicates

The existence symbol ε allows us to represent objects o that optionally appear in the model by setting $\mathcal{I}_P(\varepsilon)(o) = 1/2$. In contrast, objects with $\mathcal{I}_P(\varepsilon)(o) = 1$ surely appear. Uncertain equality $\mathcal{I}_P(\sim)(o, o) = 1/2$ of an object o with itself denotes *multi-objects* that can represent multiple concrete model objects. In contrast, objects with $\mathcal{I}_P(\sim)(o, o)$ stand for single concrete model objects.

Example 4.10. For the metamodel of Figure 2(a), $\{\text{Train}, \text{Segment}, \text{Turnout}\} \subseteq \Sigma$ are unary class predicates, and $\{\text{location}, \text{occupiedBy}, \text{connectedTo}, \text{straight}, \text{divergent}\} \subseteq \Sigma$ are binary predicates.

Figure 5 shows a partial model $P = \langle \mathcal{O}_P, \mathcal{I}_P, \mathcal{S}_P \rangle$ conforming to the MoDeS3 metamodel. Objects are drawn as boxes with the values of the interpretations $\mathcal{I}_P(C_i)$ of the class symbols written inside, while edges are drawn as arrows labelled with the relation symbols R_j and the equality symbol \sim . Solid edges correspond to 1 logic values, dashed edges are $1/2$ logic values, and 0 logic values are omitted. Uncertain existence ε is shown with a dashed outline.

The switching direction of the turnouts tu_0 and tu_1 is unknown. Additionally, there are no concrete trains on the track, but a *multi-object* tr_{new} represents all trains and their potential locations. Formally, $\mathcal{I}_P(\varepsilon)(s_0) = \mathcal{I}_P(\sim)(s_0, s_0) = \mathcal{I}_P(\text{Segment})(s_0) = \mathcal{I}_P(\text{connectedTo})(s_0, tu_0) = 1$, but $\mathcal{I}_P(\text{connectedTo})(tu_0, s_1) = 1/2$. Because $\mathcal{I}_P(\varepsilon)(tr_{\text{new}}) = \mathcal{I}_P(\varepsilon)(tr_{\text{new}}, tr_{\text{new}}) = 1/2$, tr_{new} is a *multi-object* that may stand for any number of Train instances (even 0). The location of tr_{new} is also uncertain, new trains may be located on any Segment or Turnout.

The scope $\mathcal{S}_P = \{\mathbf{x}_1 \leq 3, \mathbf{x}_2 = 0\}$ is shown in the lower right corner of the figure. In Theorem 4.21, this scope will restrict the number of Train objects in the model to ensure its well-formedness.

Runtime snapshots of a system are captured by *concrete (instance) models*, which contain no uncertainty or multi-objects and truth values are restricted to 1 and 0.

Definition 4.11. A partial model $M = \langle \mathcal{O}_M, \mathcal{I}_M, \mathcal{S}_M \rangle$ is *concrete* if

- \mathcal{I}_M contains only 1 and 0 values, i.e., $\mathcal{I}_M(\sigma)(o) \in \{1, 0\}$ for all $\sigma \in \Sigma$ and $o \in \mathcal{O}_M^{\alpha(\sigma)}$;
- all objects surely exist, i.e., $\mathcal{I}_M(\varepsilon)(o) = 1$ for all $o \in \mathcal{O}_M$;
- the interpretation of the equality symbol \sim matches the usual equality of objects, i.e., for all $o_1, o_2 \in \mathcal{O}_M$ $\mathcal{I}_M(\sim)(o_1, o_2) = 1$ if $o_1 = o_2$, 0 otherwise; and
- \mathcal{S}_M is satisfiable, i.e., there is some valuation $k: \mathcal{X} \rightarrow \mathbb{Z}$ such that $k \models \mathcal{S}_M$.

Example 4.12. Figure 2(b) shows a concrete model M conforming to the MoDeS3 metamodel. Only 1 and 0 logic values appear in the interpretation. The associated scope \mathcal{S}_M has a single solution $\{\mathbf{x}_1 \mapsto 2, \mathbf{x}_2 \mapsto 0\} \models \mathcal{S}_M$.

Concrete models are obtained from partial models by a series of *refinements* [53], which add further information by setting unknown logic values $1/2$ to 1 or 0, while known 1 and 0 values remain unchanged. This is captured by the refinement relation $X \succcurlyeq Y := (X = 1/2) \vee (X = Y)$. During model generation, refinements are carried out until a concrete model is reached.

Definition 4.13. The function $abs: O_Q \rightarrow O_P$ is an *abstraction function* from the partial model $Q = \langle O_Q, \mathcal{I}_Q, \mathcal{S}_Q \rangle$ to the partial model $P = \langle O_P, \mathcal{I}_P, \mathcal{S}_P \rangle$ (written as $P \succcurlyeq_{abs} Q$) if

- for all $\sigma \in \Sigma$, $q_1, \dots, q_{\alpha(\sigma)} \in O_Q$, logic values in \mathcal{I}_Q are the refinements of the corresponding values in \mathcal{I}_P , i.e., $\mathcal{I}_P(\sigma)(abs(q_1), \dots, abs(q_{\alpha(\sigma)})) \succcurlyeq \mathcal{I}_Q(\sigma)(q_1, \dots, q_{\alpha(\sigma)})$;
- surely existing object do not disappear, i.e., for all $p \in O_P$, $\mathcal{I}_P(\varepsilon)(p) = 1$ implies that there is some $q \in O_Q$ with $abs(q) = p$; and
- solutions of \mathcal{S}_Q are also solutions of \mathcal{S}_P , i.e., $\mathcal{S}_Q \models \mathcal{S}_P$.

Q is a *refinement* of P (written as $P \succcurlyeq Q$) if $P \succcurlyeq_{abs} Q$ for some $abs: O_Q \rightarrow O_P$.

Example 4.14. The concrete model M in Figure 2(b) is a refinement of the partial model P in Figure 5: $P \succcurlyeq_{abs} M$. The abstraction function maps $abs(tr_0) = tr_{new}$ and $abs(tr_1) = tr_{new}$, i.e., newly added trains are refinements of the train multi-object. Any other object is mapped by abs to itself, i.e., the identities of the rest of the objects remained unchanged. We may also see that the directions the turnouts were set, e.g., $\mathcal{I}_P(\text{connectedTo})(tu_0, s_1) = 1/2 \succcurlyeq \mathcal{I}_M(\text{connectedTo})(tu_0, s_1) = 0$ and $\mathcal{I}_P(\text{connectedTo})(tu_0, s_2) = 1/2 \succcurlyeq \mathcal{I}_M(\text{connectedTo})(tu_0, s_2) = 1$. Moreover, $\mathcal{S}_M \models \mathcal{S}_P$.

Note that refinement is associative: if $P_1 \succcurlyeq_{abs_1} P_2$ and $P_2 \succcurlyeq_{abs_2} P_3$, then $P_1 \succcurlyeq_{abs_1 \circ abs_2} P_3$. Thus, it is possible to gradually add information during model generation with several refinement steps $P_0 \succcurlyeq P_1 \succcurlyeq P_2 \succcurlyeq \dots \succcurlyeq M$ to arrive at a concrete model M .

4.3 First-order Logic Predicates for Queries Over Graph Models

The formal definitions of metamodel and instance model enable the formulation of first-order logic (FOL) predicates, which can be evaluated as *graph queries* over the logic structure of an instance model. Informally, base predicates check either for equality or for the existence of certain objects and references of a respective type (predicate) in the underlying runtime model. Then complex predicates are derived by traditional FOL connectives (e.g., not, exists, forall, and, or).

Definition 4.15. A *first-order logic predicate* (or *query*) φ , where v_1, \dots, v_n denote free variables (not appearing in any quantifiers) of φ can be evaluated over a partial model P along a *variable binding* $Z: \{v_1, \dots, v_n\} \rightarrow O_P$ (denoted as $\llbracket \varphi \rrbracket_Z^P$) to return 1, 0 or $1/2$ as shown in Figure 6.

Because concrete models contain only 1 and 0 logic values, any predicate φ evaluates to either 1 or 0 in a concrete model. Hence, on concrete models, we can run queries and obtain their match sets without uncertainties.

Definition 4.16. In a concrete model, *predicate / query evaluation* aims to find a variable binding $Z: \{v_1, \dots, v_n\} \rightarrow O_M$ for a predicate φ that maps all free variables of the predicate to objects of M such that the predicate evaluates to *true*, i.e., $\llbracket \varphi \rrbracket_Z^M = 1$.

Definition 4.17. The *match set* of a query predicate φ with free variables v_1, \dots, v_n is the set $Matches(M, \varphi) = \{Z: \{v_1, \dots, v_n\} \rightarrow O_M \mid \llbracket \varphi \rrbracket_Z^M = 1\}$. One element in this set is called a *match*, while $M\#\varphi = |Matches(M, \varphi)|$ denotes the size of the match set.

Note that in our context, a match of a query will typically represent a violation of a well-formedness constraint of the domain or a hazardous situation with respect to a safety property.

Example 4.18. Consider the graph query φ_{CT} for the “close trains” hazard formalized as a FOL expression φ_{CT} in Figure 3(c). In the concrete model M in Figure 2(b), $\llbracket \varphi_{CT} \rrbracket_{s \mapsto s_2, e \mapsto s_3}^M = 1$. In the partial model P in Figure 5, $\llbracket \varphi_{CT} \rrbracket_{s \mapsto s_2, e \mapsto s_3}^P = 0$ due to the uncertain existence of the tr_{new} multi-object. In M , φ_{CT} has two matches $\text{Matches}(M, \varphi_{CT}) = \{\{s \mapsto s_2, e \mapsto s_3\}, \{s \mapsto s_2, e \mapsto s_3\}\}$. Therefore, $M\#\varphi_{CT} = 2$.

4.4 Well-formedness and Scope Constraints

In order to make static analysis of data-driven monitors more precise, we can add additional domain-specific information into models as constraints to exclude impossible or irrelevant runtime snapshots from consideration.

A domain metamodel is frequently complemented in practice with *well-formedness constraints* to restrict the possible relationships between domain concepts. The constraints, such as **type hierarchy**, **type compliance**, **multiplicity**, **inverse relation** and containment hierarchy constraints, can be captured by FOL predicates [55].

Additionally, numerical *scope constraints* restrict the sizes of models to conform with allocation requirement in monitor programs and guide the analysis toward models that are relevant in practice (e.g., the size of the model and the ratios between the number of objects of given types match realistic scenarios).

In this work, we are interested in the timing analysis of monitors that take well-formed models conforming to numerical constraints as their input. Therefore, we will require model to conform to *theories* formed by graph predicates.

Definition 4.19. A *theory* over a signature $\langle \Sigma, \alpha \rangle$ is a pair $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$, where

- $\Phi = \{\varphi_1, \dots, \varphi_{|\Phi|}\}$ is a finite set of graph predicates over $\langle \Sigma, \alpha \rangle$; and
- $\mathbf{r}: \Phi \rightarrow \mathcal{X}$ maps graph predicates to linear equation variables.

Definition 4.20. A concrete model $M = \langle O_M, I_M, S_M \rangle$ is *compatible* with the theory $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ (written as $M \vDash \mathcal{T}$) if $S_M \vDash \mathbf{r}(\varphi_i) = M\#\varphi_i$ for all $\varphi_i \in \Phi$.

Theories, along with refinement, enable to constrain the number of graph predicate matches via partial models. In particular, if φ_i is an *error predicate* that should never match well-formed models, the presence of exactly 0 matches should be enforced by the model scope S_P .

Proposition 4.1 ([45]). Let $P = \langle O_P, I_P, S_P \rangle$ be a partial model, $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ be a theory and $\varphi_i \in \Phi$ be a graph predicate. If $S_P \vDash \mathbf{r}(\varphi_i) \leq U$ (resp. $S_P \vDash \mathbf{r}(\varphi_i) \geq L$), then $M\#\varphi_i \leq U$ (resp. $M\#\varphi_i \geq L$) holds for all concrete models $P \succcurlyeq M$ compatible with \mathcal{T} , i.e., M satisfies the upper (resp. lower) bound imposed on the number of φ_i matches.

Example 4.21. Consider the FOL predicates

$$\varphi_{\text{Train}}(v_1) = \text{Train}(v_1), \quad \varphi_{\text{connectedTo}}(v_1, v_2) = \text{connectedTo}(v_1, v_2) \wedge \neg \text{connectedTo}(v_2, v_1),$$

and the theory $\mathcal{T} = \langle \{\varphi_{\text{Train}}, \varphi_{\text{connectedTo}}\}, \mathbf{r} \rangle$, where $\mathbf{r}(\varphi_{\text{Train}}) = \mathbf{x}_1$ and $\mathbf{r}(\varphi_{\text{connectedTo}}) = \mathbf{x}_2$.

The predicate φ_{Train} selects all `Train` instances. Thus, the linear inequality $(\mathbf{x}_1 \leq 3) \in S_P$ in the partial model P in Figure 5 corresponds to the scope constraint that there should be no more than 3 `Train` instances in the model.

