

Field-sensitive program slicing^{☆,☆☆}Carlos Galindo^a, Jens Krinke^b, Sergio Pérez^a, Josep Silva^{a,*}^a VRAIN, Universitat Politècnica de València, Camino de Vera s/n, 46022, Valencia, Spain^b CREST Centre, University College London, Gower Street, London, WC1E 6BT, United Kingdom

ARTICLE INFO

Dataset link: <http://dx.doi.org/10.5281/zenodo.10439063>

Keywords:

Program analysis

Program slicing

Composite data structures

ABSTRACT

The granularity level of the traditional program dependence graph (PDG) for composite data structures (tuples, lists, records, objects, etc.) is inaccurate when slicing their inner elements. We present the constrained-edges PDG (CE-PDG) that addresses this accuracy problem. The CE-PDG enhances the representation of composite data structures by decomposing statements into a subgraph that represents the inner elements of the structure, and the inclusion and propagation of data constraints along the CE-PDG edges allow for accurate slicing of complex data structures. Both extensions are conservative with respect to the traditional PDG, in the sense that all slicing criteria (and more) that can be specified in the PDG can be also specified in the CE-PDG, and the slices produced with the CE-PDG are always smaller or equal to the slices produced by the PDG. An evaluation of our approach shows a reduction in the size of the slices of around 10%.

1. Introduction

The *Program Dependence Graph* (PDG) (Ottenstein and Ottenstein, 1984; Ferrante et al., 1987) represents the statements of a program as a collection of nodes, and their control and data dependencies as edges. The PDG is used in *program slicing* (Silva, 2012; Tip, 1995), a technique for program analysis and transformation whose main objective is to extract from a program the set of statements, the so-called *program slice* (Weiser, 1981), that affect the values of a set of variables v at a program point p ((p, v)), which is known as the *slicing criterion* (Ottenstein and Ottenstein, 1984). Program slicing is applied in many disciplines such as software maintenance (Hajnal and Forgács, 2012), debugging (DeMillo et al., 1996), code obfuscation (Majumdar et al., 2007), and program specialization (Ochoa et al., 2005), among others.

The original PDG is not able to handle all the features that most modern programming languages offer. Therefore, several extensions and enhancements of the PDG have been proposed to represent features like arbitrary control-flow (Galindo et al., 2021; Ball and Horwitz, 1993); exception handling (Galindo et al., 2023b; Allen and Horwitz, 2003); interprocedural behavior (Clarke et al., 2002; Binkley, 1993); or concurrency (Krinke, 2003b; Chen and Xu, 2001a); among others. Nevertheless, there is still a largely unaddressed problem that is a source of imprecision and that affects all programming languages: the slicing of composite data structures. Finite composite data structures

can be atomized (Ramalingam et al., 1999) and then sliced as usual, however, infinite or recursive data structures cannot be atomized and slicing them is therefore imprecise.

In this paper, we propose a general method that solves the problem of accurately representing and slicing any composite data structure, even if it is recursive (infinite data structures can be also sliced) or if it is collapsed and expanded again (we solve the *slicing pattern matching* problem (Silva et al., 2012), which is explained in Section 2). The key ideas are (i) to expand the PDG with new nodes to precisely represent the subexpressions of the data structures, and (ii) to introduce the concept of *constrained edges*: we label the PDG edges with information about the data structures so that this information can be used at slicing time to know exactly which edges should be traversed. We call the new resulting graph the *Constrained-Edges PDG* (CE-PDG). Finally, (iii) we provide a new slicing algorithm that takes advantage of constrained edges, limiting the traversal when necessary, and obtaining more accurate slices in the presence of composite data structures.

The main goal of our technique is to handle recursive data structures in combination with pattern matching to make program slicing field-sensitive in order to improve the accuracy of slicing programs with composite data structures.

This paper is an extended version of work (Galindo et al., 2022) presented at the 20th International Conference on Software Engineering and Formal Methods (SEFM 2022). In this version, we extend

[☆] This work has been partially supported by grant PID2019-104735RB-C41 funded by MCIN/AEI/10.13039/501100011033 and by the Generalitat Valenciana under grant CIPROM/2022/6 (FassLow). Carlos Galindo was partially supported by the Spanish Ministerio de Universidades under grant FPU20/03861.

^{☆☆} Editor: Antonio Filieri.

* Corresponding author.

E-mail addresses: cargaji@vrain.upv.es (C. Galindo), j.krinke@ucl.ac.uk (J. Krinke), serperu@dsic.upv.es (S. Pérez), jsilva@dsic.upv.es (J. Silva).

```

1 foo() {
2   struct S {
3     string name;
4     int age;
5   };
6   S person = {"John",36};
7   int maxAge = person.age;
8   std::cout << maxAge;
9 }

1 enum Light {
2   Red = 0,
3   Yellow = 1,
4   Green = 2
5 }
6 void Main() {
7   Light pass = Light.Yellow
8   | Light.Green;
9   Console.WriteLine(pass);
10 }

(a) Records (C++)
(b) Enums (C#)

1 void foo(int arg){
2   int[] nums = {2,arg,27};
3   int x = nums[2];
4   System.out.println(x);
5 }

1 class Person:
2   def __init__(self,name,age):
3     self.name = name
4     self.age = age
5   p1 = Person("John", 36)
6   print(p1.age)

(c) Arrays (Java)
(d) Objects (Python)

```

Fig. 1. Slicing composite data structures (slicing criterion underlined, minimal slice in black).

the technical results with new formalizations, theoretical results, and their proofs. Moreover, we provide extended examples, showing that, e.g., other state-of-the-art techniques, such as atomization, cannot solve the addressed problem. The slicing algorithm proposed has been revised and is more efficient and precise than the one presented before. All experiments have been repeated to evaluate the new performance, which has improved by one order of magnitude.

The rest of the paper is structured as follows: Section 2 demonstrates the problems in slicing composite data structures. Section 3 gives a short introduction to program slicing with basic definitions. Section 4 presents the CE-PDG and how it is used for slicing. Section 5 presents an implementation and an empirical evaluation of the proposed technique. It is followed by a discussion of related work (Section 6) and conclusions (Section 7).

2. Slicing composite data structures

In this section, we show the inaccuracy problems caused by the traditional PDG when it is used to slice programs with complex data structures.

It is important to remark that these problems can be studied and solved at the level of the PDG (i.e., for intraprocedural programs). Because we can present the fundamental ideas and solutions of field-sensitive slicing at this level, we omit the more complex representation in the *System Dependence Graph* (SDG) (Horwitz et al., 1990) (i.e., for interprocedural programs). In this way, we keep the presentation easier to understand, avoiding the complexity introduced by the SDG (procedure calls, input/output edges, summary edges...). Of course, an extension of our work for the SDG is possible and will increase the precision of our technique by propagating dependencies throughout procedures.

Example 1 (PDG's Composite Data structures). Consider the four fragments of code with different data structures shown in Fig. 1. We are interested in the values computed at the slicing criterion (the underlined variable in blue). The only part of the code that can affect the slicing criterion (i.e., the minimal slice) is colored in green. Nevertheless, the slice computed with the traditional PDG contains the whole program in the four cases.

The complexity of these structures, together with the lack of granularity of the PDG (each node of the PDG represents a statement) results in a lack of accuracy when slicing these structures. Slicing algorithms cannot remove unnecessary inner components of statements using only

```

1 foo(X,Y) ->      1 foo(X,Y) ->      1 foo(X,Y) ->
2 {A,B} = {X,Y},  2 {A,B} = {X,Y},  2 {A,B} = {X,Y},
3 Z = {[8],A},    3 Z = {[8],A},  3 Z = {[8],A},
4 {[C],D} = Z.    4 {[C],D} = Z.    4 {[C],D} = Z.

(a) Original Program  (b) PDG Slice  (c) Minimal Slice

```

Fig. 2. Slicing Erlang tuples (slicing criterion underlined, slice in black).

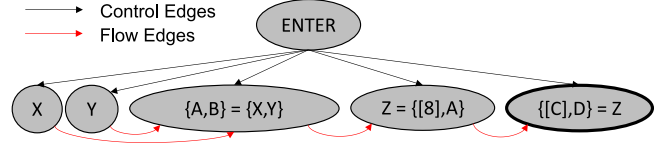


Fig. 3. Erlang program and corresponding PDG.

the information of the PDG. Unfortunately, this is only one part of the imprecision, because all variables wrongly captured by the slice trigger a snowball effect, effectively including in the slice all the parts of the program that potentially influence them. For instance, in Fig. 1(c), including in the slice variable `arg` implies also including in the slice each expression that is passed as argument to `foo`.

In some cases, it is possible to solve the situation with a program transformation (Korel and Laski, 1988, 1990; Binkley and Gallagher, 1996). For instance, in Fig. 1(a) we could replace `person = {"John",36};` by: `person.name = "John"; person.age = 36;`. Or, similarly in Fig. 1(c), we could replace `nums = {2,arg,27};` by `nums[0] = 2; nums[1] = arg; nums[2] = 27;`. This transformation, called *atomization* (Ramalingam et al., 1999), decomposes data structures in simpler assignments for each of their components. It uses the qualified name `person.age` or the indexed array `nums[0]` as the name of an independent variable (Binkley and Gallagher, 1996). Note that only finite data structures can be atomized. An alternative approach uses the AST nodes of a program as PDG nodes (Silva et al., 2012). Unfortunately, despite solving several problems, these approaches are also imprecise because they are unable to resolve the most problematic constructs: recursive (infinite) data types and pattern matching. Let us illustrate the problem induced by pattern matching with an example.

Example 2 (Pattern Matching). Consider the fragment of Erlang code in Fig. 2(a), where we are interested in the values computed at variable `C` (the slicing criterion is `<4,C>`). The only part of the code that can affect the values at `C` (i.e., the minimal slice) is colored in green in Fig. 2(c). Nevertheless, the slice computed with the PDG (shown in Fig. 2(b)) contains the whole program. This is again a potential source of more imprecisions outside this function because it wrongly includes in the slice the parameters of function `foo` and, thus, in calls to `foo` their arguments and the code in which they depend are also included.

Consider Fig. 3, which represents the PDG of the code in Fig. 2, where the slicing criterion is the node marked with a bolded border. In the example, a whole data structure (the tuple `{[8],A}`) has been collapsed to a variable (`Z`) and then expanded again (`{[C],D}`). Therefore, the list `[C]` depends on the list `[8]`. Nevertheless, the traditional PDG represents the equality as a whole, making `[C]` flow dependent on `Z`, and in turn, `Z` flow depends on `A`. Because flow dependence is transitive, slicing the PDG wrongly infers that `C` depends on `A` (`A` is wrongly included in the slice for `C`), and this lack of precision is propagated to parameter `X`.

This problem worsens in the presence of recursive data types. For instance, trees or objects (consider a class `A` with a field of type `A`, which produces a potentially infinite data type) can prevent the slicer from knowing statically what part of the collapsed structure is needed. Späth et al. (2019, pp. 2–3) present an interesting discussion and example about this problem.

3. Background: The program dependence graph

This section briefly summarizes the theoretical background of the PDG that is needed to keep the paper self-contained. Extended explanations about the applications of the PDG can be found in program slicing surveys (Tip, 1995; Binkley and Gallagher, 1996; Silva, 2012). Readers already familiar with the PDG can skip this section.

Given a program, two kinds of dependencies can be defined over it to construct the PDG (Ferrante et al., 1987): control dependence and flow dependence (aka data dependence).

Definition 1 (Control Dependence). Let G be a CFG (Control-Flow Graph). A node n_1 post-dominates a node n_2 in G if all paths in the CFG from n_2 to the Exit node of the CFG traverse n_1 . Node n_1 is *control dependent* on node n_2 if and only if n_1 post-dominates one but not all of n_2 's successors.

Definition 2 (Flow Dependence). A node n_2 is *flow dependent* on a preceding node n_1 if (1) n_1 defines a variable x , (2) n_2 uses x , and (3) there exists a control-flow path from n_1 to n_2 where x is not defined.

Definition 3 (Program Dependence Graph). Given a procedure p , its *Program Dependence Graph* (PDG) is a graph $G = (N, E)$, where there is a node in N to represent each statement in p ; and E is a set of edges that represent all control and flow dependencies between the nodes in N .

4. Constrained-edges program dependence graph

This section introduces the CE-PDG, in which the key idea is to expand all PDG nodes where a composite data structure is defined or used. This expansion augments the PDG with a tree representation for composite data structures. We describe how this structure is generated, and we introduce a new kind of dependence edge used to build this tree structure. For this, we formally define the concepts of *constraint* and *constrained edge*, describe the different types of constraint, and how they affect the graph traversal in the slicing process.

4.1. Extending the PDG

Fig. 2(b) shows that PDGs are not accurate enough to differentiate the elements of composite structures. For instance, the whole statement in line 4 is represented by a single node, so it is not possible to distinguish the data structure $\{A, B\}$ nor its internal subexpressions. This can be solved by transforming the PDG into a CE-PDG. The transformation consists of three steps.

Step 1. The first step is to decompose all nodes that contain composite data structures so that each component is represented by an independent node. As in most ASTs, we represent data structures with a tree-like representation (similar to the one used in object-oriented programs to represent objects in calls (Liang and Harrold, 1998; Walkinshaw et al., 2003)). The decomposition of PDG nodes into CE-PDG nodes is straightforward from the AST. It is a recursive process that unfolds the composite structure by levels, i.e., if a subelement is another composite structure, it is recursively unfolded until the whole syntax structure is represented in the tree. The CE-PDG only unfolds data types to the level they are represented in the source code, thus unfolding is always finite (unlike atomization). In contrast to the PDG nodes (which represent complete statements), the nodes of this tree structure represent expressions. Therefore, we need a new kind of edge to connect these intra-statement nodes. We call these edges *structural edges* because they represent the syntactical structure of the program.

Definition 4 (Structural Edge). Let $G = (N, E)$ be a CE-PDG where N is the set of nodes and E is the set of edges. Given two CE-PDG nodes $n, n' \in N$, there exists a *structural edge* $n \rightarrow n'$ if and only if:

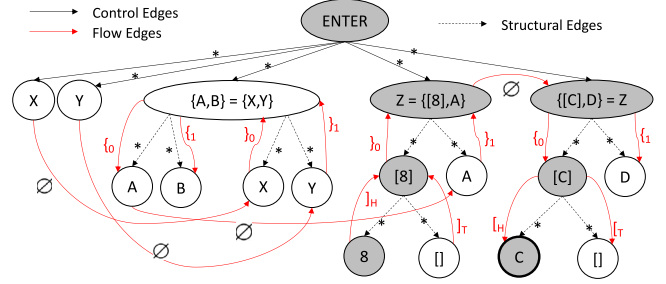


Fig. 4. CE-PDG of the code in Fig. 2.

- n contains a data structure for which n' is a subcomponent, and
- $\forall n'' \in N : n \rightarrow n' \wedge n' \rightarrow n'' \rightarrow n \not\rightarrow n''$.

Structural edges point to the components of a composite data structure, composing the inner skeleton of its abstract syntax tree. More precisely, each field in a data type is represented with a separate node that is connected to the PDG node that contains the composite data structure. For instance, the structural edges of the CE-PDG in Fig. 4 represent the tuples of the code in Fig. 2. The second condition of the definition enforces the tree structure as otherwise “transitive” edges could be established. For example, without the second condition a structural edge between $\{[C], D\} = Z$ and C could exist.

Step 2. The second step is to identify the flow dependencies that arise from the decomposition of the data structure. Clearly, the new nodes can be variables that flow-depend on other nodes, so we need to identify the flow dependencies that exist among the new (intra-statement) nodes. They can be classified according to two different scenarios: composite data structures being (i) defined and (ii) used. In Fig. 2 we have a definition (line 4), a use (line 3) and a definition and use in the same node (line 2). The explicit definition of a whole composite data structure (e.g., a tuple in the left-hand side of an assignment, see line 4) always defines every element inside it, so the values of all subelements depend on the structure that immediately contains them. Hence, the subexpressions depend on the structure being defined (i.e., flow edges follow the same direction as structural edges. See $\{[C], D\} = Z$ in Fig. 4). Conversely, the structure being used depends on its subexpressions (i.e., flow edges follow the opposite direction than structural edges. See $Z = \{[8], A\}$ in Fig. 4). Additionally, because the decomposition of nodes augments the precision of the graph, all flow edges that pointed to original PDG nodes that have been decomposed, now point to the corresponding node in the new tree structure. An example of a flow edge that has been moved due to the decomposition is the flow edge between the new A nodes. In the original PDG, this flow edge linked the nodes $\{A, B\} = \{X, Y\}$ and $Z = \{[8], A\}$.

Step 3. The last step to obtain the CE-PDG is labeling the edges with constraints that are later used during the slicing phase. The idea is that the slicing algorithm traverses the edges and collects the labels in a stack that is used to decide what edges should be traversed and what edges should be ignored. We call the new labeled edges *constrained edges* because the labels act as constraints for the graph traversal.

For the sake of simplicity, and without loss of generality, we distinguish between tuples and functional (algebraic) lists. The position in a tuple is indicated with an integer, while the position in a list is indicated with head (H) or tail (T). The case of objects, records, or any other structure can be trivially included by just specifying the position with the name of the field.

Definition 5 (Constraint). A constraint C is a label defined with the following grammar:

$$\begin{aligned} C &::= \emptyset \mid * \mid \text{Access} \\ \text{Access} &::= \text{Tuple} \mid \text{List} \\ \text{Tuple} &::= \{ \text{int} \mid \} \text{int} \\ \text{List} &::= [\text{pos} \mid] \text{pos} \\ \text{Pos} &::= H \mid T \end{aligned}$$

where int is a positive integer.

The meaning of each kind of constraint is the following:

Empty Constraint ($n \xrightarrow{\emptyset} n'$). It specifies that an edge can always be traversed by the slicing algorithm.

Asterisk Constraint ($n \xrightarrow{*} n'$). It also indicates that an edge can always be traversed; but it ignores all the collected restrictions so far, which means that going forward, the whole data structure is needed. This kind of constraint is the one used in control and structural edges, which are traversed ignoring the previous constraints collected.

Access Constraint ($n \xrightarrow{\text{op position}} n'$). It indicates that an element is the position -th component of another data structure that is a tuple if $\text{op} = \text{Tuple}$ or a list if $\text{op} = \text{List}$. op also indicates whether the element is being defined (“{”, “[”) or used (“}”, “]”).

Example 3 (Labeling Edges). All edges in Fig. 4 are labeled with constraints. Because B is the second element being defined in the tuple {A, B}, the constraint of the flow dependence edge that connects them is $\{_1$. Also, because 8 is the head in the list [8], the constraint of the flow dependence edge that connects them is $]_H$.

At this point, the reader can see that the constraints can be used to accurately slice the program in Fig. 2(a). In the CE-PDG (Fig. 4), the slicing criterion (C) is the head of a list (indicated by the constraint $]_H$), and this list is the first element of a tuple. When traversing backwards the flow dependencies, we do not want the whole Z, but only the head of its first element (i.e., the cumulated constraints $]_H \{_0$). Then, when we reach the definition of Z, we find two flow dependencies ([8] and A). But looking at their constraints, we exactly know that we want to traverse first $\}_0$ and then $]_H$ to reach the 8. So far, no structural edge is traversed during the slice in the above example. Structural edges represent the syntactical structure of composite data structures in trees. Therefore, after a structural edge has been traversed, no flow edges are allowed to be traversed (note that every structural edge has a corresponding flow edge in the same or opposite direction). The slice computed in this way is composed of the gray nodes, and it is exactly the minimal slice in Fig. 2(c).

The CE-PDG is a generalization of the PDG because the PDG is a CE-PDG where all edges are labeled with empty constraints (\emptyset). In contrast, all edges in the CE-PDG are labeled with different constraints:

- Structural and control edges are always labeled with asterisk constraints.
- Flow edges for definitions inside a data structure are labeled with opening ($\{, [$) access constraints.
- Flow edges for uses inside a data structure are labeled with closing ($\},]$) access constraints.
- The remaining data edges are labeled with empty constraints.

The behavior of access constraints and asterisk constraints in the graph traversal is further detailed in the next section, where we also formalize the slicing algorithm that performs the traversal of the CE-PDG.

$$P ::= CO$$

$$\begin{aligned} C &::= \{ \text{int} \mid \} \text{int} \mid \text{pos} \mid \emptyset \mid * \mid P \mid \epsilon & W &::= O' \\ R &::= \{ \text{int} \mid \} \text{int} \mid \text{pos} \mid \emptyset \mid R \mid \epsilon & O' &::= \{ \text{int} \mid \} \text{int} \mid \text{pos} \mid O' \mid \epsilon \\ O &::= \{ \text{int} \mid \} \text{int} \mid \text{pos} \mid \emptyset \mid O \mid * \mid P \mid \epsilon \end{aligned}$$

(a) Realizable paths grammar

(b) Stack words

Fig. 5. Grammars defining allowed constraints ($\text{pos} \in \{H, T\}$ and $\text{int} \in \mathbb{Z}$).

4.2. Constrained traversal

In this section, we show how constraints can improve the accuracy of the slices computed with the CE-PDG. In order to represent the paths of the CE-PDG that can be traversed, we use a grammar. The label of an edge can be seen as a terminal. Therefore, by traversing the edges we build words. But not all edges can be traversed; paths are only realizable when the word induced by the path belongs to a language for which the grammar is shown in Fig. 5(a). The realizable path grammar is modeled after the grammars by Repts et al. (1995). In this grammar, P is the initial symbol, C , R , and O represent sequences that contain closing, resolved, and opening constraints, respectively. \emptyset and $*$ stand for empty and asterisk constraints, respectively.

The key point of this grammar is *resolved* constraints. A resolved constraint is an opening constraint followed by the *complementary* closing constraint (e.g., $\{_2$ followed by $\}_2$). The paths of the CE-PDG that can be traversed are formed by any combination of closing constraints followed by opening constraints. Any number of empty constraints (\emptyset) can be placed along the path. On the other hand, asterisk constraints ($*$) always ignore any constraints already collected. Therefore, after traversing an asterisk constraint, the paths that can be traversed are the same as if no constraint was previously collected. The reason behind this behavior is later detailed in Example 6.

Example 4 (Traversing Constraints). Consider Fig. 4 again and the slicing criterion (C). To compute the slice we traverse edges backwards, tracing a path from C to 8 formed by the following sequence of constraints: $]_H \{_0 \emptyset \}_0$, which can be derived from the grammar in Fig. 5(a):

$$\begin{aligned} P &\xrightarrow{P \rightarrow CO} CO \xrightarrow{C \rightarrow \epsilon} O \xrightarrow{O \rightarrow RO} RO \xrightarrow{O \rightarrow \epsilon} R \xrightarrow{R \rightarrow [H]_H} [H]_H R \\ &\xrightarrow{R \rightarrow \{_0\}_0} [H]_H \{_0\}_0 R \xrightarrow{R \rightarrow \emptyset R} [H]_H \{_0\} \emptyset R \xrightarrow{R \rightarrow \epsilon} [H]_H \{_0\} \emptyset \}_0 R \end{aligned}$$

Opening access constraints can be seen as queries that can be solved along the path traversal. Let us explain this view with an example. Consider Fig. 2 again and assume that we are interested in the value of variable A in line 2. In this case, A is being defined. Since A is inside a particular position of the composite structure (position 0 inside the tuple {A, B}), it can only receive a value from the same position in an analogue composite structure. This may happen in the same statement or in any other previous one. For this reason, we metaphorically “open” a query to find the value of variable A, and indicate it with the corresponding structure and position symbols ($\{_0$). When we reach the statement level, we notice that the expression giving value to the whole data structure is an analogue data structure (the right-hand side of the equality, tuple {X, Y}). Thus, the expression that gives value to variable A must be inside this data structure. Because we are looking for a specific position with an open query ($\{_0$), we can choose the element we are interested in between the possible ones (X, reachable through $\}_0$), “closing” this query. As a result, we reach a state where the pending query has been “resolved”, and we can now focus on another open query, if any. To sum up, with this point of view, each variable definition contained in a composite structure is considered as the opening of a flow dependence query, and each variable use inside an analogous structure as the closing of this flow dependence query.

Table 1

Processing edges with a stack. x and y are positions (*int* or *H/T*). \emptyset and $*$ are empty and asterisk constraints, respectively. S is a stack, \perp the empty stack.

Case	Input stack	Edge constraint	Output stack
(1)	S	\emptyset	S
(2)	S	$\{x \text{ or } l_x$	$S\{x \text{ or } S\{x$
(3)	\perp	$\}_x \text{ or } l_x$	\perp
(4)	$S\{x \text{ or } S\{x$	$\}_x \text{ or } l_x$	S
(5)	$S\{x \text{ or } S\{x$	$\}_y \text{ or } l_y$	<i>error</i>
(6)	S	$*$	\perp

Then, every time a query is successfully closed we say that the opened query has been resolved.

To ensure that only realizable paths are visited, the slicing algorithm uses a stack to store the word generated by traversing the CE-PDG. When a node is selected as the slicing criterion, the algorithm starts from this node with an empty stack (\perp) and accumulates constraints with each edge traversed. Only opening constraints impose a restriction on the symbols that can be pushed onto the stack: when an opening constraint is on the top of the stack, the only closing constraint accepted to build a realizable word is its complementary closing constraint. Therefore, the only information necessary to determine whether an edge can be traversed is the sequence of non-resolved (opening) constraints at the top of the stack. They form the words that remain in the stack when a path is traversed (see the grammar in Fig. 5(b)).

Example 5 (Evolution of the Stack). To produce the derivation shown in Example 4 we start from the slicing criterion with an empty stack that is filled and emptied during the traversal:

$$\perp \xrightarrow{l_H} [l_H] \xrightarrow{l_0} [l_H l_0] \xrightarrow{\emptyset} [l_H l_0] \xrightarrow{l_0} [l_H l_0 l_0] \xrightarrow{l_H} \perp$$

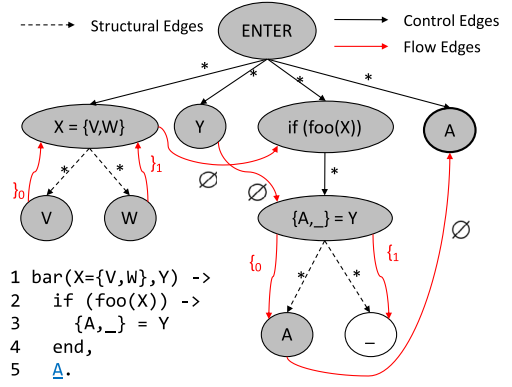
All words in the stack are only formed from opening access constraints, as defined by the grammar in Fig. 5(b).

Table 1 shows how the stack is updated in all possible situations. The constraints are collected or resolved depending on the last constraint added to the word (the one at the top of the *Input stack*) and the new one to be treated (column *Edge Constraint*). All cases shown in Table 1 can be summarized in four different situations:

- **Traverse constraint (cases 1 and 3):** The edge is traversed without modifying the stack.
- **Collect constraint (case 2):** The edge can be traversed by pushing the edge's constraint onto the stack.
- **Resolve constraint (cases 4 and 5):** There is an opening constraint at the top of the stack and an edge with a closing constraint that matches it (case 4), so the edge is traversed by popping the top of the stack; or they do not match (case 5), so the edge is not traversed.
- **Ignore constraints (case 6):** Traversing the edge empties the stack.

Previous examples have shown the behavior of access and empty constraints. Example 6 shows a program in which asterisk constraints are necessary to select all the nodes that are needed.