The predicate $\varphi_{\text{connectedTo}}$ selects `connectedTo` links that do not have a corresponding link in the reverse direction. Since railway tracks can be traversed in both directions, we enforce a symmetric `connectedTo` relation by a well-formedness constraint encoded as $(\mathbf{x}_2 = 0) \in S_P$ in P .

The concrete model M in Figure 2(b) conforms to the theory \mathcal{T} . As shown in Proposition 4.1, M obeys the scope and well-formedness constraints prescribed in S_P , since $P \succcurlyeq M$.

4.5 Model Generation

Automated synthesis of domain-specific graph models has been actively researched in the field of model-based software engineering [7, 23, 55]. Hereby, we revisit some core concepts.

Definition 4.22. A model generation task is of the form $\max \sum_{\mathbf{x}_j \in \mathcal{X}} c_j \cdot \mathbf{x}_j$ subject to $\langle P_{\text{init}}, \mathcal{T} \rangle$, where

- $\Sigma = \{C_1, \dots, C_{m_C}, R_1, \dots, R_{m_R}, \varepsilon, \sim\}$ is a metamodel;
- the partial model $P_{\text{init}} = \langle O_{P_{\text{init}}}, \mathcal{I}_{P_{\text{init}}}, \mathcal{S}_{P_{\text{init}}} \rangle$ over Σ is the *initial partial model*;
- $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ is a theory of well-formedness and scope constraints over Σ ; and
- $\sum_{\mathbf{x}_j \in \mathcal{X}} c_j \cdot \mathbf{x}_j$ is the *objective*.

Definition 4.23. The solutions of the model generation task are concrete models that are refinements of the initial partial model and are compatible with the theory:

$$\text{solutions}(\Sigma, P_{\text{init}}, \mathcal{T}) = \{M \mid M \text{ is a concrete instance of the metamodel } \Sigma, P_{\text{init}} \succcurlyeq M, M \models \mathcal{T}\}.$$

According to Proposition 4.1, such concrete models M satisfy the numerical constraints given in scope $\mathcal{S}_{P_{\text{init}}}$ that restrict the number of matches $M \# \varphi_i$ of graph predicates $\varphi_i \in \Phi$. If no structural information is available about the sought models, we may set $O_{P_{\text{init}}} = \{\text{new}\}$ and $\mathcal{I}_{P_{\text{init}}}(C_i)(\text{new}) = \mathcal{I}_{P_{\text{init}}}(R_j)(\text{new}, \text{new}) = \mathcal{I}_{P_{\text{init}}}(\varepsilon)(\text{new}) = \mathcal{I}_{P_{\text{init}}}(\sim)(\text{new}, \text{new}) = 1/2$ for all $C_i, R_j \in \Sigma$ to obtain the maximally uncertain initial partial model with a single multi-object. Otherwise, P_{init} contains the known parts of the model and multi-objects may serve as placeholders for objects to be added.

Definition 4.24. In the model generation task $\max \sum_{\mathbf{x}_j \in \mathcal{X}} c_j \cdot \mathbf{x}_j$ subject to $\langle P_{\text{init}}, \mathcal{T} \rangle$, the *cost function* $g(M) = \max \sum_{\mathbf{x}_j \in \mathcal{X}} c_j \cdot \mathbf{x}_j$ subject to \mathcal{S}_M of the model generation task can be computed by solving an ILP problem. A concrete model $M \in \text{solutions}(\Sigma, P_{\text{init}}, \mathcal{T})$ is an *optimal solution* of the task if $g(M') \leq g(M)$ for all $M' \in \text{solutions}(\Sigma, P_{\text{init}}, \mathcal{T})$, i.e., its cost is maximal.

Our work relies on the model generator presented in [45, 55] which was proved to be complete and sound in [61]. Informally, it is able to derive all concrete (instance) models in a domain (up to a designated size defined by the scopes) which satisfy the constraints by exploring a *state space* of possible partial models along refinements.

5 TIMING ANALYSIS OF QUERY-BASED MONITORS

Estimating the WCET of query-based monitors is a highly complex task which involves multiple classic challenges of timing analysis. The runtime model of the system is a continuously changing data structure that captures an up to date snapshot of the underlying running system. Hence, it is not sufficient to analyze execution time on a single input model, but all models possible at runtime must be considered.

However, the space of possible models is enormous. For example, in a metamodel with 3 reference types, there may be up to $2^{3 \cdot 25 \cdot 25} = 2^{1875}$ models with 25 objects. Thus, explicit enumeration of graph models is intractable, which necessitates the use of abstractions.

Another major challenge is that *query execution time is heavily data-dependent*, i.e., the same control flow of a query program may have substantially different run times based upon the structural characteristics of the underlying graph model. Assuming some constraints on model size (e.g., capped by available memory) and some general restrictions on model scope (e.g., there are more segments than trains in any real model), a key open challenge is *how to provide a model where the execution time of a particular query program is maximal*. In this work, we provide *witness models* that maximize an estimate of the execution time, which aids in WCET analysis and in identifying bottlenecks in query execution.

Moreover, a single model may be represented in memory in several isomorphic ways (Section 3.2.2). During the runtime evolution of the graph, a particular snapshot might be reached in

Table 2. WCET analysis approaches for data-driven runtime monitor programs

	Inputs	Outputs
	Low-level analysis inputs	
CL	HW description + Query program	WCET estimate
	Value analysis inputs	
VAL	HW description + Query program + Concrete model + Memory image	WCET estimate for single memory image
	Domain-specific high-level analysis inputs	
DS _M	HW description + Query program + Concrete model	WCET est. for single model
DS _Σ	HW description + Query program + Constraints	WCET est. for all valid models + Witness model
DS _P	HW description + Query program + Partial model + Constraints	WCET est. for all refinements + Witness model

any of its possible in-memory representations. Thus, WCET estimation even for a single concrete input model must tackle the dependency of execution paths on the data representations. As a single model of n objects has $n!$ possible in-memory representations even if we only consider inserting the objects into a single continuous linear array of n elements, explicit enumeration is again intractable.

5.1 Comparison of Timing Analysis Approaches

Table 2 illustrates the existing and proposed approaches of WCET analysis for query programs. *Classical* (CL) analysis is based on binary code of query program and the characteristics of the hardware platform, but does not consider structure and the well-formedness of the runtime models.

Value analysis (VAL) can derive more precise WCET estimates for executing a query on a single memory image (comprised of a single concrete model). However, it is unable to consider equivalent in-memory representations of the same concrete model (i.e., different parts of the model allocated to different spatial locations), or to cover all possible consistent concrete models, thus it is unsuitable for the analysis of data-driven monitors.

To alleviate this issue, we propose three *domain-specific* (DS) WCET analysis methods for data-driven monitors. We introduce the concept of *witness models*, which are consistent models that are *feasible inputs of the graph query program* and *maximize the WCET estimate for all models within the given scope*. They serve as representative data to calculate WCET for *any* model within the scope.

- First, we estimate WCET for a *single concrete model* (DS_M). The estimate is valid for all in-memory representations of a given concrete model M .
- In the next case, the set of possible runtime snapshots is specified with *metamodel* $\langle \Sigma, \alpha \rangle$ along with well-formedness and scope constraints (DS_Σ). This WCET estimate is valid for all possible runtime snapshots within the memory limits of the system, i.e., for all consistent instances of the metamodel up to the size specified by the scope constraints.
- Thirdly, an *initial partial model* P may specify the set of possible runtime snapshots (DS_P) including static (known and concrete) and dynamic (uncertain at design time) parts of the runtime model. The WCET estimate is valid for all possible refinements $P \succcurlyeq M$ of P .

Figure 7 sketches the *model space* of runtime graph models (represented with dots), i.e., the set of all input models. Possible changes made to a model at runtime (depicted as arrows) result in a new model. To obtain a safe and tight WCET estimate for query programs, we make some assumptions about realistic (and consistent) models captured in the form of a *model scope*. If an initial partial model P is provided, the analysis is further restricted to its (valid) refinements, thus inconsistent

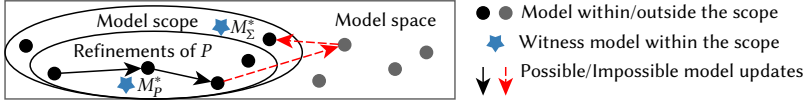


Fig. 7. Classification of query input models and model updates from the perspective of WCET analysis

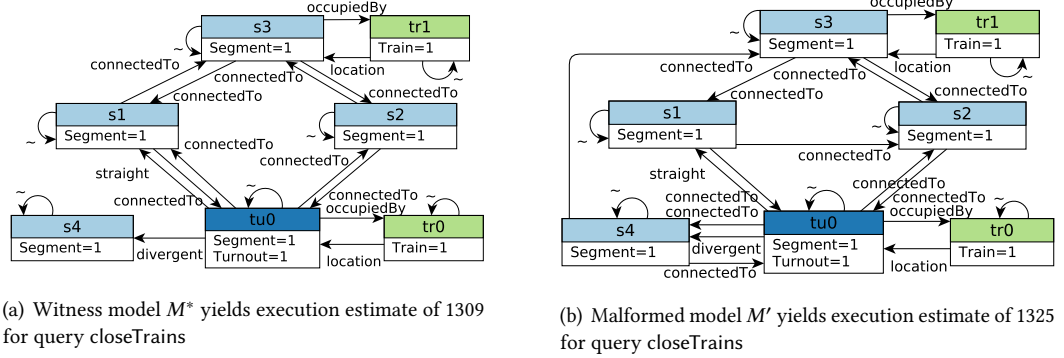


Fig. 8. Illustrating model generation problems for witness models

models are considered to be unrealistic. The witness model M_Σ^* for the consistent instances of the metamodel and M_P^* for the refinements of P are depicted as blue stars in Figure 7.

The witness models M_P^* may aid in iteratively refining the model scope. If the partial model corresponds to a situation that is impossible at runtime, it indicates that the model scope was specified in a too general way. We may exclude such situations by refining the partial model P [61]. However, care must be taken to avoid excluding feasible inputs and overfitting the WCET estimate. If the witness model is a feasible input, it may be inspected to study the characteristics and bottlenecks of the graph query program.

Example 5.1. Figure 8(a) shows the witness model M^* for the WCET of the `closeTrains` query for well-formed models with up to 7 objects in total (as model scope), out of which up to 2 are Train instances. The corresponding WCET estimate is 1309 systicks. The model M' in Figure 8(b) has the same number of elements, but with a higher execution time estimate of 1325 systicks. However, M' lies outside the model scope, because it is malformed due to non-symmetric `connectedTo` references (e.g., `s1` is `connectedTo` `s2` but not vice versa).

Classical (CL) WCET estimation techniques cannot exclude M' from the analysis and they would return a higher WCET estimate, while our novel DS_Σ technique can restrict the analysis to the model scope to return the correct estimate of 1309 systicks along with the witness M^* .

5.2 Architectural Overview

Figure 9 presents the high-level description of our design time tasks to obtain a WCET estimate in the DS_M , DS_Σ and DS_P scenarios. The high-level inputs of the process include the *query specification*, the target *hardware description*, and the *well-formedness and scope constraints* of the domain.

First, a *query plan* (A) is constructed from the query specification, based on which the *query program* (B) is generated according to Section 3.3. Our approach is complementary to IPET-based WCET estimators and leverages the results of high- and low-level analysis in form of the CFG (possibly after some loop unrolling) and the corresponding linear program (C).

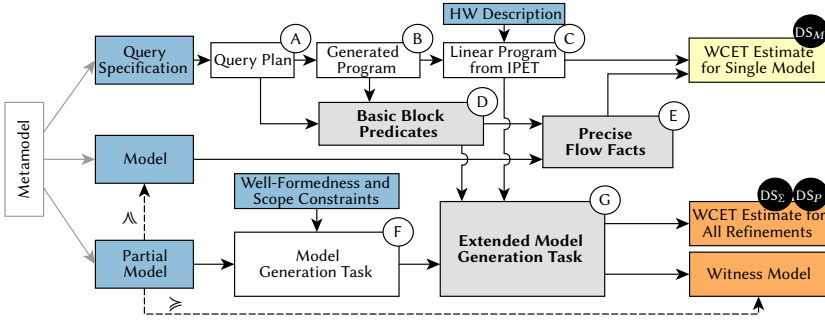


Fig. 9. Workflow of WCET estimation for query-based monitors

In case of WCET estimation for a concrete model M (DS_M), M is also provided as an input. Based on the query plan and the generated monitor code, *basic block predicates* (D) are derived, whose matches in the concrete model M correspond to executions of basic blocks in the monitor program. We leverage these matches to construct *precise flow facts* (E) for IPET analysis in Section 5.3. The resulting flow facts and WCET estimate consider all possible in-memory representations of M .