Example 6 (The Need for Asterisk Constraints). Consider the Erlang function and its associated CE-PDG in Fig. 6, where variable A in line 5 is the slicing criterion. Consider also the table in Fig. 6, where the most relevant steps of the backward traversal of the graph are shown. Each row represents the traversal of one edge, except the initial row, which represents the selection of the first node (the slicing criterion). In the table, column **Step** represents the number of the current step during the traversal. Different alternative paths are shown with letters (*a*, *b*, *c*). Columns **Start Node** and **Node Reached** represent the start and end



Step	Start Node	Edge	Node Reached	Stack
0	-	-	(5,A)	\perp
1	(5,A)	flow	\emptyset (3,A)	\perp
2	(3,A)	flow	$\{0\}$ (3,{A,_} = Y)	$\{0\}$
3a	(3,{A,_} = Y)	flow	\emptyset (1,Y)	$\{0\}$
4a	(1,Y)	control	$*$ (1,ENTER)	\perp
3b	(3,{A,_} = Y)	control	$*$ (2, if(foo(X)))	\perp
4b	(2, if(foo(X)))	flow	\emptyset (1,X={V,W})	\perp
5b	(1,X={V,W})	flow	$\{0\}$ (1,V)	\perp
5c	(1,X={V,W})	flow	$\{1\}$ (1,W)	\perp

Fig. 6. Erlang function, associated CE-PDG, and slice step by step.

nodes (*line,expression*) when traversing the current edge, **Edge** represents the type and constraint of the traversed edge; and **Stack** represents the stack computed after traversing the edge.

After reaching $\{A, _ \} = Y$ (step 2), the stack contains the opening constraint $\{0\}$, and there are two possible paths: (a) a flow path to the parameter variable Y (step 3a), and (b) a control path to the *if* condition (step 3b). Let us focus on the second path to show the necessity of asterisk constraints. When we traverse the control edge, all the constraints stacked due to the traversal of previous flow edges must be dropped from the stack (case 6 in Table 1). The reason is simple: when we reach a statement by a control edge, we are no longer interested in the value of the uses of variables that the traversal has accumulated in the stack, but in the value of the variables used in this controller statement. The fact is that keeping the previous stack constraints may result in erroneous slices. For instance, consider a scenario where, if we do not empty the stack in step 3b, we would reach the $X=\{V,W\}$ statement with the stack $\{0\}$, and the traversal would only reach the first element of the tuple (V) traversing $\{0\}$. Therefore, W would be never included in the slice because it can only be reached traversing the constraint $\{1\}$ that does not match the constraint of the stack. In contrast, emptying the stack in step 3b when traversing the control edge forces the slice to correctly include both V and W . Note also that the constraint $\{0\}$ collected in step 2 is not entirely useless. It is still used in the flow path to Y (we only want the first component of Y).

4.3. Subsumed constraints

The inclusion of constraints in the PDG's edges induces interesting properties that we can take advantage of. An important property is that the constraints of some stacks subsume the constraints of other stacks. The subsumed constraints can be ignored, thus, improving the efficiency and the precision in the slicing traversal.

When we reach a node of a data structure with a non-empty stack, we are explicitly requiring a part of the data structure represented by the stack. For instance, if we reach a list with the constraint $[_T]$ we require the *whole* tail of the list, but not the head. Therefore, if we later reach the same data structure with a stack $[_H l_T]$, which requires only

the second element of the list (the head of the tail), we can ignore the second stack because it is subsumed by the previous one.

In general, the number of cumulated constraints in a stack when we reach a data structure indicates the depth where the element that we are searching for is located.

Example 7. Consider the tuple $\{\{1, 2\}, \{3, 4\}\}, \{\{5, 6\}, \{7, 8\}\}$, that we reach with the following four stacks:

- \perp : If the stack is empty, we need all the elements of the data structure.
- $\{0\}$: We need the whole tuple in position 0: $\{\{1, 2\}, \{3, 4\}\}$.
- $\{1\}$: We need the whole tuple in position 1: $\{\{5, 6\}, \{7, 8\}\}$.
- $\{1, \{0\}$: We need a part of element 0, but not all. We only need the element in position 1: $\{3, 4\}$.

The following inter-relations exist: $\{0\} \subseteq \perp$, $\{1\} \subseteq \perp$, $\{1, \{0\} \subseteq \perp$, and $\{1, \{0\} \subseteq \{1, \{0\}$.

The above example reveals an important property: if we reach a node with two different stacks and *the first is a suffix of the second*, then the second stack can be ignored. This relationship imposes an order among stacks.

Definition 6 (Stack Ordering). Given two stacks S and S' , $S' \subseteq S$ if and only if S is a suffix of S' .

This ordering is a partial order because it is reflexive ($S \subseteq S$), antisymmetric ($S_1 \subseteq S_2 \wedge S_2 \subseteq S_1 \iff S_1 = S_2$) and transitive ($S_1 \subseteq S_2 \wedge S_2 \subseteq S_3 \implies S_1 \subseteq S_3$).

We can use this property to sometimes stop the traversal of the CE-PDG: The traversal is stopped when a node is reached with a stack, and this node has been already reached by another previous stack that is a suffix of the new one. Note that this includes the case where a stack is empty because the empty stack is always a suffix of any other stack. Hence, when a node is traversed with the empty stack, it must not be traversed again.

4.4. The slicing algorithm

Algorithm 1 illustrates the process to slice the CE-PDG. It works similarly to the standard algorithm (Reps et al., 1994), traversing backwards all edges from the slicing criterion and collecting nodes to form the final slice. The algorithm uses a work list with the states that must be processed. A state represents the (backward) traversal of an edge. It includes the node reached, the current stack, and the sequence of already traversed edges (line 7). In every iteration, the algorithm processes one state. First, it collects all edges that target the current node (function `GETINCOMINGEDGES` in line 8). If the previously traversed edge is structural, we avoid traversing flow edges (lines 10–12) and only traverse structural or control dependence edges. The reason for this is that structural edges are only traversed to collect the structure of a data type so that the final slice is syntactically correct (for instance, to collect the tuple to which an element belongs). Flow edges are not further traversed to avoid collecting irrelevant dependencies of the structural parent. The function `PROCESSCONSTRAINT` checks the existence of a loop (reaching an already traversed edge) during the slicing traversal and implements Table 1 to produce the new stack generated by traversing the edge to the next node (line 13). If the edge cannot be traversed according to Table 1 (line 14), or because the new stack is subsumed (see Section 4.3) by one in the previously processed states or the work list (line 16), then the reachable node is ignored. Otherwise, the node is stored together with the new stack (line 18). Finally, the state is added to a list of processed states, used to avoid multiple evaluations of the same state (line 19), the work list is updated (removing any element that would be subsumed by new items — lines 20–22, and then adding

Algorithm 1 Intraprocedural slicing algorithm for CE-PDGs

Input: The slicing criterion node n_{sc} .

Output: The set of nodes that compose the slice.

```

1: function SLICINGALGORITHMINTRA( $n_{sc}$ )
2:    $slice \leftarrow \emptyset$ ;  $processed \leftarrow \emptyset$ 
3:    $workList \leftarrow \{(n_{sc}, \perp, [])\}$ 
4:   while  $workList \neq \emptyset$  do
5:     select some  $state \in workList$ ;
6:      $newItems \leftarrow \{\}$ 
7:      $\langle node, stack, traversedEdges \rangle \leftarrow state$ 
8:     for all  $edge \in GETINCOMINGEDGES(node)$  do
9:        $\langle node', type, \_ \rangle \leftarrow edge$ 
10:      if GETLASTEDGETYPE(traversedEdges) = structural then
11:        if  $type = flow$  then
12:          continue for all
13:         $newStack \leftarrow PROCESSCONSTRAINT(stack, edge, traversedEdges)$ 
14:        if  $newStack \neq error$  then
15:          if  $\nexists (node', s, t) \in processed \cup workList$ 
16:            |  $s$  is suffix of  $newStack$  then
17:             $w \leftarrow \langle node', newStack, traversedEdges \cup \{edge\} \rangle$ 
18:             $works \leftarrow works \cup \{w\}$ 
19:           $processed \leftarrow processed \cup \{state\}$ 
20:          for all  $(n, s, t) \in works$  do
21:             $workList \leftarrow workList \setminus \{(n', s', t') \in workList \mid$ 
22:               $n = n' \wedge s \text{ is suffix of } s'\}$ 
23:           $workList \leftarrow workList \cup works$ 
24:           $slice \leftarrow slice \cup \{node\}$ 
25:      return  $slice$ 

26: function PROCESSCONSTRAINT( $stack, edge, traversedEdges$ )
27:    $\langle \_, \_, constraint \rangle \leftarrow edge$ 
28:    $\langle op, position \rangle \leftarrow constraint$ 
29:   if  $constraint = AsteriskConstraint$  then ▷ Case 6
30:     return  $\perp$ 
31:   if  $edge \in traversedEdges$  then ▷ Check cycles
32:      $loop \leftarrow FINDLOOP(traversedEdges)$  ▷ Check loops
33:     if  $loop \neq error \wedge ISINCREASINGLOOP(loop, edge)$  then
34:       return  $\perp$  ▷ Inc. loop replaced by * constraint
35:   if  $constraint = EmptyConstraint$  then ▷ Case 1
36:     return  $stack$ 
37:   else if  $op = \{ \vee op = \}$  then ▷ Case 2
38:     PUSH(constraint, stack)
39:   else if  $stack = \perp$  then ▷ Case 3
40:     return  $\perp$ 
41:   else if  $op = \}$   $\wedge TOP(stack) = \{\langle \_, position \rangle\}$  then ▷ Case 4
42:     POP(stack)
43:   else if  $op = \}$   $\wedge TOP(stack) = \{\langle \_, position \rangle\}$  then ▷ Case 4
44:     POP(stack)
45:   else ▷ Non-matching closing constraint ▷ Case 5
46:     return  $error$ 
47:   return  $stack$ 

```

said items — line 23), and finally the current node is included in the slice (line 24).

The function `PROCESSCONSTRAINT` computes a new stack for all possible types of constraint, taking into account Table 1 and Section 4.5 (loop detection and avoidance). It checks first for an asterisk constraint and returns an empty stack (line 29). Then, the condition in line 31 checks the existence of a cycle (reaching an already traversed edge) during the slicing traversal. The function `FINDLOOP` (line 32) returns the shortest suffix of the sequence of traversed edges that form the last loop, while function `ISINCREASINGLOOP` (line 33), whose rationale is extensively explained in Section 4.5, consequently empties the stack when needed. If no dangerous loop is detected, the function checks for each of the cases in Table 1, pushing or popping elements from the stack accordingly (lines 35–47).

Example 8 (Applying Algorithm 1). Consider function `foo` in the code of Fig. 2(a) again, the selected slicing criterion ($\langle 4, C \rangle$), and its CE-PDG, shown in Fig. 4. The slicing process starts from the node that represents the slicing criterion (the expanded representation of the CE-PDG allows us to select C , the bold node, inside the tuple structure, excluding the rest of the tuple elements). Algorithm 1 starts the traversal of the graph with an empty stack (\perp). The evolution of the stack after traversing each flow edge is the following:

$$\perp \xrightarrow{l_H} [H] \xrightarrow{l_0} [H] \xrightarrow{\emptyset} [H] \xrightarrow{l_0} [H] \xrightarrow{l_H} \perp$$

Due to the traversal limitations imposed by row 5 in Table 1, node A is never included in the slice because the following transition is not possible: $[H] \xrightarrow{l_1} \text{error}$.

The slicing algorithm will also traverse the structural edges reaching the traversed nodes and generate new states in the worklist with empty stacks due to the asterisk constraint, however, no flow dependence edge is traversed after a structural edge (lines 10–12) and therefore despite the node $Z = \{[8], A\}$ being encountered with an empty stack, the flow edge to A is not traversed.

As already noted, the resulting slice provided by Algorithm 1 is exactly the minimal slice shown in Fig. 2(c).

4.5. Dealing with loops

In static slicing we rarely know the values of variables (they often depend on dynamic information), so we cannot know how many iterations will be performed in a program loop¹ (see the programs in Fig. 7, where the value of `max` is unknown). For the sake of completeness, we must consider any number of iterations, thus program loops are often seen as potentially infinite. Program loops produce cycles in the PDG. Fortunately, the traversal of cycles in the PDG is not a problem, since every node is only visited once. In contrast, the traversal of a cycle in the CE-PDG could produce a situation in which the stack grows infinitely (see Fig. 7(d)² and its CE-PDG in Fig. 9), generating an infinite number of states. Fortunately, not all cycles produce this problem³: To keep the discussion precise, we need to formally define when a cycle in the CE-PDG is a *loop*.

Definition 7 (Loop). A cyclic flow dependence path $P = n_1 \xleftarrow{c_1} n_2 \dots \xleftarrow{c_n} n_1$ is a *loop* if P can be traversed more than once with an initially empty stack (\perp) following the rules of Table 1.

Example 9 (Cycles Vs. Loops). Cycles that are not loops are not dangerous because the cycle's edges constraints prevent us to traverse them infinitely. The code in Fig. 7(a) contains a cycle that is not a loop because it can be traversed only once with an empty stack. This can be observed in Fig. 8, where we have the flow dependence cycle: $(6, x) \xleftarrow{l_0} (6, a) \xleftarrow{\emptyset} (5, a) \xleftarrow{l_0} (5, d) \xleftarrow{\emptyset} (4, d) \xleftarrow{l_1} (4, x) \xleftarrow{\emptyset} (6, x)$. But this is not a loop because no matter the stack we enter the cycle with, when l_1 is pushed on the stack, the cycle cannot be entered again, as the l_0 constraint cannot be matched against the top of the stack.

There exist three kinds of loops:

```

1 read(max);
2 x = init_tuple();
3 for(int i=0; i<max; i++){
4   {e,d} = x;
5   {a,b} = d;
6   x = {a,5};
7 }
8 {c,d} = x;
9 print(c);

```

(a) A cycle but not a loop

```

1 read(max);
2 read(b);
3 x = init_tuple();
4 for(int i=0; i<max; i++){
5   a = {x,i};
6   x = {a,b};
7 }
8 {c,d} = x;
9 print(c);

```

(b) Loop: decreasing stack size

```

1 read(max);
2 x = init_tuple();
3 for(int i=0; i<max; i++){
4   {a,b} = x;
5   x = {a+i,b+i};
6 }
7 {c,d} = x;
8 print(c);

```

(c) Loop: same stack size

```

1 read(max);
2 x = init_tuple();
3 for(int i=0; i<max; i++){
4   {e,d} = x;
5   {a,b} = d;
6   x = {a,b};
7 }
8 {c,d} = x;
9 print(c);

```

(d) Loop: increasing stack size

Fig. 7. Slicing flow-dependence cycles in the CE-PDG (slicing criterion underlined, slice in black).

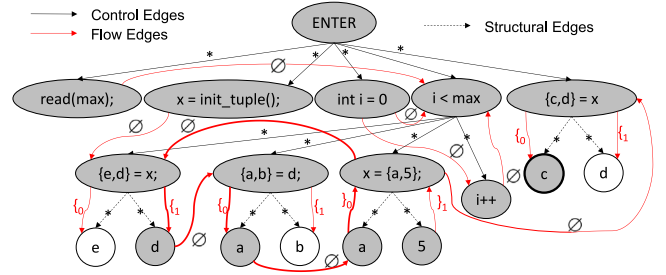


Fig. 8. CE-PDG of Fig. 7(a). The cycle is represented with bold edges.

1. Loops that decrease the size of the stack in each iteration (Fig. 7(b)) can only produce a finite number of states because the stack will eventually become empty. Such loops can be traversed collecting the elements specified by the stack, without a loss of precision.
2. Loops that maintain the stack's size constant in each iteration (Fig. 7(c)) are also not a problem because traversing the loop multiple times does not generate new states. Again, they can be traversed as many times as required by the stack, without a loss of precision.
3. Loops that increase the size of the stack in each iteration (Fig. 7(d)) produce an infinite number of states because the stack grows infinitely. This can be observed in the CE-PDG (Fig. 9) of the code in Fig. 7(d), where we can see the increasing loop: $(6, x) \xleftarrow{l_1} (6, b) \xleftarrow{\emptyset} (5, b) \xleftarrow{l_1} (5, d) \xleftarrow{\emptyset} (4, d) \xleftarrow{l_1} (4, x) \xleftarrow{\emptyset} (6, x)$

We formally define a special kind of loop which is the only one with potential danger: the *increasing loop*.

Definition 8 (Increasing Loop). A loop L is an *increasing loop* if the number of opening constraints along L is greater than the number of closing constraints.

Identifying increasing loops. To detect the increasing loops (those that can grow the stack infinitely) we have designed the pushdown automaton (PDA) shown in Fig. 10. The input of this automaton is the sequence of constraints that form a dependence cycle. The PDA tracks

¹ Note the careful wording in this section, where we distinguish between “program loops” (while, for...), “cycles” (paths in the PDG that repeat a node), and “loops” (repeated sequence of nodes during the graph traversal).

² It is easier to see how the stack changes by reading the code backwards from the slicing criterion.

³ The interested reader has a developed example for each kind of loop, which includes their CE-PDGs, in the technical report at <https://mist.dsic.upv.es/techreports/2023/02/case-studies-field-sensitive.pdf>.

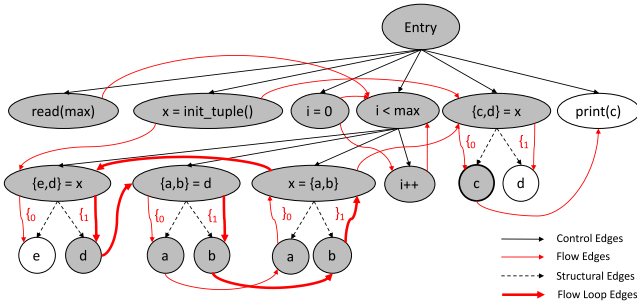
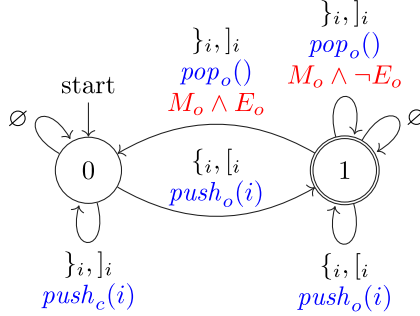


Fig. 9. CE-PDG of Fig. 7(d). The increasing loop is represented with bold edges.



$$M_o \equiv \text{top}(S_o) = i$$

$$E_o \equiv \text{len}(S_o) = 1$$

Fig. 10. Pushdown automaton to recognize increasing loops.

two stacks (for closing and opening constraints) across two states. Initial state 0 represents the case where all opening constraints of the sequence are balanced by the corresponding closing constraint. When a closing constraint is reached, the PDA pushes the constraint into the closing stack (push_c). When an opening constraint is processed, the PDA pushes the opening constraint into the opening stack (push_o) and moves to state 1. Final state 1 represents the case where an opening constraint has been processed but not balanced yet. In state 1, when a closing constraint that matches a previous opening constraint (condition M_o) is processed, we pop the opening constraint from the stack (pop_o). If the popped element of the opening stack is the last element of the stack (condition E_o), the PDA returns to state 0. Finally, if a path is accepted by this automaton, the path forms an increasing loop if and only if:

the reversed stack S_c is a prefix of S_o and they are not equal.

The rationale of this condition is that it ensures that, in each iteration, there are more opening constraints (those in S_o) than closing constraints (those in S_c), and all the closing constraints close some but not all opening constraints (because they are a prefix), thus the number of opening constraints grows infinitely. Note that asterisk constraints (*) do not appear in the PDA because they cannot appear in an increasing loop (an asterisk constraint empties the stack and thus the same state would be repeated).

Example 10 (Detecting an Increasing Loop). Consider the dependence cycle formed from lines 4, 5, and 6 of Fig. 7(d) (see its CE-PDG in Fig. 9): $(6, x) \xrightarrow{1} (6, b) \xrightarrow{\emptyset} (5, b) \xrightarrow{1} (5, d) \xrightarrow{\emptyset} (4, d) \xrightarrow{1} (4, x) \xrightarrow{\emptyset} (6, x)$, which contains the word: $\{1\}\emptyset\{1\}\emptyset\{1\}\emptyset$.

Now, if we parse this word with the PDA we produce the sequence of states shown in Fig. 11. The final state is an accepting state, and the

reverse of S_c (1) is a prefix of S_o (1, 1) (but they are not equal), so this path corresponds to an increasing loop. Moreover, the PDA also detects that this loop adds $\{1$ (the remainder of S_o once the prefix is removed) to the stack in every iteration.

More formally, an increasing loop $n_1 \xleftarrow{C_1} n_2 \xleftarrow{C_2} \dots \xleftarrow{C_n} n_1$ can be identified because $C_1 C_2 \dots C_n$ belongs to the language induced by the PDA and the two final stacks computed with the PDA (S_c and S_o) satisfy that $\text{reverse}(S_c)$ is a prefix of S_o and $\text{reverse}(S_c) \neq S_o$.

Detecting loops. When a loop is detected, the PDA can determine whether it is increasing. Fortunately, it is not necessary to preprocess the CE-PDG to detect the loops. The loops can be detected at slicing time during the traversal (with Algorithm 1). The strategy proposed in Galindo et al. (2022) to detect a loop and prevent traversing it infinitely was:

1. To traverse the CE-PDG backwards following all the edges until we traverse the same edge twice.
2. If the same node is reached with the same stack twice, then the traversal of this node is stopped.
3. Otherwise, use the PDA to determine whether the loop is increasing.
4. If the loop is increasing, the constraint of the traversed edge is considered as an asterisk constraint (which empties the stack). This ensures that a loop is only traversed at most twice (the second time that the stack is emptied, the state is repeated) and slicing is ensured to terminate.

The new algorithm (Algorithm 1) modifies point 2: the traversal of a node is stopped if this node has been previously reached *with a stack that is a suffix of the current stack* (see Section 4.3). This change significantly reduces the number of times that the PDA is triggered as can be seen in Example 11.

Example 11 (Increased Efficiency of Algorithm 1). Consider the CE-PDG in Fig. 13 and let c be the slicing criterion. A traversal with the approach of Galindo et al. (2022) is shown in the following table:

Step	Start node	Edge	Node reached	Stack
0	—	—	c	\perp
1	c	flow $\{0$	$\{c,d\}=x$	$\{0$
2	$\{c,d\}=x$	flow \emptyset	x	$\{0$
3	x	flow $\{0$	$\{x,a\}=b$	$\{0\{0$
4	$\{x,a\}=b$	flow \emptyset	$b=x$	$\{0\{0$
5	$b=x$	flow \emptyset	x	$\{0\{0$
6	x	flow $\{0 \rightsquigarrow *$	$\{x,a\}=b$	\perp
7	$\{x,a\}=b$	flow $\emptyset \rightsquigarrow *$	$b=x$	\perp
8	$b=x$	flow $\emptyset \rightsquigarrow *$	x	\perp
9	x	flow $\{0 \rightsquigarrow *$	$\{x,a\}=b$	\perp

In step 6, the edge $x \leftarrow \{x, a\} = b$ is visited for the second time and, thus, the PDA is triggered to determine whether the loop is increasing, which it is. Therefore, the edge's constraint is ignored, and an asterisk constraint is applied instead (emptying the stack). The same happens in steps 7, 8, and 9. In step 9, for the first time, a node $\{x, a\} = b$ is visited twice with the same stack and, thus, the traversal terminates.

In contrast, if we use Algorithm 1, steps 6–9 are never executed; and the PDA is never triggered. This happens because, in step 5, we visit node x for the second time with the stack $\{0\{0$. The first visit to x (step 2) was done with the stack $\{0$, which is a suffix of the new stack. Therefore, the traversal terminates at step 5.

Skipping steps 6–9 also skip the traversal of $x = \text{init_tuple}()$ with the empty stack (see step 7) and subsequent traversals from this node with the empty stack.

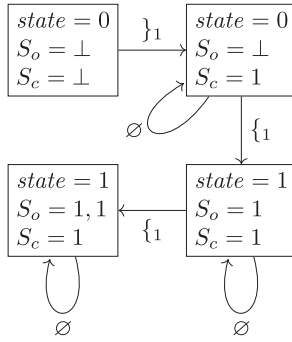


Fig. 11. States produced by the PDA in Fig. 10 with the word $\}_1\emptyset\{1\emptyset\}_1\emptyset$.

Not only the efficiency, but also the precision of Algorithm 1 is augmented with respect to Galindo et al. (2022). This is explained in Example 12.

Example 12 (Increased Precision of Algorithm 1). Continuing Example 11, after step 4, where we reach node $b = x$ with the stack $\{0\}_0$, we can exit the loop and reach node $x = \text{init_tuple}()$ with the stack $\{0\}_0$. This means that we only want the first element of the tuple assigned to x . Nevertheless, with the approach of Galindo et al. (2022), after step 7, we can exit the loop again and reach the same node ($x = \text{init_tuple}()$) with a different stack (the empty stack). Unfortunately, this means that we want the whole tuple, which is clearly imprecise.

Only increasing loops can produce non-termination. For this reason, Algorithm 1 detects loops (lines 31–32) and checks whether they are increasing with function `isIncreasingLoop` (line 33). This function uses the PDA of Fig. 10 to determine whether the loop is increasing and in such a case the stack is emptied (line 34), i.e., the traversal continues unconstrained.

At this point, the reader might be wondering whether the PDA is actually needed, because Example 12 showed that even with an increasing loop, the PDA was not used. However, with some loops, the PDA is necessary. This is illustrated in Example 13.

Example 13 (Necessity of the PDA). Consider again the CE-PDG in Fig. 13 and this time let d be the slicing criterion. A traversal with Algorithm 1 is shown in the following table:

Step	Start node	Edge	Node reached	Stack
0	–	–	(d)	\perp
1	(d)	flow $\{_1$	$(\{c,d\}=x)$	$\{_1$
2	$(\{c,d\}=x)$	flow \emptyset	(x)	$\{_1$
3	(x)	flow $\{_0$	$(\{x,a\}=b)$	$\{_1\{_0$
4	$(\{x,a\}=b)$	flow \emptyset	$(b=x)$	$\{_1\{_0$
5	$(b=x)$	flow \emptyset	(x)	$\{_1\{_0$
6	(x)	flow $\{_0 \rightsquigarrow *$	$(\{x,a\}=b)$	\perp
7	$(\{x,a\}=b)$	flow $\emptyset \rightsquigarrow *$	$(b=x)$	\perp
8	$(b=x)$	flow $\emptyset \rightsquigarrow *$	(x)	\perp
9	(x)	flow $\{_0 \rightsquigarrow *$	$(\{x,a\}=b)$	\perp

In this case, the algorithm would never stop the traversal because there is never a previous stack that is a suffix of the current stack. To stop the traversal we need the PDA. In step 6, we visit for the first time a previously visited edge (we find a cycle). At this point, the PDA is triggered and it determines that this is an increasing loop. Therefore, the PDA forces the change of the loop edges' constraints to $*$. This makes us stop the traversal at step 9 because the same node $(\{x,a\}=b)$ is visited twice with the same stack (\perp) (one is a suffix of the other).

```

1 ...
2 for(int i=0; i<max; i++){
3   b = x;
4   {x,a} = b;
5 }

```

(a) Increasing loop

	$b = x;$	$b_1 = x_1;$	$b_1 = x_1;$
	$b = x;$	$b_2 = x_2;$	$b_2 = x_2;$
$\{x,a\} = b;$	$x = b_1;$	$x = b_1;$	$x_1 = b_{11};$
$a = b_2;$	$a = b_2;$	$a = b_2;$	$x_2 = b_{12};$
			$a = b_2;$

(b) The first steps of atomization for the increasing loop

Fig. 12. Infinite unfolding of atomization.

Theorem 1 ensures termination of the whole slicing process.

Theorem 1 (Termination of Slicing). Let $G = (N, E)$ be a CE-PDG and let $n_{sc} \in N$ be a slicing criterion for G . Algorithm 1 terminates when it slices G with respect to n_{sc} .