For WCET estimation for any valid instance of a metamodel $\langle \Sigma, \alpha \rangle$ (DS_Σ), the initial partial model P_{init} is constructed according to Section 4.5. When a partial model P is already provided as input (DS_P), it replaces P_{init} as the initial partial model. Hence, along with the \mathcal{T} of well-formedness and scope constraints, we obtain a *model generation task* (F), whose solutions are partial models within the analyzed model scope. We incorporate the basic block predicates (D) into an extended theory \mathcal{T}' in Section 5.4, which forms an *extended model generation task* (G) for witness model generation along with the linear program (C). *Witness models* are systematically generated using a graph solver [45, 55] as solutions of such tasks along refinements $P_{\text{init}} \geq M$ of the initial partial model. The cost associated with the witness model M^* , which is a solution of the IPET linear program (C) extended with domain-specific flow facts (E), is a safe and tight WCET estimate.

5.3 Approximating Execution Time with Graph Predicates

To derive precise flow facts for WCET analysis of a graph query program with concrete input model M and characterize its data-dependent execution time, we construct a *basic block predicate* ψ_{bb} for each basic block $bb \in BB$ of the query program. Free variables of ψ_{bb} correspond to program variables (bound by for loops). Due to the structure of the code generated from the query plans (Section 3.3.3), each execution of bb corresponds to a match $Z \in \text{Matches}(M, \psi_{bb})$ of ψ_{bb} in M .

For loop headers, we construct an additional $\psi'_{bb_{\ell,h}}$ where the matches of $\psi'_{bb_{\ell,h}}$ represent executions of the loop ℓ where the loop condition holds, while the matches of $\psi_{bb_{\ell,h}}$ correspond to the executions where the loop exits.

Algorithm 2 takes a data-driven monitor program generated from a graph query and constructs the basic block predicates. In addition to the set of basic blocks BB , the algorithm requires traceability information *lineTrace* (that connects basic blocks to source code lines) and *planTrace* (that connects source code lines to **extend** and **check** constraints in the query plan). The *lineTrace* is extracted from the IPET analysis tool (based on debug information in the compiled executable), while *planTrace* is the output of the query code generator.

As state-of-the-art WCET analysis tools [3] do not recommend analyzing programs compiled with advanced optimizations, we did not assess programs that use optimization. Therefore, source line, as well as **extend** and **check** constraint information in *lineTrace* remains valid after compilation.

Algorithm 2: Basic block predicate construction

```

1 Function DerivePredicates( $BB, lineTrace, planTrace$ ) is
2    $\Psi = \emptyset$ 
3   for  $bb \in BB$  do
4      $\psi_{bb} = 1$ 
5     Let  $ln_b - ln_e$  be the source lines associated with  $bb$  in
        $lineTrace$ 
6     for if and for statements  $st$  containing  $ln_b - ln_e$  do
7        $\psi_{bb} = \psi_{bb} \wedge \text{StatementToLogic}(st, planTrace)$ 
8     Let  $v_1, \dots, v_m$  be the free variables of  $\psi_{bb}$ 
9      $\Psi = \Psi \cup \{\psi_{bb}\}$ 
10    if  $bb$  is the header of the loop  $\ell$  then
11       $\psi'_{bb} = \psi_{bb} \wedge \text{StatementToLogic}(\ell, planTrace)$ 
12      Let  $v_1, \dots, v_{m+1}$  be the free variables of  $\psi'_{bb}$ 
13       $\Psi = \Psi \cup \{\psi'_{bb}\}$ 
14  return  $\Psi$ 

```

Algorithm 3: Translate search plan to logic

```

1 Function StatementToLogic( $st, planTrace$ ) is
2   Determine the constraint implemented by  $st$  from
        $planTrace$ 
3   if  $st$  implements extend  $\exists v_i : C(v_i)$  then
4     return  $C(v_i)$ 
5   else if  $st$  implements extend  $\exists v_j : R(v_i, v_j)$  then
6     return  $R(v_i, v_j)$ 
7   else if  $st$  implements check  $C(v_i)$  then
8     return  $C(v_i)$ 
9   else if  $st$  implements check  $R(v_i, v_j)$  then
10    return  $R(v_i, v_j)$ 
11  else if  $st$  implements check  $v_i = v_j$  then
12    return  $v_i = v_j$ 
13  else if  $st$  implements check  $\neg C(v_i)$  then
14    return  $\neg C(v_i)$ 
15  else if  $st$  implements check  $\neg R(v_i, v_j)$  then
16    return  $\neg R(v_i, v_j)$ 
17  else if  $st$  implements check  $\neg(v_i = v_j)$  then
18    return  $\neg(v_i = v_j)$ 

```

Algorithm 4: Precise flow fact construction for a single concrete model M

```

1 Function PreciseFlowFacts( $BB, lineTrace, planTrace, CFG = \langle V, E, s, t, w, tr \rangle, \bar{f}, M$ ) is
2    $S_{flow} = \emptyset, \Psi = \text{DerivePredicates}(BB, lineTrace, planTrace)$ 
3   for  $bb \in BB$  do
4     if  $bb$  is a loop header then  $S_{flow} = S_{flow} \cup \{\sum_{e=(n_1, n_2) \in E, tr(n_1)=bb} \bar{f}(e) = M\#\psi_{bb} + M\#\psi'_{bb}\}$ ;
5     else  $S_{flow} = S_{flow} \cup \{\sum_{e=(n_1, n_2) \in E, tr(n_1)=bb} \bar{f}(e) = M\#\psi_{bb}\}$ ;
6  return  $S_{flow}$ 

```

However, the following algorithms can be extended to support compiler optimizations, as long as a compiled basic block still corresponds to a single constraint and the control flow remains structured (comprised on loops and conditionals).

In line 4, ψ_{bb} is initialized to true, which has a single (trivial) match in any model to reflect that blocks not implementing any query plan constraints will be executed exactly once. Then, in line 5, we traverse $lineTrace$ to extract the source lines corresponding to bb . The loop in lines 6–7 processes all **if** and **for** statements enclosing the source lines for bb . As a result, ψ_{bb} becomes the conjunction of atomic predicates, which correspond to the query plan constraints implemented by the processed statements. Lastly, in lines 10–13, if bb is a loop header, we also add the atomic predicate corresponding to the loop itself to obtain ψ'_{bb} , which characterizes executions when the loop condition holds.

Algorithm 3 implements translation of **for** and **if** statements to atomic logical predicates. The algorithm traverses $planTrace$ to process the corresponding query plan constraint. For **extend** constraints (usually associated with **for** loops), the existential quantifier \exists is removed so that the constraint introduces a new free variable to ψ_{bb} . **Check** constraints (associated with **if** statements) are returned as-is, because all their variables are already introduced by some enclosing **extend** operation. Thus, for a basic block bb enclosed by m statements with **extend** constraints will have ψ_{bb} with free variables v_1, \dots, v_m . If bb is a loop header, ψ'_{bb} has an additional v_{m+1} free variable.

Example 5.2. Listing 4 shows the generated source code of the `closeTrains` graph query, while Figure 4 show the corresponding CFG (without any loop unrolling). The `lineTrace` information is depicted as line numbers next to the CFG nodes, and comments above the control structures contain `planTrace`. The basic block bb_9 corresponds to the loop header in line 17. Collecting the

Algorithm 5: Witness generation task construction for a partial model P

```

1 Function WitnessGenerationProblem( $BB, lineTrace, planTrace, CFG = \langle V, E, s, t, w, tr \rangle, P = \langle OP, IP, SP \rangle, \mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ ) is
2   Construct the IPET analysis  $\max \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot \mathbf{x}_i$  subject to  $\mathcal{S}_{IPET}$  for  $CFG$  with  $\mathbf{f}: E \rightarrow \mathcal{X}$  (w.l.o.g.  $\text{Ran } \mathbf{f} \cap \text{Ran } \mathbf{r} = \emptyset$ )
3    $\mathbf{r}' = \mathbf{r}, \mathcal{S}_{merge} = \emptyset, \Psi = \text{DerivePredicates}(BB, lineTrace, planTrace)$ 
4   for  $bb \in BB$  do
5     if  $bb$  is a loop header then
6       Let  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$  be two fresh variables not yet appearing in  $\text{Ran } \mathbf{f} \cup \text{Ran } \mathbf{r}'$ 
7        $\mathbf{r}' = \mathbf{r}' \cup \{\psi_{bb} \mapsto \mathbf{x}, \psi'_{bb} \mapsto \mathbf{x}'\}, \mathcal{S}_{merge} = \mathcal{S}_{merge} \cup \{\mathbf{x} + \mathbf{x}' - \sum_{e=(n_1, n_2) \in E, tr(n_1)=bb} \mathbf{f}(e) = 0\}$ 
8     else
9       Let  $\mathbf{x} \in \mathcal{X}$  be a fresh variable not yet appearing in  $\text{Ran } \mathbf{f} \cup \text{Ran } \mathbf{r}'$ 
10       $\mathbf{r}' = \mathbf{r}' \cup \{\psi_{bb} \mapsto \mathbf{x}\}, \mathcal{S}_{merge} = \mathcal{S}_{merge} \cup \{\mathbf{x} - \sum_{e=(n_1, n_2) \in E, tr(n_1)=bb} \mathbf{f}(e) = 0\}$ 
11   $\mathcal{S}_{P'} = \mathcal{S}_P \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge}, P' = \langle OP, IP, SP' \rangle, \Phi' = \Phi \cup \Psi, \mathcal{T}' = \langle \Phi', \mathbf{r}' \rangle$ 
12  return  $\max \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot \mathbf{x}_i$  subject to  $\langle P', \mathcal{T}' \rangle$ 

```

query plan constraints from the control structures enclosing line 17 with Algorithm 2, we find that

$$\begin{aligned} \psi_{bb_9} &= \text{Train}(t) \wedge \text{location}(t, s) \wedge \text{connectedTo}(s, m) \\ \psi'_{bb_9} &= \text{Train}(t) \wedge \text{location}(t, s) \wedge \text{connectedTo}(s, m) \wedge \text{connectedTo}(m, e). \end{aligned}$$

In the concrete model M^* in Figure 8(a), executions of bb_9 are represented by the 4 matches $\text{Matches}(M, \psi_{bb_9}) = \{\{t \mapsto \text{tr}0, s \mapsto \text{tu}0, m \mapsto s1\}, \{t \mapsto \text{tr}0, s \mapsto \text{tu}0, m \mapsto s2\}, \{t \mapsto \text{tr}1, s \mapsto s3, m \mapsto s1\}, \{t \mapsto \text{tr}0, s \mapsto s3, m \mapsto s2\}\}$ of ψ_{bb_9} , as well as the 8 matches of ψ'_{bb_9} obtained by extending each ψ_{bb_9} match by the two possible segments connectedTo the value of m as the value of the variable e . Each match describes the values of the program variables when entering bb_9 .

Algorithm 4 constructs precise domain-specific flow facts for a concrete model M . In addition to the basic blocks BB and the traceability information, the algorithm reads the control flow graph $CFG = \langle V, E, s, t, w, tr \rangle$ and the function $\mathbf{f}: E \rightarrow \mathcal{X}$ associating CFG edges with linear equation variables. Line 2 initializes the empty system of linear equations \mathcal{S}_{flow} and constructs the basic block predicates Ψ . Leveraging the CFG traceability function $tr: E \rightarrow BB$, expressions $\sum_{e=(n_1, n_2) \in E, tr(n_1)=bb} \mathbf{f}(e)$ are built, which represent the number of times a basic block bb is executed (Section 4.1). For a loop header, this number is equal to the number of ψ_{bb} and ψ'_{bb} matches in M (line 4), while for other blocks, only matches of ψ_{bb} are counted (line 5).

The resulting set of linear equations \mathcal{S}_{flow} serve as flow facts in IPET analysis. More precisely, by incorporating \mathcal{S}_{flow} into the analysis, we may obtain a safe and tight estimate for the execution time of a graph query program on the concrete model M .

Proposition 5.1. Let τ be the execution time of the query program q on the concrete model M ,

$$CL = \max \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot \mathbf{x}_i \text{ subject to } \mathcal{S}_{IPET}, \quad DS_M = \max \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot \mathbf{x}_i \text{ subject to } \mathcal{S}_{IPET} \cup \mathcal{S}_{flow},$$

where CL is the classical IPET estimated obtained from q , and DS_M is the domain-specific estimate with flow facts derived from M using Algorithm 4. Then $\tau \leq DS_M \leq CL$. (See proof in Appendix A.)

5.4 Witness Generation of Worst-Case Execution Time

To estimate the WCET of some query program q over a set of models, we specify the *model scope* of interest as the $\text{solutions}(P, \mathcal{T})$ of a model generation task. We construct an extended model generation problem $\max \sum_{\mathbf{x}_i \in \mathcal{X}} c_i \cdot \mathbf{x}_i$ subject to $\langle P', \mathcal{T}' \rangle$ in Algorithm 5, where $P \succcurlyeq P'$ extends P with the results of IPET analysis and \mathcal{T}' incorporates the basic block predicates obtained from q . We use the notation $\text{Ran } \mathbf{f}$ to denote the *range* (possible values) of the function \mathbf{f} .