The proof of this theorem can be found in Section 4.7.

The reader could think that it would be a good idea to identify all increasing loops at CE-PDG construction time. Unfortunately, finding all cycles has an average complexity $\mathcal{O}(N^2EL)$, where L is the number of cycles. The maximum number of cycles in a graph is exponential ($L = 2^N$), and thus the worst-case complexity is $\mathcal{O}(2^N)$ (Gongye et al., 2020). Our approach avoids the problem of finding all loops. We just treat them on demand, when they are found by the slicing algorithm (i.e., we do not search for loops, we just find them during the CE-PDG traversal). So we only process those loops found in the slicing process; and processing a loop has a linear cost (in the worst case $\mathcal{O}(N)$, if the loop includes all program statements).

4.6. Recursive data structures: constrained traversal versus atomization

In this section we show that our technique is not only an alternative to atomization but a fundamental improvement that solves an important problem that cannot be solved with atomization. Atomization cannot handle recursive data types such as lists, trees, linked lists (often implemented via a node that contains a value and a reference to another node), etc. The fundamental problem is that atomization would need to infinitely unroll the recursive data type.

Example 14 (Recursive Unfolding Cannot Be Atomized). Consider the program in Fig. 12(a), in which a recursive tuple x is unfolded in a loop. The inside of the loop cannot be atomized, as the loop contains two instructions: $b = x$ and $\{x,a\} = b$. Atomization would convert structures to multiple assignments. However, the depth of x and thus b is unknown.

Consider the atomization sequence shown in Fig. 12: A first step (second column) is to split the pattern matching assignment into two separate assignments for x and a . However, b_1 and b_2 are not defined, so we must split b 's assignment. The result (third column) still has undefined variables (x_1 and x_2). Defining those would require splitting b_1 (last column), resulting in a situation analogous to the first split, starting an infinite loop.

In contrast, the CE-PDG represents each field explicitly, performing the unfolding process as many times as required during the slicing traversal and thus, yielding to the correct slice, as can be seen in Fig. 13.

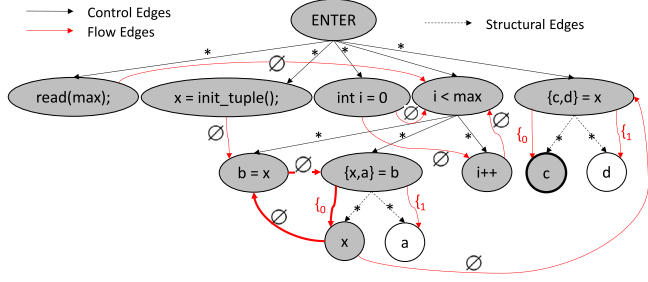


Fig. 13. The CE-PDG for a program (Fig. 12(a)) with recursively unfolding data structures.

4.7. Proof of Algorithm 1's termination

In this section, we prove [Theorem 1](#) by showing that the traversal of all non-increasing loops is terminating with any initial stack. And the same happens with all loops with an asterisk constraint. Because emptying the stack is equivalent to traversing an asterisk constraint, then emptying the stack of increasing loops is enough to ensure that their traversal is also terminating with any initial stack.

In [Section 4.5](#) we define loops and increasing loops ([Definitions 7 and 8](#)), the remaining kinds of loops can be defined in a similar fashion:

Definition 9 (Decreasing Loop). A loop L is a *decreasing loop* if the number of closing constraints along L is greater than the number of opening constraints.

Definition 10 (Balanced Loop). A loop L is a *balanced loop* if the number of closing and openings constraints along L is the same.

It is important to note here that, according to [Definitions 8, 9, and 10](#), a loop can only be decreasing, balanced, or increasing.

We also need to prove some properties of a loop that are needed to prove the theorem. These properties are captured by the following lemmas.

Lemma 1. Let L be a non-increasing loop and S be a stack. The traversal of L with S is terminating.

Proof. Non-increasing loops can be decreasing loops ([Definition 9](#)) or balanced loops ([Definition 10](#)). We prove each case separately.

Decreasing loops. First of all, for any initial stack, a first traversal of all the edges of L must be possible to explore this scenario. Otherwise, there cannot be a traversal loop due to unmatched constraints. Moreover, according to [Definition 9](#) the path of L contains more closing constraints than opening constraints. When slicing the graph, and before iterating into a decreasing loop, there are two stack possibilities: the empty stack or a stack with a sequence of opening constraints at the top (see [Fig. 5\(b\)](#)):

Empty stack (\perp). Since the constraint in L can be in any order, we need to consider two different scenarios:

1. There is no suffix of L where the number of opening constraint is greater than the number of closing constraints (e.g., $\{\{0, \{0, \}0, \}0, \}0\}$). Since L is a decreasing loop, in this scenario we can summarize L as a sequence of n closing constraints. In this case, according to case 3 of [Table 1](#), pushing any number j of closing constraint to an empty stack results in the empty stack.

$$\perp \xrightarrow{\{l_i, \forall i \in 1..j\}^*} \perp \quad \text{and} \quad \perp \xrightarrow{\{l_i, \forall i \in 1..j\}^*} \perp$$

Then, the same node is reached twice with the same stack (\perp) and a second traversal is not processed as indicated by line 22 in [Algorithm 1](#).

2. There is a suffix of L where the number of opening constraints is greater than the number of closing constraints (e.g., $\{\{0, \}0, \}0, \{0, \}0\}$). In this scenario, there are always closing constraints that are consumed by the initial empty stack (\perp) according to case 3 of [Table 1](#). On the other hand, the final part of L pushes more opening constraints than closing constraints and the stack calculated after the first traversal of the loop is not empty. As L is a decreasing loop, all the constraints pushed in the first traversal are consumed at the beginning of the second traversal. Finally, the final part of L generates the same sequence of positive constraints at finishing the second traversal. Therefore, as the same node is reached twice with the same stack, the traversal stops.

$$\perp \xrightarrow{L} S \xrightarrow{L^*} S$$

Stack with a sequence of opening constraints ($\{C_{O_1}, \dots, C_{O_n}\}$).

Cases 4 and 5 of [Table 1](#) represent the two possibilities that can occur when pushing a closing constraint to a non-empty stack: the closing constraint balances the opening constraint at the top of the stack or it fails.

- According to case 5, if the closing constraint does not balance the opening constraint at the top of the stack, the traversal is aborted by an error. Hence, the traversal is finite.

$$[\dots, \{l_i\}] \xrightarrow{\{l_j, \forall i \neq j\}} \text{error} \quad \text{and} \quad [\dots, \{l_i\}] \xrightarrow{\{l_j, \forall i \neq j\}} \text{error}$$

- On the contrary, case 4 shows that if the closing constraint balances the opening constraint at the top of the stack, the constraint is popped and the traversal continues. Since the number of elements of the stack is finite (an infinite stack would have been forever growing in an increasing loop), popping each opening constraint on the top of the stack with a closing constraint of the loop sequence will result in an empty stack at some point of the traversal. Note that this point of the traversal is not necessarily after traversing the whole loop L , but can be after traversing any intermediate edge of the loop. As the processing of any decreasing loop L with an empty stack (\perp) has already been proved terminating, then we can state that the traversal is finite.

$$[\{l_1, \dots, \{l_n\}\}] \xrightarrow{L^m} S \xrightarrow{L^*} S$$

Balanced loops. With the same reasoning as in the decreasing loops, we differentiate two scenarios according to the initial form of the stack:

Empty stack (\perp). Since the constraint in L can be in any order, we need to consider the same two scenarios we did before:

- There is no suffix of L where the number of opening constraint is greater than the number of closing constraints (e.g., $\{\{0, \{0, \}0, \}0, \}0\}$). Since L is a balanced loop, in this scenario we can summarize L as an empty sequence of constraints. The result of iterating into a balanced loop of this type is the same stack, since each iteration of the loop itself balances the opening constraints with their complementary closing constraints. Hence, the initial node of the graph will be reached twice with the same stack, and the traversal will be stopped as indicated by line 22 in [Algorithm 1](#).

$$\perp \xrightarrow{\emptyset^*} \perp$$

- There is a suffix of L where the number of opening constraints is greater than the number of closing constraints (e.g., $\{\{0, \}0, \}0, \{0, \}0\}$). In this scenario, there are always closing constraints that are consumed by the initial empty stack (\perp)

according to case 3 of Table 1. On the other hand, the final part of L pushes more opening constraints than closing constraints and the stack calculated after the first traversal of the loop is not empty. As L is a balanced loop, all the constraints pushed in the first traversal are consumed at the beginning of the second traversal. Finally, the final part of L generates the same stack when the traversal of L finishes. Therefore, as the same node is reached twice with the same stack, the traversal stops.

$$\perp \xrightarrow{L} S \xrightarrow{L^*} S$$

Stack with a sequence of opening constraints ($[C_{O_1}, \dots, C_{O_n}]$).

Cases 4 and 5 of Table 1 represent the two possibilities that can occur when pushing a closing constraint to a non-empty stack: the closing constraint balances the opening constraint at the top of the stack or it fails.

- According to case 5, if the closing constraint does not balance the opening constraint at the top of the stack, the traversal is aborted by an error. Hence, the traversal is finite.

$$[\dots, \{i\}] \xrightarrow{\{j, \forall i \neq j\}} \text{error} \text{ and } [\dots, \{i\}] \xrightarrow{\{j, \forall i \neq j\}} \text{error}$$

- On the contrary, case 4 shows that if the closing constraint balances the opening constraint at the top of the stack, the constraint is popped and the traversal continues. Since the number of opening and closing constraints in L is the same, if any constraint from the stack is consumed by the initial closing constraints it will be restored later by the corresponding opening constraint contained in L . Then, the stack obtained after traversing L will be the same. As a result, since the same node is reached twice with the same stack, the traversal finishes.

$$[\{1, \dots, \{n\}] \xrightarrow{L^*} [\{1, \dots, \{n\}]$$

Finally, note that if we consume every opening constraint during the traversal of the closing constraints at the beginning of L , we will find ourselves in the previous scenario where the stack was empty and, thus, the traversal finishes too. \square

Lemma 2. Given a CE-PDG increasing loop L . There exists a stack S for which it is possible to infinitely traverse L with S .

Proof. Following the same reasoning that we did in some particular decreasing loops, each iteration of an increasing loop can be summarized as a sequence of opening constraints of the form $[\{1, \dots, \{n\}]$ independently of the sequence of constraints. Case 2 of Table 1 shows that opening constraints can always be traversed independently of the top of the stack. For this reason, the loop can be infinitely traversed generating an infinite stack. \square

$$S \xrightarrow{\{i\}^*} S ++ [\{1, \dots, \{n\}, \{1, \dots, \{n\}, \dots]$$

Lemma 3. Given a CE-PDG increasing loop L , a stack S that allows Algorithm 1 to iterate at least once into it, and a loop L' which corresponds to L but replacing any access constraint by an asterisk constraint, then L' is not an increasing loop and it is not possible to infinitely traverse L' with S .

Proof. According to Definition 8, a cyclic flow dependency path is an increasing loop if the sequence of constraints generated by traversing it belongs to the language induced by the PDA in Fig. 10. The PDA cannot reach the final state if an asterisk constraint exists in L , thus, it cannot be an increasing loop. Moreover, the traversal of an asterisk

constraint described in case 6 of Table 1 always results in an empty stack (\perp). Therefore, if an asterisk constraint is included in an edge of L , then the second time that this edge is traversed the same node will be reached again with the same stack (\perp). Therefore, a second traversal is never done as indicated by line 22 in Algorithm 1. \square

Finally, we can prove Theorem 1:

Proof. The traversal of any sequence of nodes that is not a cycle (i.e., that does not represent a loop in the program) trivially terminates. Only loops can produce non-termination in Algorithm 1. But all loops are detected by the algorithm in line 22. According to Lemma 1, the traversal of all non-decreasing loops always terminates. On the other hand, as shown in Lemma 2 increasing loops can produce non-termination. However, all of them are detected by the PDA in Fig. 10. When an increasing loop is detected by the algorithm the stack is emptied.

We know by case 6 in Table 1 that including an asterisk constraint in a path is equivalent to emptying the stack. Therefore, according to Lemma 3 it is not possible to infinitely traverse the increasing loops found in the traversal made by the algorithm. Hence, Algorithm 1 always terminates. \square

4.8. Properties of the CE-PDG

In this section, we prove two important properties of the CE-PDG. We prove that it represents an improvement: slices are equal or smaller in size when compared to the PDG. We also prove that the slices produced contain all relevant elements.

Theorem 2 (Precision of the CE-PDG). Let P be a program, let $G = (N, E)$ and $G' = (N', E')$ be their corresponding PDG and CE-PDG. Let each edge in E be labeled with an empty constraint (\emptyset) and each edge in E' be labeled according to Section 4.1. Let sc represent each possible slicing criterion in the PDG, and SLICE be a function that applies Algorithm 1 for the given graph and slicing criterion.

$$\forall sc \in N \mid \text{SLICE}(G', sc) \subseteq \text{SLICE}(G, sc)$$

According to Theorem 2, the slices produced by the CE-PDG are at least as precise as the ones produced by the PDG. To prove Theorem 2, we first prove two lemmas:

Lemma 4 (CE-PDG Node Equivalence). Let $G = (N, E)$ be a PDG, and $G' = (N', E')$ be a CE-PDG. Let $E'_s \subset E'$ be the set of structural edges in G' , let E'^*_s be the reflexive and transitive closure of E'_s and let $\text{SUBTREE}(n) = \{n\} \cup \{m \mid (n, m) \in E'^*_s\}$.

$$\forall n' \in N' \mid \exists! n \in N \mid n' = n \vee n' \in \text{SUBTREE}(n)$$

Proof. Trivial, by construction of the CE-PDG. The CE-PDG has the following two properties:

1. $\forall n \in N . n \in N'$
2. $\forall n' \in N', n \notin N . n' \in \text{SUBTREE}(n''), n'' \in N$

Additionally, each node in the PDG is unique because each represents a different statement (even if two statements have the same code). Therefore, a node can either be common to both graphs (given 1), or represents a member of a data structure in a given node (given 2), fulfilling either the first or the second part of the lemma. \square

We use Lemma 4 to create a function $\mathcal{N} : N' \rightarrow N$, defined as $\mathcal{N}(x) = y \in N \mid x = y \vee x \in \text{SUBTREE}(y)$, given a node from a CE-PDG, returns its matching node in the corresponding PDG.