In line 2, the algorithm builds an IPET integer program based on the provided *CFG*. Then, in line 3, we invoke Algorithm 2 to obtain the basic block predicates Ψ . The extended theory $\mathcal{T} = \langle \Phi', \mathbf{r}' \rangle$ is comprised of the predicates Φ from the original theory $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ and the basic block predicates Ψ . Function \mathbf{r}' extends \mathbf{r} by assigning a variable $\mathbf{r}'(\psi_{bb})$ for each basic block predicate ψ_{bb} and a variable $\mathbf{r}'(\psi'_{bb})$ for each loop header predicate ψ'_{bb} . By Definition 4.20, in any concrete model M compatible with \mathcal{T}' , we have $\mathcal{S}_M \models \mathbf{r}'(\psi_{bb}) = M\#\psi_{bb}$ for each basic block $bb \in BB$ and $\mathcal{S}_M \models \mathbf{r}'(\psi'_{bb}) = M\#\psi'_{bb}$ for each loop header in addition to any constraints prescribed by the original theory \mathcal{T} . Without loss of generality, we assume that newly assigned variables are fresh (i.e., do not appear in either the original partial model scope \mathcal{S}_P or in the IPET analysis \mathcal{S}_{IPET}) and the variables of \mathcal{S}_{IPET} are distinct from \mathcal{S}_P .

Lines 7 and 10 of Algorithm 5 correspond to lines 4 and 5 of Algorithm 4. However, we use the associated variables $\mathbf{r}'(\psi_{bb})$ and $\mathbf{r}'(\psi'_{bb})$ instead of the raw match counts $M\#\psi_{bb}$ and $M\#\psi'_{bb}$, so that the linear equations \mathcal{S}_{merge} hold for any concrete model $M \models \mathcal{T}$. Lastly, in line 11, we assemble the extended model generation task and the partial model $P' = \langle \mathcal{O}_P, \mathcal{I}_P, \mathcal{S}_{P'} \rangle$, where $\mathcal{S}_{P'}$ contains the original scope \mathcal{S}_P , the IPET linear equations \mathcal{S}_{IPET} , and the \mathcal{S}_{merge} equations that merge the original model generation task and the IPET analysis together. The objective on the extended model generation task coincides with that of the IPET linear program.

The resulting extended model generation task provides a safe and tight estimate for the query program WCET with theory \mathcal{T} .

Proposition 5.2 (Safety and tightness). Let $\tau(M)$ be the execution time of a query program q on a concrete model M , P be partial model, \mathcal{T} be a theory, and

$$CL = \max_{\mathbf{x}_i \in \mathcal{X}} \sum c_i \cdot \mathbf{x}_i \text{ subject to } \mathcal{S}_{IPET}, \quad DS_P = \max_{\mathbf{x}_i \in \mathcal{X}} \sum c_i \cdot \mathbf{x}_i \text{ subject to } \langle P', \mathcal{T}' \rangle,$$

where CL is the classical IPET estimated obtained from q , and DS_P is the domain-specific estimate based on the extended graph generation problem form Algorithm 5. Then $\tau(M) \leq DS_P \leq CL$ for all $M \in solutions(P, \mathcal{T})$. (See proof of propositions in Appendix A.)

Optimal solutions are *witness models*, which maximize the domain-specific of WCET for concrete refinements of the input partial model P compatible with the theory \mathcal{T} . As the witness model M^* is in the model scope, it is a feasible (as opposed to spurious) input of the query program.

Proposition 5.3 (Witness model). Let $DS_M(M)$ be the domain-specific WCET estimate of a query program q obtained by Algorithm 4 for a concrete model M , DS_P be the domain-specific WCET estimate of q for a partial model P and theory \mathcal{T} by Algorithm 5, and M^* be the witness model for the WCET of q , i.e., the optimal solution of DS_P . Then $M^* \in solutions(P, \mathcal{T})$ and $DS_M(M) \leq DS_M(M^*) = DS_P$ for all $M \in solutions(P, \mathcal{T})$.

Moreover, refinements of the partial model $P \succcurlyeq Q$ may be used to tighten the WCET estimate by reducing the model scope under discussion.

Proposition 5.4 (Tightening by refinement). Let $DS_P(P, \mathcal{T})$ denote the domain-specific WCET estimate of a query program q for a partial model P and theory \mathcal{T} obtained by Algorithm 5 and $P \succcurlyeq Q$. Then $DS_P(Q, \mathcal{T}) \leq DS_P(P, \mathcal{T})$. In particular, if $P = P_{init}$ is the initial partial model for a metamodel $\langle \Sigma, \alpha \rangle$ from Section 4.5, then we may see that the WCET estimate for any partial model conforming to the metamodel is at least as tight as the DS_Σ estimate for the metamodel.

Example 5.3. Figure 10 shows a simplified execution of the graph generator. Suppose that Algorithm 5 has output an extended graph generation task

$$\max 250 \cdot \mathbf{x}_2 \text{ subject to } \langle P', \mathcal{T}' \rangle,$$

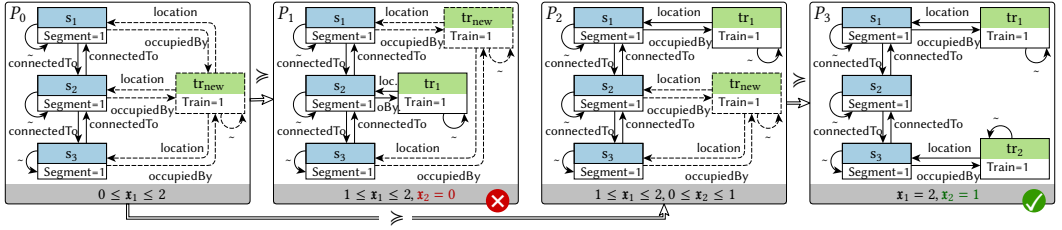


Fig. 10. Execution of the graph generation task $\max 250 \cdot \mathbf{x}_2$ subject to $\langle P', \mathcal{T}' \rangle$, where $P' = P_0$, $\mathcal{T}' = \langle \Phi', \mathbf{r}' \rangle$, $\Phi' = \{\varphi_{\text{Train}}, \varphi_{\text{CT}}\}$, $\varphi_{\text{Train}} = \text{Train}(v_1)$, φ_{CT} is as in Figure 3(c), $\mathbf{r}'(\varphi_{\text{Train}}) = \mathbf{x}_1$, and $\mathbf{r}'(\varphi_{\text{CT}}) = \mathbf{x}_2$.

where $P' = P_0$ as shown in Figure 10. The extended theory is $\mathcal{T}' = \langle \Phi', \mathbf{r}' \rangle$, there $\Phi' = \{\varphi_{\text{Train}}, \varphi_{\text{CT}}\}$, $\varphi_{\text{Train}} = \text{Train}(v_1)$, φ_{CT} is the predicate from the `closeTrains` query from Figure 3(c), $\mathbf{r}'(\varphi_{\text{Train}}) = \mathbf{x}_1$, and $\mathbf{r}'(\varphi_{\text{CT}}) = \mathbf{x}_2$.

The inequalities $[0 \leq \mathbf{x}_1 \leq 2]$ in conjunction with the theory \mathcal{T}' prescribe a *type scope* of between 0 and 2 `Train` instances. The objective function $g(M) = \max_{k \in S_M} 250 \cdot k(\mathbf{x}_2)$ corresponds to each match of φ_{CT} (i.e., `closeTrains`) taking 250 clock ticks to calculate.

Two non-isomorphic partial models P_1, P_2 can be obtained from $P_0 = P'$. In P_1 , the generator placed a train tr_1 on the middle segment s_2 of the track. Therefore, the multi-object tr_{new} represents at most one additional train ($1 \leq \mathbf{x}_1 \leq 2$). The φ_{CT} predicate cannot match ($\mathbf{x}_2 = 0$), since it is impossible to place a new train on both s_1 and s_3 . Any possible concrete refinement $P_1 \succcurlyeq M$ of P_1 has an objective value $g(M) = 250 \cdot 0 = 0$.

If we place a new train on s_1 , we obtain P_2 . Placing a train on s_3 results in an isomorphic model, so it is sufficient to only consider P_2 instead. Here, there can be 0 or 1 matches φ_{CT} ($0 \leq \mathbf{x}_2 \leq 1$). Indeed, if we place an additional train on s_3 , we obtain the concrete model $M^* = P_3$ with a single match of φ_{CT} ($\mathbf{x}_2 = 1$). The corresponding objective value is $g(M^*) = 250 \cdot 1 = 250$. No larger objective value is possible by any concrete refinement of P_1 , since $S_{P_1} \models \mathbf{x}_2 = 0$. Hence we may discard P_1 along with its potential refinements, and output M^* as the witness model obtained as the optimal solution of model generation task.

Assuming that the (simplified) objective function g is the domain-specific WCET estimate of the query program, $g(M^*) = 250$ is our WCET estimate, which execution time bound is expected to be reached (according to the low-level IPET analysis) when executing the query program over the witness model M^* as input.

6 EVALUATION

We conducted experiments to address the following research questions related to the WCET of query programs. For each research question, we investigate the scenario (a) DS_Σ , where only the metamodel and the relevant well-formedness and scope constraints are known; and (b) DS_P , when an initial partial model (describing a track layout but not its runtime state) is also provided.

RQ1 How difficult is it to find witness models?

RQ2 How safe and tight are WCET estimates w.r.t. existing approaches and real execution times?

RQ3 How does query program complexity impact the overestimation of computed WCET bounds?

RQ1 aims at determining whether our model generation based approach can find the witness model in practical time and whether it constitutes an improvement over random search. The rest of the experiments study the quality of WCET bounds, which is a key factor in the applicability of our approach. In particular, RQ2 attempts to compare our computed WCET estimates with the state of

the art in challenging settings with partial runtime information, while RQ3 presents increasingly challenging query programs to our approach.

6.1 Evaluation Overview and Setup

6.1.1 Queries. To address these research questions, we use graph queries from the domain of the MoDeS3 CPS demonstrator [63]. This demonstrator uses high-level runtime monitoring rules captured as graph queries, and showcases synthesized monitoring programs executing these queries over the runtime graph model of the underlying running system. Our experiments focus only on query evaluation, and updates to the runtime model are out of scope for the current paper. Therefore, we ran the query programs on various snapshots of runtime graph models. We evaluated the following queries adopted from [8]:

- **Close trains (ct):** This is the query introduced in the running example of Section 3.3.
- **End of siding (eos):** This query finds trains that are dangerously close (one segment distance) to an end of the track.
- **Misaligned turnout (mt):** Pairs of trains and turnouts are the objectives of this query, where the train would derail if it reached the turnout because it is switched in a different direction.
- **Train locations (tl):** A simple query to find pairs of trains and segments that describe the locations of each train.

The calculation of query search plans is out of scope of the current paper, but they were created and optimized based on the typical model statistics of runtime model snapshots in the MoDeS3 system. For example, the search plan presented in Table 1 is the one used by the program executing the query Close trains.

6.1.2 WCET algorithms and WCET tools. To compare the results produced by our WCET estimation approaches DS_{Σ} and DS_P with estimates produced by other tools, we used the commercial aiT [21] (version 20.10i) and the open-source OTAWA [3] (version V1.2.0) tools. For aiT, we used a *high precision* configuration with pipeline-level analysis and full (up to the determined loop bound) loop unrolling, as well as a *low precision* configuration with only basic block-level analysis and no unrolling. To incorporate the results of low-level analysis into DS_{Σ} and DS_P , we extracted the IPET linear equations from the low-level configuration of aiT manually, as no facility was available for automatic export or accessing the high precision system of linear equations directly. We also extracted the IPET linear equations from OTAWA, which have BB execution context information (paths of length two).

6.1.3 Graph models. In the following, we describe how we obtained a variety of models to assess the impact of models with different characteristics on query evaluation times.

Using the metamodel in the MoDeS3 case study, we generated *witness models* M_{Σ}^* for each monitoring query and for both low-level analyses (aiT, OTAWA) such that the query is estimated to have the longest possible execution time according the low-level analysis. For all of these models, we used the same model scope inspired by the railway domain: up to 20% of the objects can be Trains and up to 20% of the objects can be Turnouts. The rest of the objects are Segments; we capped the maximum number of objects at 25. The resulting models are syntactically valid and they can represent a realistic railway system thanks to the domain-specific well-formedness constraints.

To obtain a *realistic model* M_{real} , we manually captured a detailed runtime model snapshot of MoDeS3 that is similar to the one presented in Figure 2(b) with a total of 25 objects. Then, we removed all Train objects from M_{real} and unset all turnout directions, and used the resulting (partial) track layout P to find specific placements of trains and switching of turnouts such that the run times of the queries are maximized on the generated M_P^* witness models.