Lemma 5 (CE-PDG Edge Equivalence). Let $G = (N, E)$ be a PDG, $G' = (N', E')$ be a CE-PDG, and $n'_a, n'_b \in N'$ be two nodes in G' .

$$\forall (n'_a, n'_b) \in E' \mid (\mathcal{N}(n'_a), \mathcal{N}(n'_b)) \in E \vee \mathcal{N}(n'_a) = \mathcal{N}(n'_b)$$

Proof. The condition will be true if any of the \vee operands is true, so we can analyze them separately.

There are two different possibilities for the condition $(\mathcal{N}(n'_a), \mathcal{N}(n'_b)) \in E$ to hold:

1. The edge exists between the same nodes in the CE-PDG and in the PDG: $n'_a, n'_b \in N$.
2. The edge exists in the PDG and its source and/or target in the CE-PDG points to a data structure node n' , with $n' \notin N$.

On the other hand, there is one possibility for the condition $\mathcal{N}(n'_a) = \mathcal{N}(n'_b)$ to hold:

3. The edge does not exist in the PDG, but connects a data structure element to its parent or vice-versa, both nodes represents the same statement $n \in N$, with $\mathcal{N}(n'_a) = n$ and $\mathcal{N}(n'_b) = n$.

The CE-PDG features three kinds of edges, all of them backed by one of the previous three cases:

- Control edges (case 1): these edges remain unchanged in the CE-PDG.
- Flow edges: these can be classified according to their label (which can be an access or empty constraint).
 - Empty constraints (cases 1 and 2): these edges connect a variable definition to a usage. If the source and target are not members of a data structure, they remain unchanged (case 1). Otherwise, the source and/or target points to a data structure element (case 2).
 - Access constraints (case 3): these edges connect a member of a data structure to its parent or vice-versa.
- Structural edges (case 3): these edges connect a node that contains a data structure to a member of that data structure.

Thus, all edges present in the CE-PDG either correspond to a edge in the PDG (cases 1 and 2) or represent a reflexive edge (case 3), which has no effect on program slicing. \square

We can now prove [Theorem 2](#):

Proof. The PDG and CE-PDG are equivalent node by node ([Lemma 4](#)) and edge by edge ([Lemma 5](#)), so the slices produced by them would be equivalent. However, the PDG is labeled only with empty constraints and does not contain structural edges. For this reason, in some cases, the constraints present in the CE-PDG (see case 5 in [Table 1](#) or line 45 in [Algorithm 1](#)) and the combination of structural and flow edges (lines 10–12 in [Algorithm 1](#)) do not allow the traversal to continue, making CE-PDG slices smaller. Thus, every slice produced by the CE-PDG will be equal to or smaller than the equivalent slice produced by the PDG. \square

For the next theorem we need to provide a formal definition of complete slice:

Definition 11 (Complete Slice). Let P be a program, $G = (N, E)$ its PDG, and $n \in N$ a slicing criterion for P . A set of nodes $S \subseteq N$ is a *complete slice* if and only if the sequence of values produced in the variables of n when P is executed is a prefix of the sequence of values produced in n when S is executed.

Lemma 6. Let P be a program. Let $G = (N, E)$ be the PDG of P , $G' = (N', E')$ be its corresponding CE-PDG, and $sc \in N'$ be a slicing criterion. Not traversing flow edges after structural edges during the computation of $\text{SLICE}(G', sc)$ does not affect the completeness of said slice.

Proof. In the CE-PDG, a structural edge is only traversed by the slicing algorithm ([Algorithm 1](#)) after reaching a node inside a data structure. This node can either be a *definition* or a *usage*. We prove that the traversal restriction applied to the CE-PDG ([Algorithm 1](#), lines 10–12) generates complete slices in both situations:

- **The reached node ($n' \in N'$) is a definition.** Let $n = \mathcal{N}(n')$ be the CE-PDG node with the whole data structure that contains n' . According to the CE-PDG construction algorithm ([Section 4.1](#)), there is a path p_s from n to n' formed by structural edges (labeled with asterisk constraints), but also another parallel path p_f from n to n' formed by flow edges (labeled with their associated opening constraints). Both paths are traversed by [Algorithm 1](#) and both of them reach the same nodes in the data structure. If we reach n' from sc with a stack S and we traverse p_s from n' , then the stack will be emptied (see [Table 1](#), case (6)) with the first structural edge; thus traversing flow edges after the structural edges would be done with an empty stack (this would cause the algorithm to lose the context accumulated in the stack, which, in turn, would probably lead it to collect nodes that do not affect the slicing criterion). However, if we traverse p_f we will reach the same nodes in the data structure but with a properly updated stack S' (see [Table 1](#), case (2)). Moreover, [Algorithm 1](#) traverses all flow edges from the nodes in p_f with the proper stacks. Therefore, the traversal of flow edges (with an empty stack) after traversing structural edges is unnecessary and can be ignored, preserving completeness.
- **The reached node ($n' \in N'$) is a use.** Let $n = \mathcal{N}(n')$ be the CE-PDG node with the whole data structure that contains n' . According to the CE-PDG, when n' is a use, there is only one path p_s from n to n' , formed by structural edges. Traversing p_s is only necessary to include in the slice the hierarchical structure of n' inside the data structure. Therefore, there is no need to traverse flow edges after reaching a node with structural edges because the value of n' cannot depend on the nodes reachable through structural edges. In fact, if the value of n' depends on a node n_0 , then there must exist a flow edge (n_0, n') in the graph, and this edge is traversed by [Algorithm 1](#). Thus, the traversal of flow edges after traversing structural edges is unnecessary and can be ignored, preserving completeness. \square

Theorem 3 (Completeness of the CE-PDG). Let P be a program. Let G be the PDG of P , G' be its corresponding CE-PDG, and $sc \in N'$ be a slicing criterion. If $\text{SLICE}(G, \mathcal{N}(sc))$ is a complete slice, then $\text{SLICE}(G', sc)$ is also a complete slice.

Proof. The CE-PDG slicing algorithm contains two traversal limitations compared to the standard PDG slicing algorithm : (i) flow edges are not traversed after a structural edge is traversed and (ii) the traversal limitation imposed by access constraints. As proved in [Lemma 6](#), limitation (i) does not affect completeness. On the other hand, to prove limitation (ii) we consider all possible situations that concern the management of access constraints. We divide them into two different cases:

1. $\text{SLICE}(G, \mathcal{N}(sc))$ **does not include any node that contains a data structure.** This case is trivially proved using [Lemmas 4](#) and [5](#). The nodes without data structures and the edges connecting them in the PDG remain unchanged in the CE-PDG. Then, during the algorithm, the same edges are traversed and the same nodes are included in the slice in both graphs ($\text{SLICE}(G, \mathcal{N}(sc)) = \text{SLICE}(G', sc)$). Thus, since the PDG slice is complete, the CE-PDG slice is also complete.
2. $\text{SLICE}(G, \mathcal{N}(sc))$ **includes at least one node with a data structure.** Inside this case, we consider two different scenarios:

(I). No element included in the slice has been defined inside a data structure (e.g. $\{A, B\} = \underline{X}$ or $\underline{X} = \{A, B\}$). According to [Table 1](#), only opening constraints (associated to definitions inside data structures) are pushed into the stack (case (2)). Opening constraints at the top of the stack are the ones that can limit the traversal during the CE-PDG slicing algorithm (see case (5) in [Table 1](#)). However, since no opening constraint is collected during the traversal (because all elements included in the slice have been defined outside a data structure), we

reach a scenario similar to the previous case. The elements outside data structures are connected by flow and control dependencies in the same way in both the PDG and CE-PDG (as stated in Lemma 5 (case 1)). Additionally, since closing constraints can be traversed while the stack is empty, the elements inside data structures which are used are also included in the slice, ensuring the CE-PDG to include the same code as the PDG. Thus, the PDG and the CE-PDG include the same code in the slice, resulting in the completeness of the CE-PDG slice.

(li). At least one element in the slice has been defined inside a data structure (e.g. $\{A, B\} = X$). In this case, an opening constraint is pushed into the traversal stack, enabling the possibility to limit the traversal. When slicing the CE-PDG after pushing one (or more) opening constraints, there are two possible situations:

- (a) No closing constraints are reached during the same CE-PDG traversal inside the procedure. In this case the traversal is not limited by the CE-PDG, and the slice includes only those elements in the data structures that contain the definitions included, which are reached by flow dependences. In this case, the CE-PDG makes possible to exclude the elements defined in the same data structure that are not reached by flow edges, which do not influence the slicing criterion. Thus, $\text{SLICE}(G', sc) \subseteq \text{SLICE}(G, \mathcal{N}(sc))$, being both complete.
- (b) Other closing constraints are reached during the same CE-PDG traversal inside the procedure. We prove this case by induction on the depth d of the definition inside its data structure. In the proof, S represents the value of the stack before reaching this node, O_i represents the opening constraint with index i , and C_i represents the closing constraint with the same index i . \square

(Base case: $d = 1$). First, in the case where the definition is at depth level 1, the traversal collects the constraint at the flow edge (O_i) and pushes it into the stack (S, O_i). When reaching the data structure with closing constraints, only the edge with the complementary constraint (C_i) is traversed (case (4) in Table 1). As a result, the element required at the data structure with uses is the one at the same position that the one at the data structure with definitions. Since no more elements are required to compute the value of the slicing criterion, and according to the flow dependences defined at the CE-PDG (which are extracted from the flow dependences of the PDG), the CE-PDG slice is complete.

(Induction hypothesis: $d = n$). We assume as the induction that completeness holds for definitions inside data structures located at depth level n , which generates a sequence of opening constraints S_O .

(Inductive case: $d = n + 1$). We prove that the theorem holds for any expression with depth $n + 1$. When a definition node inside a data structure at depth $n + 1$ is reached with an initial stack S , the traversal of the flow edge with constraint O_j pushes it to the slice resulting in the stack S, O_j . Then, the traversal reaches depth n and, according to the induction hypothesis it includes to the stack a sequence of n opening constraints (S_O). As a result the stack results in the new stack S, O_j, S_O . Considering the induction hypothesis, when the algorithm reaches a data structure that contains a use, it pops this sequence of constraints reaching depth level n and preserving completeness (then, the stack becomes S, O_j again). Hence, we are in the same situation as the base case, where only the edge with the complementary constraint (C_j) is traversed according to case (4) in Table 1, generating a complete slice.

5. Implementation and empirical evaluation

Comparing our implementation against other slicers is not the best way to assess the proposed stack extension to the PDG because we would find big differences in the PDG construction time, slicing time, and slicing precision due to differences in the libraries used, different treatments for syntax constructs such as list comprehensions, guards, etc. Therefore, we would not be able to assess the specific impact of

the stack on the slicer's precision and performance. The only way to do a fair comparison is to implement a single slicer that is able to build and slice the PDG with and without constraints.

All the algorithms and ideas described in this paper have been implemented in a slicer for Erlang called e-Knife. e-Knife can produce slices based on either the PDG or the CE-PDG. Thus, it allows us to know exactly the additional cost required to build and traverse the constraints, and the extra precision obtained by doing so. e-Knife is a Java program with 13355 LOC (excluding comments and empty lines). It is an open-source project and is publicly available.⁴ The exact version used to produce the data in this work and all benchmarks has been archived on Zenodo (Galindo et al., 2023a).

Additionally, anyone can slice a program via a web interface,⁵ without the need to build the project locally. Large or very complex programs may run into the memory and time limitations that are in place to avoid abuse.

To evaluate e-Knife, we used Bencher,⁶ a program slicing benchmark suite for Erlang. All the benchmarks were interprocedural programs, so we have created a new intraprocedural version of them (by inlining functions). This intraprocedural version has been made publicly available (every benchmark in Bencher has a link to its intraprocedural version). To evaluate the techniques proposed throughout this work, we have built both graphs (PDG and CE-PDG) for each of the intraprocedural benchmarks. Then, we sliced both graphs with respect to all possible slicing criteria,⁷ which guarantees that there is no bias in the selection of slicing criteria.

The benchmarks were run on a Intel Xeon E-2136 CPU running Debian Linux 11 with 32 GB of DDR4 RAM available. All processes were stopped, except for *init* and *sshd*. We strictly followed the methodology proposed by Georges et al. (2007). Each program's graph was built 1001 times, and the graphs were sliced 1001 times per criterion. To ensure real independence, the first iteration was always discarded (to avoid the influence of dynamically loading libraries to physical memory, data persisting in the disk cache, etc.). From the 1000 remaining iterations, we retained a window of 10 measurements when steady-state performance was reached, i.e., once the coefficient of variation (CoV, the standard deviation divided by the mean) of the 10 iterations falls below a preset threshold of 0.01 or the lowest CoV if no window reached it. It is with these 10 iterations that we computed the average time taken by each operation (building each graph or slicing each graph w.r.t. each criterion).

The results of the experiments performed are summarized in Table 2. The two columns (PDG, CE-PDG) display the average time required to build each graph. Building the CE-PDG, as in the PDG, is a quadratic operation; and the inclusion of labels in the edges is a linear operation w.r.t. the amount of nodes in the graph. Thus, building the CE-PDG is only slightly slower than its counterpart. The other columns are as follows (average values are w.r.t. all slicing criteria):

DSS: the number of data structure access constraints in the CE-PDG. It is a metric to measure the amount and size of the composite data structures.

Function: the name of the function where the slicing criterion is located.

#SCs: the number of slicing criteria in that function.

PDG, CE-PDG: the average time required to slice the corresponding graph.

⁴ <https://mist.dsic.upv.es/git/program-slicing/e-knife-erlang>

⁵ <https://mist.dsic.upv.es/e-knife-constrained/>

⁶ <https://mist.dsic.upv.es/bencher/>

⁷ Each variable use or definition in all functions that contain complex data structures.

Table 2Summary of experimental results for *Bench*, comparing the PDG (without constraints) to the CE-PDG (with constraints).

Program	Graph generation			Slice					
	PDG	CE-PDG	DSs	Function	#SCs	PDG	CE-PDG	Slowdown	Red. size
bench1A.erl	3230.08 ms	3233.42 ms	515	getLast/2	26	59.12 μ s	269.26 μ s	5.68 \pm 1.38	14.88 \pm 3.16%
				getNext/3	174	228.10 μ s	1014.41 μ s	4.62 \pm 0.29	13.09 \pm 1.57%
				getStringDate/1	11	22.56 μ s	164.10 μ s	8.62 \pm 2.05	8.67 \pm 3.88%
				main/1	57	800.96 μ s	1644.59 μ s	2.21 \pm 0.20	37.14 \pm 6.97%
bench3A.erl	37.53 ms	37.54 ms	4	tuples/2	22	30.36 μ s	100.91 μ s	3.36 \pm 0.20	5.46 \pm 2.03%
bench4A.erl	53.22 ms	53.26 ms	20	main/2	31	64.19 μ s	163.78 μ s	2.65 \pm 0.19	20.79 \pm 5.38%
bench5A.erl	32.87 ms	32.90 ms	8	lists/2	18	44.67 μ s	136.81 μ s	3.12 \pm 0.17	6.51 \pm 2.02%
bench6A.erl	236.33 ms	236.42 ms	37	ft/2	34	65.87 μ s	249.55 μ s	4.11 \pm 0.28	8.71 \pm 2.45%
				ht/2	16	20.37 μ s	77.71 μ s	4.07 \pm 0.25	10.79 \pm 3.69%
bench9A.erl	117.30 ms	117.45 ms	16	main/2	18	136.84 μ s	210.77 μ s	1.47 \pm 0.07	1.38 \pm 1.04%
bench11A.erl	15.10 ms	15.13 ms	6	lists/2	16	35.92 μ s	108.67 μ s	2.92 \pm 0.33	6.47 \pm 2.15%
bench12A.erl	1042.13 ms	1042.94 ms	103	add/4	26	78.11 μ s	273.66 μ s	5.02 \pm 1.13	15.38 \pm 4.13%
				from_ternary/2	9	17.28 μ s	86.14 μ s	6.18 \pm 1.84	3.56 \pm 2.61%
				main/3	39	75.28 μ s	196.66 μ s	4.02 \pm 0.79	8.43 \pm 6.19%
				mul/3	21	41.43 μ s	143.00 μ s	4.95 \pm 1.17	2.74 \pm 1.27%
				to_ternary/2	13	57.37 μ s	126.01 μ s	4.13 \pm 1.91	1.02 \pm 1.32%
bench14A.erl	2300.65 ms	2301.05 ms	75	main/2	81	73.50 μ s	284.47 μ s	3.89 \pm 0.34	8.62 \pm 2.46%
bench15A.erl	1182.25 ms	1182.66 ms	44	main/4	71	165.03 μ s	268.53 μ s	3.84 \pm 0.75	1.72 \pm 1.53%
bench16A.erl	159.12 ms	159.26 ms	16	word_count/5	36	59.85 μ s	130.93 μ s	2.67 \pm 0.23	7.72 \pm 2.68%
bench17A.erl	44.94 ms	45.05 ms	8	mug/3	19	41.23 μ s	92.11 μ s	2.43 \pm 0.22	5.59 \pm 3.02%
bench18A.erl	49.41 ms	49.52 ms	8	mbe/2	19	58.36 μ s	117.12 μ s	2.28 \pm 0.23	7.38 \pm 4.59%
Totals and averages for set A					757	164.67 μ s	491.93 μ s	3.88 \pm 0.50	11.27 \pm 2.91%
bench1B.erl	2640.65 ms	2643.99 ms	493	main/1	273	1538.26 μ s	6648.83 μ s	4.24 \pm 0.23	24.02 \pm 1.93%
bench2B.erl	73.61 ms	73.62 ms	2	main/2	17	79.86 μ s	168.92 μ s	3.07 \pm 0.70	0.43 \pm 0.58%
bench3B.erl	35.16 ms	35.17 ms	4	tuples/2	18	51.44 μ s	133.12 μ s	2.57 \pm 0.07	4.33 \pm 1.21%
bench4B.erl	26.49 ms	26.52 ms	20	main/2	39	89.57 μ s	240.93 μ s	2.66 \pm 0.17	13.11 \pm 3.79%
bench5B.erl	18.31 ms	18.34 ms	8	lists/2	11	54.50 μ s	136.21 μ s	2.50 \pm 0.10	6.88 \pm 0.85%
bench6B.erl	54.35 ms	54.42 ms	25	tuples/2	42	46.94 μ s	129.73 μ s	3.01 \pm 0.22	8.25 \pm 1.62%
bench8B.erl	87.29 ms	87.39 ms	16	main/2	42	204.09 μ s	607.50 μ s	3.06 \pm 0.18	0.73 \pm 0.68%
bench9B.erl	34.89 ms	34.95 ms	10	main/2	17	187.93 μ s	294.62 μ s	1.51 \pm 0.06	1.16 \pm 0.85%
bench10B.erl	97.22 ms	97.39 ms	18	main/1	35	263.35 μ s	656.71 μ s	2.85 \pm 0.26	2.23 \pm 1.15%
bench11B.erl	12.95 ms	12.98 ms	8	lists/2	13	45.68 μ s	118.99 μ s	2.56 \pm 0.13	8.02 \pm 2.08%
bench12B.erl	294.70 ms	295.25 ms	79	main/3	88	917.43 μ s	2692.51 μ s	3.23 \pm 0.40	2.61 \pm 2.67%
bench13B.erl	27.08 ms	27.11 ms	4	main/0	22	135.10 μ s	253.26 μ s	1.92 \pm 0.11	0.48 \pm 0.39%
bench14B.erl	147.19 ms	147.48 ms	49	main/2	51	112.06 μ s	371.92 μ s	3.21 \pm 0.44	13.34 \pm 4.49%
bench15B.erl	217.26 ms	217.49 ms	38	main/4	65	286.95 μ s	495.93 μ s	2.70 \pm 0.41	8.78 \pm 2.84%
bench16B.erl	102.53 ms	102.62 ms	16	word_count/5	40	131.26 μ s	349.56 μ s	3.11 \pm 0.26	4.14 \pm 1.53%
bench17B.erl	57.63 ms	57.64 ms	8	mug/3	19	166.01 μ s	352.78 μ s	2.11 \pm 0.09	4.96 \pm 2.38%
bench18B.erl	62.95 ms	63.03 ms	8	mbe/2	19	273.12 μ s	629.87 μ s	2.34 \pm 0.11	0.05 \pm 0.10%
Totals and averages for set B					811	704.73 μ s	2735.10 μ s	3.30 \pm 0.26	11.79 \pm 2.05%
Totals and averages					1568	444.00 μ s	1652.14 μ s	3.54 \pm 0.38	11.54 \pm 2.47%

Table 3Summary of experimental results for the *erlsom* library, comparing PDG to the CE-PDG.

Program	LOC	Graph generation			Slice				
		PDG	CE-PDG	DSs	#SCs	PDG	CE-PDG	Slowdown	Red. size
erlsom.erl	182	0.279 s	0.280 s	100	172	29.48 ms	49.69 ms	1.68 \pm 0.55	5.85 \pm 8.17%
erlsom_add.erl	28	0.224 s	0.224 s	34	60	46.87 ms	81.74 ms	1.69 \pm 0.84	4.23 \pm 12.19%
erlsom_compile.erl	693	25.662 s	25.665 s	1002	1335	57.09 ms	102.70 ms	1.71 \pm 0.58	7.00 \pm 10.72%
erlsom_example_value.erl	270	1.515 s	1.518 s	299	349	35.59 ms	49.64 ms	1.51 \pm 0.78	11.31 \pm 12.18%
erlsom_lib.erl	939	6.729 s	6.736 s	754	881	87.35 ms	159.15 ms	1.67 \pm 0.66	11.41 \pm 8.81%
erlsom_parse.erl	893	84.991 s	84.995 s	891	1224	134.94 ms	232.84 ms	1.67 \pm 0.48	14.21 \pm 7.62%
erlsom_parseXsd.erl	504	52.615 s	52.629 s	1608	50	20.00 ms	34.38 ms	1.65 \pm 0.39	0.24 \pm 3.67%
erlsom_pass2.erl	605	15.463 s	15.467 s	695	1166	52.07 ms	81.51 ms	1.60 \pm 0.59	10.47 \pm 9.81%
erlsom_sax.erl	74	0.818 s	0.818 s	60	111	44.70 ms	71.36 ms	1.59 \pm 0.63	6.96 \pm 9.04%
erlsom_sax_latin1.erl	1140	18.396 s	18.401 s	733	1564	54.80 ms	96.75 ms	1.69 \pm 0.68	9.69 \pm 9.35%
erlsom_sax_latin9.erl	1140	18.527 s	18.532 s	734	1559	55.02 ms	96.84 ms	1.69 \pm 0.68	9.84 \pm 9.46%
erlsom_sax_lib.erl	160	0.591 s	0.593 s	272	261	43.76 ms	92.83 ms	1.97 \pm 0.61	9.00 \pm 7.38%
erlsom_sax_list.erl	1140	18.703 s	18.711 s	1249	1806	60.69 ms	135.42 ms	1.91 \pm 1.64	10.77 \pm 8.09%
erlsom_sax_utf8.erl	1140	21.556 s	21.561 s	736	1600	56.13 ms	98.26 ms	1.69 \pm 0.68	9.51 \pm 9.42%
erlsom_sax_utf16be.erl	1140	21.521 s	21.526 s	736	1600	55.80 ms	97.75 ms	1.68 \pm 0.68	9.51 \pm 9.42%
erlsom_sax_utf16le.erl	1140	20.414 s	20.420 s	748	1664	66.19 ms	109.91 ms	1.66 \pm 0.68	9.21 \pm 9.42%
erlsom_simple_form.erl	126	2.589 s	2.590 s	201	201	48.19 ms	71.45 ms	1.39 \pm 0.83	15.36 \pm 9.50%
erlsom_type2xsd.erl	190	1.831 s	1.832 s	289	324	84.39 ms	300.61 ms	2.17 \pm 1.66	9.29 \pm 8.00%
erlsom_ucs.erl	158	3.356 s	3.358 s	222	309	94.11 ms	130.81 ms	1.43 \pm 0.58	12.02 \pm 9.20%
erlsom_write.erl	606	16.591 s	16.598 s	719	1010	118.65 ms	207.84 ms	1.66 \pm 0.64	11.51 \pm 9.14%
erlsom_writeHrl.erl	240	1.446 s	1.447 s	201	294	46.68 ms	73.63 ms	1.51 \pm 0.87	14.16 \pm 12.53%
ucs.erl	351	4.075 s	4.078 s	211	351	53.06 ms	65.13 ms	1.30 \pm 0.61	13.90 \pm 12.52%
Totals and averages	12859	15.359 s	15.363 s	12494	17891	67.37 ms	121.40 ms	1.69 \pm 0.88	10.24 \pm 9.45%

Slowdown: the average additional time required (with 95% error margins), when comparing the CE-PDG with the PDG. For example, on average (last row), the computation of each slice is 3.54 times slower in the CE-PDG.

Red. Size: the average reduction in the sizes of slices (with 95% error margins).⁸ It is computed as $(A - B)/A$ where A is the size (number of AST nodes) of the slice computed with the standard (field-insensitive) algorithm and B is the size (number of AST nodes) of the slice computed with the field-sensitive algorithm (Algorithm 1). This way of measuring the size of the slices is much more precise and fair. LOC is not proper because it can ignore the removal of subexpressions. PDG/CE-PDG nodes is also not a good solution because the CE-PDG includes nodes and edges not present in their PDG counterparts, therefore they are incomparable.

The averages shown at the bottom of the table are the averages of all slicing criteria and not the averages of each function's average.

The first 13 benchmarks (set A) are benchmarks with complex data structures but without cycles, while the rest of the benchmarks (set B) do contain cycles. In set A, each slice produced by the CE-PDG is around four times slower. However, this has little impact, as each slice consumes just hundreds of milliseconds. As can be seen in each row, generating the graph is at least 3 orders of magnitude slower than slicing it. This increase in time is offset by the average reduction of the slices, which is 11.27%. This reduction goes up to 37.14% in function `main/1` from `bencher1A`, as it contains complex data structures that can be efficiently sliced with the CE-PDG. Therefore, our technique reduces the size of the slices by $11.27 \pm 2.91\%$ at almost no cost (only a few μ s).

If we consider programs with cycles (set B), due to the analysis of loops, the slowdown is around three to four times slower (the slowdown is 3.30) and the reduction in the size of the slices is 11.79 ± 2.05 . If we consider all benchmarks, our technique reduces the size of the slices by 11.54% with a slowdown of 3.54. This is a good result: for many applications (e.g., debugging) reducing the suspicious code over 11.54% with the cost of increasing the slicing time by only a few milliseconds is a good trade-off to make.

Regarding the effect of data structures on the results, we can see that benchmarks with more data structures have the potential to produce more precise slices with the CE-PDG.

5.1. Slicing larger programs

Although `Bencher` is useful because it contains challenging slicing programs specifically designed to test program slicers, and because it helps determine whether a slicer is complete by comparing the slices produced by a given implementation against a known minimal slice. However, the size of its benchmarks is rather small, with files ranging from 12 to 114 LOC. To test the efficiency of the CE-PDG on larger programs with more complex functions, we have performed an experimental evaluation of the `erlsom` Erlang library, which implements an XML parser.⁹ This library contains 13k LOC across 22 source files. The evaluation followed the same methodology and process as `Bencher`'s evaluation, generating over 4k graphs and analysing over 35 million slices. The results can be seen on Table 3, which is structured like Table 2.