Fig. 11. Query execution times on fully random models and realistic models

To assess the execution times of the query programs on random models, we generated models conforming to the MoDeS3 metamodel with up to a total of 25 objects. Due to the large space of possible graph models, representative sampling from the model space is an open question [33, 54]. Nevertheless, we generated 250 models with the EMF random model generator¹ (**Rand**) with up to 5 Turnouts and up to 5 Trains, but none of them represents a railway setting that can occur because they all violated well-formedness constraints due to the completely random construction.

We also generated 250 models with the VIATRA Generator (**VG**) without an optimization objective, which satisfy all well-formedness and scope constraints used for generating witness models. However, the state exploration heuristics of the generator may lead to a biased sample.

6.1.4 Hardware setup. We use the Infineon Relax Lite Kit-V1 Board² to execute the query programs. This board has an XMC4500 F100-K1024 microcontroller and it is driven by a 120MHz system clock. This microcontroller is considered to be a mature industrial microcontroller and has an ARM Cortex-M4 core. For the present evaluation, the instruction cache on the device is not used as our primary focus is on the impact of domain-specific information about high-level program flow rather than microarchitectural effects.

The bare-metal query programs are compiled with GCC compiler for ARM version 7.2.1 with `-O0` and `-g3` flags in debug mode. These programs run on the microcontroller while no other tasks (e.g., interrupts) are running. We rely on the cycle counter feature of the Data Watchpoint and Trace Unit in the device to extract the execution times of each query using a debugger. The embedded code used for the experiments as well as compiler and other configurations are available online³.

6.2 Evaluation Results

6.2.1 Research Question 1 (a). We investigate if a witness model for a query can be obtained from simpler graph generation approaches, and we do this by measuring the execution times of queries

¹<https://github.com/atlanmod/mondo-atlzoo-benchmark>

²<http://www.infineon.com/xmc-dev>

³<https://imbur.github.io/cps-query/>

over various models. Our results are presented in Figure 11(a). The run times over models by **VG** is captured by the green boxes, while the orange ones show the run times over models by **Rand**. Each query was evaluated on the same two sets of models. Additionally, the respective query execution time over each witness model M_{Σ}^* is added to these figures for comparison, where M_{Σ}^* is the witness model generated using the objective function built from the low-level analysis results of aiT. Moreover, the run time over the hand-crafted M_{real} model is also presented.

Findings. For each consistent model considered, queries exhibited the longest observed execution times on their respective witness models. In fact, for two queries, eos and mt, the execution time on the witness is longer than the maximum measured execution time over any other consistent model, which highlights the importance of our witness model generation technique.

Maximum run times over models generated by **Rand** can be both higher and lower than on witness models. For example, query ct takes 2% shorter over a random model than over its witness model, but the random model does not represent a realistic railway. On the contrary, query eos takes at least 8% longer to complete on the witness model than on any model generated by **Rand**. *Therefore, computing safe and tight WCET estimates of queries which execute over well-formed models (1) is infeasible by collecting run times over random models, and (2) necessitates finding witness models by employing sophisticated model generation approaches.*

6.2.2 Research Question 1 (b). All possible refinements of P constitute a potentially large model space where finding M_p^* can be challenging and requires the graph solver to apply suitable abstractions for optimization. In our case, there are $3^5 \cdot \sum_{i=0}^5 \binom{20}{i} = 5\,273\,100$ models (Total row in Figure 11(b)) in the space of well-formed concrete refinements of the initial partial model P containing the (selected) track layout, because each of the 5 turnouts can be in three different states, and there can be up to 5 trains which must be located on different segments. Thus, explicit enumeration of all models is possible for such a track layout, although it is computationally expensive.

As described in Section 6.1.3, we used the graph generator to add trains to a model with an empty track layout such that the expected query run times are maximized. Figure 11(b) presents the number of states explored by the graph generator compared to the space of all refinements.

Findings. The number of states visited by the generator increases with the complexity of the query. For the most complex query ct, the generator was able to find the model with the highest estimated WCET after visiting 19% of the model space. For the least complex tl query, it visited only 144 states, which allowed the witness generation to finish almost instantly. In conclusion, the generator explores a fraction of the state space, making it more favorable than explicit enumeration.

6.2.3 Research Question 2. Our goal is to compare the computed WCETs obtained from different tools with our own techniques. For this RQ, we restrict our investigation to the scenario DS_{Σ} where only the metamodel and the well-formedness and scope constraints are known (i.e., case **(a)**), because only our technique but not the baseline tools (aiT, OTAWA) support processing a partial model as in DS_p (i.e., case **(b)**). Table 3 shows the WCET estimates for the 4 queries along with measured execution time (expressed in systicks) over the respective witness model.

Findings. In the case of ct, our WCET estimation approach produces estimates 14% tighter than the one by aiT (low precision analysis). It is also important to point out that even without context-sensitive BB timings, our low precision approach provides only 3% higher estimates than aiT's high precision mode, which indicates that it is able to automatically identify infeasible paths in the program based on high-level domain-specific information. For OTAWA, improvements of the WCET estimate achieved in two cases: ct has a 23%, while eos has a 2% tighter estimate. For the rest of the queries, the analysis yields the same results as aiT low precision mode or OTAWA.

Table 3. Query code complexity, measured execution time, and WCET estimates in systicks

Query	CC	Exec. time over M_{Σ}^*	DS_{Σ}		aiT		OTAWA
			w/aiT	w/OTAWA	low precision	high precision	
Close trains	7	2652	3133	3430	3563	3038	4210
End of siding	6	1395	1757	1820	1757	1477	1860
Misaligned t.	5	939	1097	1370	1097	987	1370
Train locations	3	489	592	695	592	507	695

Table 4. Query code complexity, measured execution time, and WCET estimates with a partial model

Query	CC	Exec. time over M_P^*	DS_P		aiT with partial memory image	
			w/aiT	w/OTAWA	low precision	high precision
Close trains	7	2544	3079	3338	3706	3091
End of siding	6	1275	1554	1636	1813	1523
Misaligned t.	5	969	1097	1370	1065	950
Train locations	3	504	592	695	586	506

Therefore, WCET estimates by DS_{Σ} were at least as tight as those obtained by low-level IPET analysis. Thus, domain-specific analysis can improve WCET estimates while simultaneously synthesizing witness models to study query program behavior. Conceptually, it would be possible to formulate more precise DS_{Σ} estimates by incorporating low-level analysis results from the high precision mode of aiT as shown in Section 5.4, but such equations cannot be obtained from aiT.

6.2.4 Research Question 3 (a). With this RQ, we look at the impact of query complexity on the computed WCET bounds, so that we can give recommendations on where our approach offers the greatest benefits. The execution times of queries over M_{Σ}^* in Table 3 provide a lower bound to the actual WCET (i.e., the longest possible execution time of the program over inputs which represent well-formed models in the model scope), while the CC columns shows query cyclomatic complexity. Since the actual WCET of the program is unknown (but it must lay between the measured execution time and the WCET estimates produced by the analyses), we use the measured execution time over witness models as the baseline when discussing overestimation in WCET estimates.

Findings. The biggest visible advantage of DS_{Σ} is in the case of the most complex query ct: the overestimation is 18% with BB timings from aiT, while the aiT low precision analysis computes a 34% higher value. In other cases, it produces the same result as aiT, with overestimates being between 16% (query mt) and 26% (query eos). We come to the same conclusion using BB timings from OTAWA, although these timings are slightly more conservative. The high precision analysis available in aiT is able to leverage the microarchitectural properties and thus provide the most precise estimates with the overestimation being 14% (observed for query ct). The overestimation increases with CC of the query code only in the case of high precision aiT analysis.

In general, DS_{Σ} computes a safe WCET bound and additionally provides a witness model. Moreover, it is able to discover additional infeasible paths the WCET estimate for the most complex query, thus provide a tighter estimate.

6.2.5 Research Question 3 (b). The goal of this RQ is to conclude if providing an initial graph model with elements of the graph model known upfront can lower the WCET, thus provide an execution time estimate for M_P^* lower than the one computed for M_{Σ}^* . The execution times of queries over M_P^* are presented in Table 4 similarly to Table 3. The estimates computed by DS_P are valid and safe for

any possible in-memory representation of the refinements of P (concrete models containing the designated track layout).

Additionally, Table 4 shows estimates from aiT computed by value analysis (VAL) over the initial track layout model provided to aiT as a *partial memory image*, where the unknown parts of P (locations of trains and the directions of turnouts) are left uninitialized. While WCET estimates based on this input are *not safe*, as they do not consider all possible in-memory representations of the refinements of P , we still added these numbers to Table 4, because they correspond to the most likely usage of existing WCET analysis tools with partial input.

Findings. Comparing DS_Σ and DS_P , we discover that providing an initial model tightens WCET estimates for the two more complex queries, but does not change the estimate for the two less complex ones. This is due to the nature of the initial model and query search plans. There is no arrangement of trains on the initial track layout such that the WCET estimate from DS_Σ for ct and eos run times is reached, whereas the evaluation of mt and tl depends less on the track layout.

Additional observations. For mt and tl, observed query run times on the witness models M_P^* were higher than on M_Σ^* , which can be attributed to the placement of data in memory as mentioned in Section 3.2.2. Nevertheless, they were still within the WCET estimates from both DS_Σ and DS_P .

Unexpectedly, the partially specified data yields higher WCET estimates by both analysis modes of aiT for the two more complex queries, ct and eos, when compared with the estimates in Table 3. Thus, for these queries, it is not possible to tighten WCET estimates for partial input data even with the caveat that all objects in the partial input are statically allocated. We have reported our observation to the developers of aiT at AbsInt GmbH, and they have confirmed that the discrepancy in the estimates is due to the differences in the placement of data in the two binaries. Note that these analyses solve a less general challenge compared to DS_P since they only consider one possible physical layout of the partial data in memory (and not all possible in-memory layouts), thus they are highlighted in gray in Table 4.

On the other hand, aiT high-precision mode produces a *lower estimate* based on the provided initial model than what is measured over M_P^* (highlighted in red in Table 4). The underlying reason for this is that the tool relies on the exact placement of data in memory rather than the abstract graph model the data encodes. Eventually, in M_P^* , the model objects were stored in the memory in different order which resulted in a higher runtime (see note about data placement in Section 3.2.2).

In general, providing partial input data allows for tightening WCET estimated for complex queries even in cases where value analysis with partial input data is not safe. For less complex queries, supplying a partial memory image can tighten the results of value analysis considerably, but requires committing to a specific in-memory representation of the partial data statically.

6.3 Threats to Validity

Internal validity. The current evaluation was performed on a device where the executing binary only included the query-based monitor, so we can assume that the measurements presented here precisely show the execution times of queries. Since the exact WCET of the program is unknown, we used the longest observed execution time to assess the overestimation of the computed WCETs. In reality, this overestimation might be lower than what is reported here, which would make our WCET estimates tighter than presented.

External validity. We carried out the evaluation using one specific hardware and compiler, thus the presented results may not generalize to other platforms. Furthermore, the presented approach is applicable to any query-based monitor generated with Algorithm 1. However, evaluation of the WCET estimation techniques using additional case-studies with query-based runtime monitors from different domains could further improve the confidence in the evaluation results.

7 CONCLUSIONS AND FUTURE WORK

In this paper, we presented a method to *provide safe and tight WCET bounds for runtime monitoring programs derived from graph queries* to enable their use in real-time systems. We provided a static WCET estimate by incorporating low-level analysis results from traditional IPET-based tools and high-level domain-specific constraints into the objective function of an advanced graph solver. In addition to a tight WCET estimate, the result also entails a witness graph model where the query-based monitoring program execution time is expected to be the longest.

We carried out extensive evaluation of our approach on an industry-grade hardware platform using a variety of graph models as inputs for query programs, and assessed the tightness of computed WCET by comparing it to the results produced by two different tools. We constructed witness models for highest estimated execution times of queries as well as random graph models as inputs for graph query programs as an attempt to showcase high execution times. While we have no formal guarantee that worst-case timing behavior is exhibited on witness models as inputs, in all our experiments, the longest execution times were always measured on such witness models.

In the short run, the proposed approach can be improved by passing the results of high-precision IPET analysis (including CFG unrolling and pipeline analysis) to the graph solver, while the evaluation of the approach should be done on a different hardware platforms as well. As a part of a long-term future research agenda, our approach could be extended to provide witness models with specific data placement in memory where the execution time equals to the WCET of the program.

ACKNOWLEDGMENTS

This work has been partially supported by NSERC RGPIN-04573-16 project and the MEDA scholarship program. The research reported in this paper and carried out at the BME has been supported by the NRDI Fund based on the charter of bolster issued by the NRDI Office under the auspices of the Ministry for Innovation and Technology. The authors would like to thank AbsInt GmbH for providing us a license for the aiT timing analyzer, and also would like to thank Martin Sicks for the extensive introduction to the tool. Similarly, we are thankful to Julien Forget and Clément Ballabriga for helping with using OTAWA.

REFERENCES

- [1] Raja Ben Abdesslem, Annibale Panichella, Shiva Nejati, Lionel C. Briand, and Thomas Stifter. 2018. Testing Autonomous Cars for Feature Interaction Failures Using Many-Objective Search. In *33rd ACM/IEEE International Conference on Automated Software Engineering*. 143–154.
- [2] Jaume Abella et al. 2015. WCET analysis methods: Pitfalls and challenges on their trustworthiness. *10th IEEE International Symposium on Industrial Embedded Systems - Proceedings (2015)*, 39–48.
- [3] Clément Ballabriga, Hugues Cassé, Christine Rochange, and Pascal Sainrat. 2010. OTAWA: An Open Toolbox for Adaptive WCET Analysis. In *LNCS*. Vol. 6399.
- [4] Clément Ballabriga, Julien Forget, and Giuseppe Lipari. 2017. Symbolic WCET computation. *ACM Trans. Embedded Comput. Syst.* 17, 2 (2017).
- [5] Ezio Bartocci et al. 2018. Specification-Based Monitoring of Cyber-Physical Systems: A Survey on Theory, Tools and Applications. In *Lectures on Runtime Verification*. 135–175.
- [6] Gordon S. Blair, Nelly Bencomo, and Robert B. France. 2009. Models@run.time. *IEEE Computer* 42, 10 (2009), 22–27.
- [7] Erwan Brottier, Franck Fleurey, Jim Steel, Benoit Baudry, and Yves Le Traon. 2006. Metamodel-based test generation for model transformations: an algorithm and a tool. In *2006 17th International Symposium on Software Reliability Engineering*. 85–94.
- [8] Márton Búr, Gábor Szilágyi, András Vörös, and Dániel Varró. 2018. Distributed graph queries for runtime monitoring of cyber-physical systems. In *LNCS*. Vol. 10802. 111–128.
- [9] Sven Burmester, Holger Giese, Martin Hirsch, and Daniela Schilling. 2004. Incremental design and formal verification with UML/RT in the FUJABA real-time tool suite. In *Proceedings of the International Workshop on Specification and Validation of UML Models for Real Time and Embedded Systems, SVERTS2004*. Citeseer.

- [10] Sven Burmester, Holger Giese, Andreas Seibel, and Matthias Tichy. 2005. Worst-case execution time optimization of story patterns for hard real-time systems. In *3rd International Fajaba Days*. 71–78.
- [11] Hugues Cassé and Pascal Sainrat. 2006. OTAWA, a framework for experimenting WCET computations. *3rd European Congress on Embedded Real-Time* January (2006), 1–8.
- [12] Kong-Rim Choi and Kyung-Chang Kim. 1996. T*-tree: a main memory database index structure for real time applications. In *3rd International Workshop on Real-Time Computing Systems and Applications*. 81–88.
- [13] Duc Hiep Chu and Joxan Jaffar. 2011. Symbolic simulation on complicated loops for WCET path analysis. In *Proc. of the 9th ACM International Conference on Embedded Software, EMSOFT'11*. IEEE, 319–328.
- [14] Maiza Claire et al. 2017. The W-SEPT Project: Towards Semantic-Aware WCET Estimation. In *17th International Workshop on Worst-Case Execution Time Analysis (WCET 2017)*, Vol. 57. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany, 9:1–9:13.
- [15] Antoine Colin and Guillem Bernat. 2002. Scope-tree: A program representation for symbolic worst-case execution time analysis. In *Proceedings 14th Euromicro Conference on Real-Time Systems. Euromicro RTS 2002*. IEEE, 50–59.
- [16] L. Cucu-Grosjean, L. Santinelli, M. Houston, C. Lo, T. Vardanega, L. Kosmidis, J. Abella, E. Mezzetti, E. Quiñones, and F. J. Cazorla. 2012. Measurement-Based Probabilistic Timing Analysis for Multi-path Programs. In *2012 24th Euromicro Conference on Real-Time Systems*. 91–101.
- [17] Wei Dou, Domenico Bianculli, and Lionel Briand. 2018. Model-Driven Trace Diagnostics for Pattern-Based Temporal Specifications. In *21th ACM/IEEE International Conference on Model Driven Engineering Languages and Systems*. 278–288.
- [18] Doron Drusinsky. 2000. The temporal rover and the ATG rover. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 1885 (2000), 323–330.
- [19] Daniel Emery. 2011. Headways on high speed lines. In *9th World Congress on Railway Research*. 22–26.
- [20] Andreas Ermedahl, Christer Sandberg, Jan Gustafsson, Stefan Bygde, and Björn Lisper. 2007. Loop bound analysis based on a combination of program slicing, abstract interpretation, and invariant analysis. In *7th International Workshop on Worst-Case Execution Time Analysis (WCET'07)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik.
- [21] Christian Ferdinand and Reinhold Heckmann. 2004. aiT: Worst-Case Execution Time Prediction by Static Program Analysis. In *Building the Information Society*, Renè Jacquart (Ed.). Springer US, Boston, MA, 377–383.
- [22] Thorsten Fischer, Jörg Niere, Lars Torunski, and Albert Zündorf. 1998. Story diagrams: A new graph rewrite language based on the unified modeling language and java. In *International Workshop on Theory and Application of Graph Transformations*. Springer, 296–309.
- [23] Franck Fleurey, Jim Steel, and Benoit Baudry. 2004. Validation in model-driven engineering: testing model transformations. In *Proceedings. 2004 First International Workshop on Model, Design and Validation, 2004*. 29–40.
- [24] Brian Gallagher. 2006. Matching structure and semantics: A survey on graph-based pattern matching. *AAAI FS 6* (2006), 45–53.
- [25] Holger Giese, Matthias Tichy, Sven Burmester, Wilhelm Schäfer, and Stephan Flake. 2003. Towards the compositional verification of real-time UML designs. *Proceedings of the ACM SIGSOFT Symposium on the Foundations of Software Engineering* (2003), 38–47.
- [26] Jan Gustafsson, Andreas Ermedahl, Christer Sandberg, and Björn Lisper. 2006. Automatic derivation of loop bounds and infeasible paths for WCET analysis using abstract execution. *Proc. of Real-Time Systems Symposium* (2006), 57–66.
- [27] Jeffery Hansen, Scott Hissam, and Gabriel A Moreno. 2009. Statistical-based wcet estimation and validation. In *9th International Workshop on Worst-Case Execution Time Analysis*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik.
- [28] Thomas Hartmann, François Fouquet, Assaad Moawad, Romain Rouvoy, and Yves Le Traon. 2019. GREYCAT: Efficient what-if analytics for data in motion at scale. *Information Systems* 83 (2019), 101–117.
- [29] Klaus Havelund. 2015. Rule-based runtime verification revisited. *Int. J. Software Tools Technol. Trans.* 17, 2 (2015), 143–170.
- [30] Klaus Havelund and Grigore Rosu. 2002. Synthesizing Monitors for Safety Properties. In *LNCS*. Vol. 2280. 342–356.
- [31] Jörg Herter and Jan Reineke. 2009. Making Dynamic Memory Allocation Static to Support WCET Analysis. In *9th International Workshop on Worst-Case Execution Time Analysis (WCET'09)*.
- [32] Wen-Chi Hou, Gultekin Ozsoyoglu, and Baldeo K. Taneja. 1989. Processing Aggregate Relational Queries with Hard Time Constraints. *SIGMOD Rec.* (1989), 10.
- [33] Ethan K. Jackson, Gabor Simko, and Janos Sztipanovits. 2013. Diversely enumerating system-level architectures. In *ACM International Conference on Embedded Software*. IEEE.
- [34] Axel Jantsch, Nikil Dutt, and Amir M. Rahmani. 2017. Self-awareness in systems on chip—A survey. *IEEE Design & Test* 34, 6 (2017), 8–26.
- [35] Jan Jürjens. 2003. Developing safety-critical systems with UML. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 2863 (2003), 360–372.
- [36] Jens Knop, Laura Kovács, and Jakob Zwirchmayr. 2013. WCET squeezing. (2013), 161.

- [37] V. P. Kozyrev. 2016. Estimation of the execution time in real-time systems. *Programming and Computer Software* 42, 1 (2016), 41–48.
- [38] S. Law and I. Bate. 2016. Achieving Appropriate Test Coverage for Reliable Measurement-Based Timing Analysis. In *2016 28th Euromicro Conference on Real-Time Systems (ECRTS)*. 189–199.
- [39] Xianfeng Li, Yun Liang, Tulika Mitra, and Abhik Roychoudhury. 2007. Chronos: A timing analyzer for embedded software. *Science of Computer Programming* 69, 1-3 (2007), 56–67.
- [40] Xiaocui Li, Zhangbing Zhou, Junqi Guo, Shangguang Wang, and Junsheng Zhang. 2019. Aggregated multi-attribute query processing in edge computing for industrial IoT applications. *Computer Networks* 151 (2019), 114–123.
- [41] Y.-T.S. Li and Sharad Malik. 1997. Performance analysis of embedded software using implicit path enumeration. *IEEE T. Comput. Aid. D.* 16, 12 (1997), 1477–1487.
- [42] Sung-Soo Lim et al. 1995. An accurate worst case timing analysis for RISC processors. *IEEE Transactions on Software Engineering* 21, 7 (1995), 593–604.
- [43] Björn Lisper. 2014. SWEET—a tool for WCET flow analysis. In *International Symposium On Leveraging Applications of Formal Methods, Verification and Validation*. Springer, 482–485.
- [44] Florian Martin, Martin Alt, Reinhard Wilhelm, and Christian Ferdinand. 1998. Analysis of loops. *Lecture Notes in Computer Science* 1383 (1998), 80–94.
- [45] Kristóf Marussy, Oszkár Semeráth, and Dániel Varró. 2020. Automated Generation of Consistent Graph Models with Multiplicity Reasoning. *IEEE Transactions on Software Engineering* (2020). <https://doi.org/10.1109/TSE.2020.3025732>
- [46] Marsha Chechik Michalis Famelis, Rick Salay. 2012. Partial models: Towards modeling and reasoning with uncertainty. In *International Conference on Software Engineering*. IEEE.
- [47] Gultekin Ozsoyoglu and Richard T. Snodgrass. 1995. Temporal and real-time databases: A survey. *IEEE Trans. Knowl. Data Eng.* 7, 4 (1995).
- [48] Christian Pek, Stefanie Manzinger, Markus Koschi, and Matthias Althoff. 2020. Using online verification to prevent autonomous vehicles from causing accidents. *Nature Machine Intelligence* 2, 9 (2020), 518–528.
- [49] Lee Pike, Alwyn Goodloe, Robin Morisset, and Sebastian Niller. 2010. Copilot: A Hard Real-Time Runtime Monitor. In *LNCIS*. Vol. 6418. 345–359.
- [50] Peter P.uschner and Anton V. Schedl. 1997. Computing Maximum Task Execution Times - A Graph-Based Approach. *Real-Time Systems* 13, 1 (1997), 67–91.
- [51] Leanna Rierson. 2017. *Developing Safety-Critical Software*. CRC Press. 22–27 pages.
- [52] Mooly Sagiv, Thomas Reps, and Reinhard Wilhelm. 2002. Parametric shape analysis via 3-valued logic. *ACM Transactions on Programming Languages and Systems* 24, 3 (2002), 193–298.
- [53] Rick Salay, Michalis Famelis, and Marsha Chechik. 2012. Language Independent Refinement Using Partial Modeling. In *FASE*. Springer.
- [54] Oszkár Semeráth, Rebeka Farkas, Gábor Bergmann, and Dániel Varró. 2020. Diversity of graph models and graph generators in mutation testing. *International Journal on Software Tools for Technology Transfer* 22, 1 (2020), 57–78.
- [55] Oszkár Semeráth, András Szabolcs Nagy, and Dániel Varró. 2018. A graph solver for the automated generation of consistent domain-specific models. In *40th International Conference on Software Engineering*. 969–980.
- [56] Michael Szvetits and Uwe Zdun. 2013. Systematic literature review of the objectives, techniques, kinds, and architectures of models at runtime. *Software & Systems Modeling* 15, 1 (2013), 31–69.
- [57] Juha Taina and Kimmo Raatikainen. 1996. Rodain: A real-time object-oriented database system for telecommunications. *International Conference on Information and Knowledge Management, Proceedings Part F129290* (1996), 10–14.
- [58] Joze Tavcar and Imre Horvath. 2019. A review of the principles of designing smart cyber-physical systems for run-time adaptation: Learned lessons and open issues. *IEEE Trans. Syst. Man Cybern. Syst.* 49, 1 (2019), 145–158.
- [59] The Eclipse Project. 2021. *Eclipse Modeling Framework*. The Eclipse Project. <http://www.eclipse.org/emf>.
- [60] Matthias Tichy, Holger Giese, and Andreas Seibel. 2006. Story diagrams in real-time software. In *Proc. of the 4th International Fijaba Days*.
- [61] Dániel Varró, Oszkár Semeráth, Gábor Szármyas, and Ákos Horváth. 2018. Towards the Automated Generation of Consistent, Diverse, Scalable and Realistic Graph Models. In *Graph Transformation, Specifications, and Nets (In Memory of Hartmut Ehrig)*.
- [62] Gergely Varró, Frederik Deckwerth, Martin Wieber, and Andy Schürr. 2015. An algorithm for generating model-sensitive search plans for pattern matching on EMF models. *Software & Systems Modeling* (2015), 597–621.
- [63] András Vörös et al. 2018. MoDeS3: Model-Based Demonstrator for Smart and Safe Cyber-Physical Systems. In *NASA Formal Methods*. 460–467.
- [64] I. Wenzel, R. Kirner, B. Rieder, and P.uschner. 2005. Measurement-based worst-case execution time analysis. In *Third IEEE Workshop on Software Technologies for Future Embedded and Ubiquitous Systems (SEUS'05)*. 7–10.
- [65] Reinhard Wilhelm et al. 2008. The worst-case execution-time problem-overview of methods and survey of tools. *Transactions on Embedded Computing Systems* 7, 3 (2008).

- [66] Cheng Xie, Beibei Yu, Zuoying Zeng, Yun Yang, and Qing Liu. 2021. Multilayer Internet-of-Things Middleware Based on Knowledge Graph. *IEEE Internet of Things Journal* 8, 4 (2021), 2635–2648.
- [67] Haitao Zhu, Matthew B. Dwyer, and Steve Goddard. 2009. Predictable runtime monitoring. *Proceedings - Euromicro Conference on Real-Time Systems 2* (2009), 173–183.

A PROOF SKETCHES

Proposition 5.1. Let τ be the execution time of the query program q on the concrete model M ,

$$CL = \max_{\mathbf{x}_i \in \mathcal{X}} \sum c_i \cdot \mathbf{x}_i \text{ subject to } \mathcal{S}_{\text{IPET}}, \quad DS_M = \max_{\mathbf{x}_i \in \mathcal{X}} \sum c_i \cdot \mathbf{x}_i \text{ subject to } \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}},$$

where CL is the classical IPET estimated obtained from q , and DS_M is the domain-specific estimate with flow facts derived from M using Algorithm 4. Then $\tau \leq DS_M \leq CL$.

PROOF SKETCH. $\tau \leq DS_M$ (*safety*): Consider any execution path π of q . Because CL is safe, there is a solution $k_\pi: \mathcal{X} \rightarrow \mathbb{Z}$ such that $k_\pi(\mathbf{f}(e)) = \pi\#e$ for all edges $e \in E$ of the CFG of q and $k_\pi \models \mathcal{S}_{\text{IPET}}$.

Consider the linear equations in $\mathcal{S}_{\text{flow}}$. We have a single linear equation for each basic block $bb \in BB$. If bb is a loop header, then every execution of bb is represented by either a match of ψ_{bb} of ψ'_{bb} . Thus, $\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \pi\#e = \sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} k_\pi(\mathbf{f}(e)) = M\#\psi_{bb} + M\#\psi'_{bb}$, which means the corresponding linear equation in $\mathcal{S}_{\text{flow}}$ holds. Otherwise, every execution of bb is represented by a match of ψ_{bb} . Thus, $\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \pi\#e = \sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} k_\pi(\mathbf{f}(e)) = M\#\psi_{bb}$, which means the corresponding linear equation in $\mathcal{S}_{\text{flow}}$ holds.

Therefore, we have $k_\pi \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}$. We also have $\tau \leq CL \leq DS_M$ using the safety of CL .

$DS_M \leq CL$ (*tightness*): Assume that $DS_M > CL$. Then there is some $k \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}$ such that $g(k) = \sum_{\mathbf{x}_i} k(\mathbf{x}_i) > CL = g(k^*)$, where $k^* \models \mathcal{S}_{\text{IPET}}$ is the optimal solution of the classical IPET integer program with $g(k^*) \geq g(k')$ for all $k' \models \mathcal{S}_{\text{IPET}}$. However, $k \models \mathcal{S}_{\text{IPET}}$ (because $\mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}} \models \mathcal{S}_{\text{IPET}}$), which means k^* cannot be optimal. Thus the assumption cannot hold. \square

Before we prove Proposition 5.2, we turn our attention to Proposition 5.3. Let $P = \langle \mathcal{O}_P, \mathcal{I}_P, \mathcal{S}_P \rangle$ be a partial model, $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ be a theory. Moreover, let $\mathcal{S}_{\text{IPET}}, \mathcal{S}_{\text{merge}}, \mathcal{S}_{P'} = \mathcal{S}_P \cup \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{merge}}$, $P' = \langle \mathcal{O}_{P'}, \mathcal{I}_{P'}, \mathcal{S}_{P'} \rangle$, $\Phi' = \Phi \cup \Psi$, $\mathcal{T}' = \langle \Phi', \mathbf{r}' \rangle$ be the IPET linear equations, the merging linear equations, the scope, the partial model, the predicates, and the theory output by Algorithm 5 for the WCET estimation of a graph query program q , respectively. Also let $g_{\text{IPET}}(k) = \sum_{e \in E} w(e) \cdot \mathbf{f}(e)$ denote the objective function of the IPET analysis (for the CFG = $\langle V, E, s, t, w, tr \rangle$ of q and the function $\mathbf{f}: E \rightarrow \mathcal{X}$), $g_{\text{witness}}(M) = \max_{k \models \mathcal{S}_M} g_{\text{IPET}}(k)$ denote the objective function of the witness generation task, and let $\mathcal{S}_{\text{flow}}(M)$ and $DS_M(M)$ denote the flow facts and the WCET estimate obtained for q by Algorithm 4 for the concrete model M . We will use the following lemmas:

Lemma A.1. If $M = \langle \mathcal{O}_M, \mathcal{I}_M, \mathcal{S}_M \rangle \in \text{solutions}(P', \mathcal{T}')$, then $M \in \text{solutions}(P, \mathcal{T})$.

PROOF SKETCH. $M \in \text{solutions}(P', \mathcal{T}')$ means that $P' \succcurlyeq M$ and $M \models \mathcal{T}'$. We will show that we also have $P \succcurlyeq M$ and $M \models \mathcal{T}$.

$P \succcurlyeq M$: Consider the abstraction function $abs: \mathcal{O}_M \rightarrow \mathcal{O}_P$ from $P' \succcurlyeq_{abs} M$. Since P has the same objects \mathcal{O}_P and interpretation \mathcal{I}_P as P' , abs will serve as the abstraction function for $P \succcurlyeq_{abs} M$, too. Because $\mathcal{S}_{P'} = \mathcal{S}_P \cup \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{merge}} \supseteq \mathcal{S}_P$ and $\mathcal{S}_M \models \mathcal{S}_{P'}$ due to $P' \succcurlyeq M$, we also have $\mathcal{S}_M \models \mathcal{S}_P$ as required by Definition 4.13.

$M \models \mathcal{T}$: Notice that $\Phi' = \Phi \cup \Psi \supseteq \Phi$. $M \models \mathcal{T}'$ means that $\mathcal{S}_M \models \mathbf{r} = M\#\varphi$ for all $\varphi \in \Phi'$. Thus $\mathcal{S}_M \models \mathbf{r} = M\#\varphi$ holds for all $\varphi \in \Phi \subseteq \Phi'$, which means $M \models \mathcal{T}$. \square

Lemma A.2. If $M \in \text{solutions}(P', \mathcal{T}')$, then $DS_M(M) \geq g_{\text{witness}}(M)$.

PROOF SKETCH. We need to show that

$$DS_M(M) = \max_{k \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)} g_{\text{IPET}}(k) \geq g_{\text{witness}}(M) = \max_{k \models \mathcal{S}_M} g_{\text{IPET}}(k),$$

for which it suffices to demonstrate that every solution $k \models \mathcal{S}_M$ of the scope associated with M is also a solution $k \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)$ of the IPET integer program, i.e., $\mathcal{S}_M \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)$. We proceed by case analysis on the linear equations $eq \in \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)$:

- If $eq \in \mathcal{S}_{IPET}$, then $\mathcal{S}_M \vDash \mathcal{S}_{P'} = \mathcal{S}_P \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge} \supseteq \mathcal{S}_{IPET} \ni eq$ implies $\mathcal{S}_M \vDash eq$.
- For each $eq = [\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = M\#\psi_{bb} + M\#\psi'_{bb}] \in \mathcal{S}_{flow}(M)$, where bb is a loop header, we have $\mathcal{S}_M \vDash \mathbf{r}'(\psi_{bb}) = M\#\psi_{bb}$ and $\mathcal{S}_M \vDash \mathbf{r}'(\psi'_{bb}) = M\#\psi'_{bb}$, because $M \vDash \mathcal{T}'$. We also have $\mathcal{S}_M \vDash \mathcal{S}_{P'} \supseteq \mathcal{S}_{merge} \ni [\mathbf{r}'(\psi_{bb}) + \mathbf{r}'(\psi'_{bb}) - \sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = 0]$. Therefore $\mathcal{S}_M \vDash eq$ follows by substitution and rearranging.
- For each $eq = [\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = M\#\psi_{bb}] \in \mathcal{S}_{flow}(M)$, where bb is not a loop header, we have $\mathcal{S}_M \vDash \mathbf{r}'(\psi_{bb}) = M\#\psi_{bb}$, because $M \vDash \mathcal{T}'$. We also have $\mathcal{S}_M \vDash \mathcal{S}_{P'} \supseteq \mathcal{S}_{merge} \ni [\mathbf{r}'(\psi_{bb}) - \sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = 0]$. Therefore $\mathcal{S}_M \vDash eq$ follows by substitution and rearranging. \square

Definition A.1. The function $rename: \mathcal{X} \rightarrow \mathcal{X}$ is a *renaming of variables* if it is bijective (i.e., there is some $rename^{-1}: \mathcal{X} \rightarrow \mathcal{X}$ such that $rename^{-1} \circ rename$ is the identity function on \mathcal{X}). A renaming of variables is *stationary w.r.t. the theory* $\mathcal{T} = \langle \Phi, \mathbf{r} \rangle$ if $rename(\mathbf{r}(\varphi)) = \mathbf{r}(\varphi)$ for all $\varphi \in \Phi$. The partial model $P^{rename} = \langle \mathcal{O}_P, \mathcal{I}_P, \mathcal{S}_P^{rename} \rangle$ is a *renaming of* $P = \langle \mathcal{O}_P, \mathcal{I}_P, \mathcal{S}_P \rangle$, where \mathcal{S}_P^{rename} is obtained from \mathcal{S}_P by replacing each variable $\mathbf{x} \in \mathcal{X}$ with $rename(\mathbf{x})$.

Lemma A.3. If $M = \langle \mathcal{O}_M, \mathcal{I}_M, \mathcal{S}_M \rangle \in solutions(P, \mathcal{T})$, then there is some $M' = \langle \mathcal{O}_{M'}, \mathcal{I}_{M'}, \mathcal{S}_{M'} \rangle \in solutions(P', \mathcal{T}')$ such that $M^{rename} \succcurlyeq M'$ and $DS_M(M) = DS_{M'}(M') = g_{witness}(M')$ for some renaming of variables $rename: \mathcal{X} \rightarrow \mathcal{X}$ stationary w.r.t. \mathcal{T} .

PROOF SKETCH. W.l.o.g. we may assume that the variable of \mathcal{S}_M as disjoint from those of $\mathcal{S}_{IPET}(M)$ and $\mathcal{S}_\Psi(M)$ and $rename$ is the identity function ($M^{rename} = M$). Otherwise, we can pick $rename$ to make the variables of \mathcal{S}_M^{rename} disjoint from those of \mathcal{S}_{IPET} , $\mathcal{S}_{flow}(M)$ and $\mathcal{S}_\Psi(M)$.

Let $\mathcal{O}_{M'} = \mathcal{O}_M$, $\mathcal{I}_{M'} = \mathcal{I}_M$, and $\mathcal{S}_{M'} = \mathcal{S}_M \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{flow}(M) \cup \mathcal{S}_\Psi(M)$, where $\mathcal{S}_\Psi(M) = \{\mathbf{r}'(\psi) = M\#\psi \mid \psi \in \Psi\}$. We will show that $M \succcurlyeq_{abs_M} M'$, where abs_M is the identity function on \mathcal{O}_M . Since $\mathcal{I}_{M'} = \mathcal{I}_M$, it suffices to show that $\mathcal{S}_{M'} \vDash \mathcal{S}_M$, which follows from $\mathcal{S}_{M'} \supseteq \mathcal{S}_M$.

To show that $M' \in solutions(P', \mathcal{T}')$, we must prove that (i) $P \succcurlyeq M'$, (ii) $M' \vDash \mathcal{T}'$, and (iii) M' is concrete. (i) Consider the refinement function $abs_P: \mathcal{O}_M \rightarrow \mathcal{O}_P$ from $P \succcurlyeq_{abs_P} M$. Since $\mathcal{O}_M = \mathcal{O}_{M'}$, $\mathcal{I}_M = \mathcal{I}_{M'}$, and $\mathcal{I}_P = \mathcal{I}_{P'}$, to demonstrate $P' \succcurlyeq_{abs_P} M$, it suffices to show that $\mathcal{S}_{M'} \vDash \mathcal{S}_{P'}$. We will show for each linear (in)equality $eq \in \mathcal{S}_{P'} = \mathcal{S}_P \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge}$ that $\mathcal{S}_{M'} \vDash eq$:

- If $eq \in \mathcal{S}_P$, then $\mathcal{S}_{M'} \vDash eq$ because $\mathcal{S}_M \vDash \mathcal{S}_P \ni eq$ and $\mathcal{S}_{M'} \vDash \mathcal{S}_M$.
- If $eq \in \mathcal{S}_{IPET}$, then $\mathcal{S}_{M'} \vDash eq$ because $\mathcal{S}_{M'} \supseteq \mathcal{S}_{IPET} \ni eq$.
- If $eq = [\mathbf{r}'(\psi_{bb}) + \mathbf{r}'(\psi'_{bb}) - \sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = 0] \in \mathcal{S}_{merge}$, where bb is a loop header, then we have $\{\mathbf{r}'(\psi_{bb}) = M\#\psi_{bb}, \mathbf{r}'(\psi'_{bb}) = M\#\psi'_{bb}\} \subseteq \mathcal{S}_\Psi(M)$ and we also have $[\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = M\#\psi_{bb} + M\#\psi'_{bb}] \in \mathcal{S}_{flow}(M)$. Then $\mathcal{S}_{M'} \supseteq \mathcal{S}_{flow}(M) \cup \mathcal{S}_\Psi(M) \vDash eq$ follows by substitution and rearranging.
- If $eq = [\mathbf{r}'(\psi_{bb}) - \sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = 0] \in \mathcal{S}_{merge}$, where bb is a not loop header, then we have $[\mathbf{r}'(\psi_{bb}) = M\#\psi_{bb}] \in \mathcal{S}_\Psi(M)$ and we also have $[\sum_{e=\langle n_1, n_2 \rangle \in E, tr(n_1)=bb} \mathbf{f}(e) = M\#\psi_{bb}] \in \mathcal{S}_{flow}(M)$. Then $\mathcal{S}_{M'} \supseteq \mathcal{S}_{flow}(M) \cup \mathcal{S}_\Psi(M) \vDash eq$ follows by substitution and rearranging.

(ii) To see that $M' \vDash \mathcal{T}'$, we also use case analysis for each predicate $\varphi \in \Phi'$:

- If $\varphi \in \Phi$, then $\mathcal{S}_{M'} \supseteq \mathcal{S}_M \vDash \mathbf{r}(\varphi) = M\#\varphi$, because $M \vDash \mathcal{T}$. This implies $\mathcal{S}_{M'} \vDash \mathbf{r}'(\varphi) = M'\#\varphi$, because $\mathbf{r}(\varphi) = \mathbf{r}'(\varphi)$, and $M\#\varphi = M'\#\varphi$ due to $\mathcal{I}_M = \mathcal{I}_{M'}$.
- If $\varphi \in \Psi$, then $\mathcal{S}_{M'} \supseteq \mathcal{S}_\Psi(M) \vDash \mathbf{r}'(\varphi) = M\#\varphi$, which implies $\mathcal{S}_{M'} \vDash \mathbf{r}'(\varphi) = M'\#\varphi$ due to $\mathcal{I}_M = \mathcal{I}_{M'}$ as above.

(iii) Since M is concrete, $\mathcal{I}_{M'} = \mathcal{I}_M$ contains only 1 and 0 logic values, and \mathcal{S}_M has some solution k . To conclude that M' is concrete, we will construct a solution k' of $\mathcal{S}_{M'}$.

Consider an execution path π of q on the input model M and the associated solution k_π of $\mathcal{S}_{\text{IPET}}$, where $k_\pi(\mathfrak{f}(e)) = \pi\#e$. Recall from Proposition 5.1 that $k_\pi \models \mathcal{S}_{\text{flow}}(M)$. Now let

$$k'(\mathfrak{x}) = \begin{cases} k_\pi(\mathfrak{x}) & \text{if } \mathfrak{x} = \mathfrak{f}(e) \text{ for some } e \in E, \\ M\#\psi & \text{if } \mathfrak{x} = \mathfrak{r}'(\psi) \text{ for some } \psi \in \Psi, \\ k(\mathfrak{x}) & \text{otherwise.} \end{cases}$$

Recall that $\{\mathfrak{r}'(\psi) \mid \psi \in \Psi\}$ (i.e., the variables appearing in $\mathcal{S}_\Psi(M)$) and $\text{Ran } \mathfrak{f}$ (i.e., the variables appearing in $\mathcal{S}_{\text{IPET}}$ and $\mathcal{S}_{\text{flow}}(M)$) are disjoint from each other and from the variables appearing in \mathcal{S}_M , so this is well-defined. Now $k' \models \mathcal{S}_{\text{IPET}}$, $k' \models \mathcal{S}_{\text{flow}}(M)$, $k' \models \mathcal{S}_\Psi(M)$, and $k' \models \mathcal{S}_M$. Thus $k' \models \mathcal{S}_{M'} = \mathcal{S}_M \cup \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M) \cup \mathcal{S}_\Psi(M)$.

Notice that $M\#\psi = M'\#\psi$ for all $\psi \in \Psi$ implies that $\mathcal{S}_{\text{flow}}(M) = \mathcal{S}_{\text{flow}}(M')$. Therefore,

$$DS_M(M) = \max_{k \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)} g_{\text{IPET}}(k) = DS_M(M') = \max_{k \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M')} g_{\text{IPET}}(k).$$

We can transform any solution $k_{\text{IPET}} \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)$ into an equivalent solution $k_{\text{witness}} \models \mathcal{S}_{M'} = \mathcal{S}_M \cup \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M) \cup \mathcal{S}_\Psi(M)$ with $g_{\text{IPET}}(k_{\text{IPET}}) = g_{\text{IPET}}(k_{\text{witness}})$ by changing the values of some variables $\mathfrak{x} \in \mathcal{X} \setminus \text{Ran } \mathfrak{f}$, since the variables appearing in $\mathcal{S}_M \cup \mathcal{S}_\Psi(M)$ are disjoint from the variables $\text{Ran } \mathfrak{f}$ appearing in $\mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M)$ and all have a zero weight in g_{IPET} . Therefore,

$$DS_M(M') = \max_{k \models \mathcal{S}_{\text{IPET}} \cup \mathcal{S}_{\text{flow}}(M')} g_{\text{IPET}}(k) = g_{\text{witness}}(M') = \max_{k \models \mathcal{S}_{M'}} g_{\text{IPET}}(k). \quad \square$$

Proposition 5.3 (Witness model). Let $DS_M(M)$ be the domain-specific WCET estimate of a query program q obtained by Algorithm 4 for a concrete model M , DS_P be the domain-specific WCET estimate of q for a partial model P and theory \mathcal{T} by Algorithm 5, and M^* be the witness model for the WCET of q , i.e., the optimal solution of DS_P . Then $M^* \in \text{solutions}(P, \mathcal{T})$ and $DS_M(M) \leq DS_M(M^*) = DS_P$ for all $M \in \text{solutions}(P, \mathcal{T})$.

PROOF SKETCH. Consider a model $M^* \in \text{solutions}(P', \mathcal{T}')$ that maximizes g_{witness} . By Lemma A.1, we also have $M^* \in \text{solutions}(P, \mathcal{T})$.

$DS_M(M^*) = DS_P$: DS_P is defined as $\max_{M \in \text{solutions}(P', \mathcal{T}')} g_{\text{witness}}(M) = g_{\text{witness}}(M^*)$. By Lemma A.2, $DS_M(M^*) \geq g_{\text{witness}}(M^*)$, which means $DS_M(M^*) \geq g_{\text{witness}}(M^*) = DS_P$. Not let us apply the construction from Lemma A.3 to M^* to obtain $M^{*'}$. We have $DS_M(M^*) = DS_M(M^{*'}) = g_{\text{witness}}(M^{*'}) \leq g_{\text{witness}}(M^*)$, which concludes the proof.

$DS_M(M) \leq DS_M(M^*)$: Assume the contrary. Then we have some $M' \in \text{solutions}(P', \mathcal{T}')$ such that $g_{\text{witness}}(M^*) = DS_M(M^*) < DS_M(M) = DS_M(M') = g_{\text{witness}}(M')$ by Lemma A.3. This leads to a contradiction, because M^* is a maximizer of g_{witness} . \square

Now we can use Proposition 5.3 to prove Proposition 5.2 as follows.

Proposition 5.2 (Safety and tightness). Let $\tau(M)$ be the execution time of a query program q on a concrete model M , P be partial model, \mathcal{T} be a theory, and

$$CL = \max_{\mathfrak{x}_i \in \mathcal{X}} \sum c_i \cdot \mathfrak{x}_i \text{ subject to } \mathcal{S}_{\text{IPET}}, \quad DS_P = \max_{\mathfrak{x}_i \in \mathcal{X}} \sum c_i \cdot \mathfrak{x}_i \text{ subject to } \langle P', \mathcal{T}' \rangle,$$

where CL is the classical IPET estimated obtained from q , and DS_P is the domain-specific estimate based on the extended graph generation problem from Algorithm 5. Then $\tau(M) \leq DS_P \leq CL$ for all $M \in \text{solutions}(P, \mathcal{T})$.

PROOF SKETCH. $\tau(M) \leq DS_P$ (*safety*): By Proposition 5.1, $\tau(M) \leq DS_M(M)$. Also, by Proposition 5.3, $DS_M(M) \leq DS_M(M^*) = DS_P$. Thus, $\tau(M) \leq DS_M(M) \leq DS_M(M^*) = DS_P$.

$DS_P \leq CL$ (*tightness*): By Proposition 5.1, $DS_M(M^*) \leq CL$. Also, by Proposition 5.3, $DS_P = DS_M(M^*)$. Thus, $DS_P = DS_M(M^*) \leq CL$. \square

Proposition 5.4 (Tightening by refinement). Let $DS_P(P, \mathcal{T})$ denote the domain-specific WCET estimate of a query program q for a partial model P and theory \mathcal{T} obtained by Algorithm 5 and $P \succcurlyeq Q$. Then $DS_P(Q, \mathcal{T}) \leq DS_P(P, \mathcal{T})$. In particular, if $P = P_{init}$ is the initial partial model for a metamodel $\langle \Sigma, \alpha \rangle$ from Section 4.5, then we may see that the WCET estimate for any partial model conforming to the metamodel is at least as tight as the DS_Σ estimate for the metamodel.

PROOF SKETCH. The model generation task objective function $g_{witness}$ coincides for both $DS_P(P, \mathcal{T})$ and $DS_P(Q, \mathcal{T})$. Thus, to show that

$$DS_P(P, \mathcal{T}) = \max_{M' \in solutions(P, \mathcal{T})} g_{witness}(M') \geq DS_P(Q, \mathcal{T}) = \max_{M' \in solutions(Q, \mathcal{T})} g_{witness}(M'),$$

we need to show that every $M' \in solutions(Q, \mathcal{T})$ also satisfies $M' \in solutions(P, \mathcal{T})$.

Consider the partial models $P' = \langle O_P, \mathcal{I}_P, \mathcal{S}_{P'} \rangle$ and $Q' = \langle O_Q, \mathcal{I}_Q, \mathcal{S}_{Q'} \rangle$ obtained by Algorithm 5 for the respective input partial models, where $\mathcal{S}_{Q'} = \mathcal{S}_Q \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge}$ and $\mathcal{S}_{P'} = \mathcal{S}_P \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge}$. The systems on linear inequalities \mathcal{S}_{IPET} and \mathcal{S}_{merge} are the same for both cases, since they do not depend on the input partial model.

We will show that $P' \succcurlyeq_{abs_Q} Q'$, where $abs_Q: O_Q \rightarrow O_P$ is the abstraction function $P \succcurlyeq_{abs_Q} Q$ from Q to P . Because P and P' , as well as Q and Q' have the same objects and interpretations, we only need to show $\mathcal{S}_{Q'} \models \mathcal{S}_{P'}$. We already have $\mathcal{S}_Q \models \mathcal{S}_P$ due to $P \succcurlyeq_{abs_Q} Q$, so we conclude $\mathcal{S}_Q \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge} \models \mathcal{S}_P \cup \mathcal{S}_{IPET} \cup \mathcal{S}_{merge}$.

Now consider any $M' \in solutions(Q', \mathcal{T})$. We have $Q' \succcurlyeq_{abs_{M'}} M'$ and $M' \models \mathcal{T}$ by Definition 4.23. Then we also have $P' \succcurlyeq_{abs_{M'} \circ abs_{Q'}} M'$ by the associativity of refinement, from which we conclude $M' \in solutions(P', \mathcal{T})$. \square