⁸ The minimum value for Red. Size is 0, even if some error margins hint at otherwise. The slices produced by the PDG can never be smaller than those produced by the CE-PDG.

⁹ The exact version used corresponds to commit hash `e00d2c7426fe7a8c49f8959da38a11d21b8b2168`, and is available for download at <https://github.com/willemdj/erlsom/archive/e00d2c7.zip>.

From the graph generation columns, we can see that the cost of building the CE-PDG is almost equal to the cost of building the PDG. As already known, the cost scales with program size, with the average growing around three orders of magnitude (from milliseconds to seconds) with respect to `Bencher`, reaching a maximum of 85 s. Any improvement on the algorithms used to compute the PDG would benefit the CE-PDG.

Regarding the time needed to slice the CE-PDG, it also grows three orders of magnitude with respect to `Bencher`, due to size. On the bright side, the increase in the size of the programs lowers the relative speed of the graphs. In `Bencher`, the CE-PDG was 3 times slower than the PDG and, in this case, the slowdown is much lower (1.69), which may indicate that this technique scales well with size.

Finally, the reduction in slice size is much more stable in `erlsom` than in `Bencher`, achieving an average of 10.24%, with a maximum of 15% in `erlsom_simple_form.erl`. However, these results have a much higher variance, probably due to a higher variance in the size of slices. In conclusion, in larger programs we achieve a reduction in the size of the slices of 10.24% at the cost of having a slowdown of 1.79 to produce slices. This slowdown is 50 ms on average.

6. Related work

Transitive data dependence analysis has been extensively studied (Reps et al., 1995; Sridharan et al., 2007). Less attention has received, however, the problem of field-sensitive data dependence analysis (Litvak et al., 2010; Ramalingam et al., 1999; Kumar K. et al., 2019; Späth et al., 2019). The existing approaches can be classified into two groups: those that treat composite structures as a whole (Lyle, 1984; Ottenstein and Ottenstein, 1984; Muchnick, 1997a; Litvak et al., 2010), and those that decompose them into small atomic data types (Korel and Laski, 1990; Agrawal et al., 1991; Binkley and Gallagher, 1996; Muchnick, 1997b; Ramalingam et al., 1999; Krinke, 2003a; Anderson et al., 2003; Graf, 2010). The latter approach is often called *atomization* or *scalar replacement*, and it basically consists of a program transformation that recursively disassembles composite structures to their primitive components. However, slicing over the decomposed structures usually uses traditional dependence graph based traversal (Krinke, 2003a; Anderson et al., 2003; Graf, 2010) which limits the accuracy. Other important approaches for field-sensitive data dependence analysis of this kind are Kumar K. et al. (2019), Späth et al. (2019), Litvak et al. (2010). Litvak et al. (2010) proposed a field-sensitive program dependence analysis that identifies dependencies by computing the memory ranges written/read by definitions/uses. Späth et al. (2019) proposed the use of pushdown systems to encode and solve field accesses and uses. Snelting et al. (2006) present an approach to identify constraints over paths in dependence graphs. Our approach combines atomization with the addition of constraints checked by pushdown systems to improve the accuracy of slicing composite data structures.

There exist approaches for the field-sensitive slicing of some specific data structures. If we refer to arrays, some static proposals consider the whole array as a variable, and each access as a definition or use of that variable (Lyle, 1984). However, this technique produces complete, but unnecessarily large program slices (Binkley and Gallagher, 1996). The PDG variant of Ottenstein and Ottenstein (1984) represents composite data types providing a node for each one of its subexpressions, and provides special *select* and *update* operators to access the elements of an array. Other static approaches rely on determining whether two statically unknown vector accesses can refer to the same memory location during runtime (Landi and Ryder, 1992; Choi et al., 1993). Some papers (Banerjee, 1988; Maydan et al., 1991; Pugh and Wonnacott, 1992) propose algorithms that demonstrate the absence of a flow dependence between array accesses under certain conditions.

Some approaches (Chen and Xu, 2001b; Graf, 2010) have been also proposed to accurately represent the inner structure of objects and the dependencies between their data members. Most object-oriented

approaches (Liang and Harrold, 1998; Walkinshaw et al., 2003; Graf, 2010; Galindo et al., 2023c) are based on the same principle: object variables and their inner data members are unfolded in a tree-like representation when used at function calls. This allows for the generation of dependencies between data members of a particular object and to accurately slice off those data members of an object that are not affecting the slicing criterion. Our representation is inspired by this tree-like structure, but with some differences. In our representation, the tree structure is connected with a new kind of edges (structural edges) instead of control edges. This allows us to apply a different slicing behavior for structural edges without interfering in the traversal restrictions given to control edges in some slicing algorithms (Kumar and Horwitz, 2002). Additionally, our tree structure is connected not only with structural edges but also with flow edges; providing a more realistic representation of the dependencies between a composite structure and all its elements.

Several works have tried to adapt the PDG for functional languages dealing with tuple structures in the process (Cheda et al., 2007; Brown, 2008; Tóth et al., 2010; Kumar K. et al., 2019). Some of them with a high abstraction level (Rodrigues and Barbosa, 2005), and other ones with a low granularity level. Silva et al. (2012) propose a new graph representation for the sequential part of Erlang called the Erlang Dependence Graph. Their graph, despite being built with the minimum possible granularity (each node in the graph corresponds to an AST node) and being able to select subelements of a given composite data structure, does not have a mechanism to preserve the dependency of the tuple elements when a tuple is collapsed into a variable; i.e., they do not solve the *slicing pattern matching* problem (for instance, they cannot solve the program in Fig. 2). In contrast, although our graph is only fine-grained at composite data structures, we overcome their limitations by introducing an additional component to the graph, the constrained edges, which allow us to carry the dependence information between definition and use even if the composite structure is collapsed in the process.

7. Conclusion

Static analyses often use a representation of the program being analyzed, and this representation strongly influences their correctness, completeness, and performance. In the particular case of intraprocedural program slicing, the standard representation used is the PDG. Unfortunately, the PDG's data dependencies are imprecise when modeling composite data structures. In particular, the information stored in the nodes of the PDG (i.e., statements) is often inappropriate when representing composite data structures.

To solve this problem, we present a generalization of the PDG called CE-PDG where (i) the inner components of the composite data structures are unfolded into a tree-like representation, providing an independent representation for their subexpressions and allowing us to accurately define intra-statement data dependencies, and (ii) the edges are augmented with constraints (constrained edges), which allows the propagation of the component dependence information through the traversal of the graph during the slicing process. As a result, the CE-PDG allows the user to select any subexpression of a data structure as the slicing criterion and it computes accurate slices for (recursive) composite data structures. It ignores irrelevant elements inside the same statement and allows for the transference of data dependence information through the compression and expansion of composite structures. An evaluation of our approach shows a slowdown of 3.88/3.30 and a reduction of the slices of 11.27%/11.79% for programs without/with cycles.

CRedit authorship contribution statement

Carlos Galindo: Writing – review & editing, Software, Investigation, Formal analysis, Visualization. **Jens Krinke:** Conceptualization, Writing – review & editing, Supervision. **Sergio Pérez:** Writing – review & editing, Methodology, Software, Investigation. **Josep Silva:** Writing – original draft, Methodology, Supervision.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

All sources used to run the experiments and produce the results are publicly available at <http://dx.doi.org/10.5281/zenodo.10439063>. Additionally, Section 5 includes links to the software and benchmarks used.

References

- Agrawal, H., DeMillo, R.A., Spafford, E.H., 1991. Dynamic slicing in the presence of unconstrained pointers. In: *Proceedings of the Symposium on Testing, Analysis, and Verification*. pp. 60–73.
- Allen, M., Horwitz, S., 2003. Slicing Java programs that throw and catch exceptions. *SIGPLAN Not.* 38 (10), 44–54.
- Anderson, P., Reps, T., Teitelbaum, T., 2003. Design and implementation of a fine-grained software inspection tool. *IEEE Trans. Softw. Eng.* 29 (8), 721–733.
- Ball, T., Horwitz, S., 1993. Slicing programs with arbitrary control-flow. In: *Proceedings of the First International Workshop on Automated and Algorithmic Debugging. AADBUG '93*, Springer-Verlag, London, UK, UK, pp. 206–222.
- Banerjee, U.K., 1988. *Dependence Analysis for Supercomputing*, Vol. 60. Springer US, USA.
- Binkley, D., 1993. Precise executable interprocedural slices. *ACM Lett. Program. Lang. Syst.* 2 (1–4), 31–45.
- Binkley, D., Gallagher, K.B., 1996. Program slicing. *Adv. Comput.* 43 (2), 1–50. [http://dx.doi.org/10.1016/S0065-2458\(08\)60641-5](http://dx.doi.org/10.1016/S0065-2458(08)60641-5).
- Brown, C.M., 2008. *Tool Support for Refactoring Haskell Programs* (Ph.D. thesis). School of Computing, University of Kent, Canterbury, Kent, UK.
- Cheda, D., Silva, J., Vidal, G., 2007. Static slicing of rewrite systems. In: *Proceedings of the 15th International Workshop on Functional and (Constraint) Logic Programming. WFLP 2006*, Elsevier ENTCS 177, pp. 123–136.
- Chen, Z., Xu, B., 2001a. Slicing concurrent Java programs. *SIGPLAN Not.* 36 (4), 41–47. <http://dx.doi.org/10.1145/375431.375420>, URL <http://doi.acm.org/10.1145/375431.375420>.
- Chen, Z., Xu, B., 2001b. Slicing object-oriented Java programs. *SIGPLAN Not.* 36 (4), 33–40. <http://dx.doi.org/10.1145/375431.375418>.
- Choi, J.-D., Burke, M., Carini, P., 1993. Efficient flow-sensitive interprocedural computation of pointer-induced aliases and side effects. In: *Proceedings of the 20th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. pp. 232–245.
- Clarke, E.M., Fujita, M., Rajan, S.P., Reps, T., Shankar, S., Teitelbaum, T., 2002. Program slicing for VHDL. *Int. J. Softw. Tools Technol. Transf.* 4 (1), 125–137. <http://dx.doi.org/10.1007/s10090100069>.
- DeMillo, R.A., Pan, H., Spafford, E.H., 1996. Critical slicing for software fault localization. *SIGSOFT Softw. Eng. Notes* 21 (3), 121–134. <http://dx.doi.org/10.1145/226295.226310>, URL <http://doi.acm.org/10.1145/226295.226310>.
- Ferrante, J., Ottenstein, K.J., Warren, J.D., 1987. The program dependence graph and its use in optimization. *ACM Trans. Program. Lang. Syst.* 9 (3), 319–349.
- Galindo, C., Krinke, J., Pérez, S., Silva, J., 2022. Field-sensitive program slicing. In: *Software Engineering and Formal Methods: 20th International Conference, SEFM 2022*, Berlin, Germany, September 26–30, 2022, *Proceedings*, Vol. 13550. Springer Nature, pp. 74–90.
- Galindo, C., Krinke, J., Pérez, S., Silva, J., 2023a. Replication Package for the Paper “Field-Sensitive Program Slicing”. Zenodo, <http://dx.doi.org/10.5281/zenodo.10439063>.
- Galindo, C., Pérez, S., Silva, J., 2021. Slicing unconditional jumps with unnecessary control dependencies. In: *Fernández, M. (Ed.), Logic-Based Program Synthesis and Transformation*. In: *Lecture Notes in Computer Science (LNCS)*, vol. 12561, Springer International Publishing, Cham, pp. 293–308.
- Galindo, C., Pérez, S., Silva, J., 2023b. Exception-sensitive program slicing. *J. Log. Algebraic Methods Program.* 130, 100832. <http://dx.doi.org/10.1016/j.jlamp.2022.100832>, URL <https://www.sciencedirect.com/science/article/pii/S2352220822000852>.
- Galindo, C., Pérez, S., Silva, J., 2023c. Program slicing of Java programs. *J. Log. Algebraic Methods Program.* 130, 100826. <http://dx.doi.org/10.1016/j.jlamp.2022.100826>, URL <https://www.sciencedirect.com/science/article/pii/S2352220822000797>.
- Georges, A., Buytaert, D., Eeckhout, L., 2007. Statistically rigorous java performance evaluation. *SIGPLAN Not.* 42 (10), 57–76. <http://dx.doi.org/10.1145/1297105.1297033>, URL <http://doi.acm.org/10.1145/1297105.1297033>.
- Gongye, X., Wang, Y., Wen, Y., Nie, P., Lin, P., 2020. A simple detection and generation algorithm for simple circuits in directed graph based on depth-first traversal. *Evol. Intell.* <http://dx.doi.org/10.1007/s12065-020-00416-6>.

- Graf, J., 2010. Speeding up context-, object- and field-sensitive SDG generation. In: 2010 10th IEEE Working Conference on Source Code Analysis and Manipulation. pp. 105–114. <http://dx.doi.org/10.1109/SCAM.2010.9>.
- Hajnal, Á., Forgács, I., 2012. A demand-driven approach to slicing legacy COBOL systems. *J. Softw. Maintenance* 24 (1), 67–82. URL <http://dblp.uni-trier.de/db/journals/smr/smr24.html#HajnalF12>.
- Horwitz, S., Reps, T., Binkley, D., 1990. Interprocedural slicing using dependence graphs. *ACM Trans. Program. Lang. Syst.* 12 (1), 26–60.
- Korel, B., Laski, J., 1988. Dynamic program slicing. *Inform. Process. Lett.* 29 (3), 155–163. [http://dx.doi.org/10.1016/0020-0190\(88\)90054-3](http://dx.doi.org/10.1016/0020-0190(88)90054-3), URL <http://www.sciencedirect.com/science/article/pii/0020019088900543>.
- Korel, B., Laski, J., 1990. Dynamic slicing of computer programs. *J. Syst. Softw.* 13 (3), 187–195. [http://dx.doi.org/10.1016/0164-1212\(90\)90094-3](http://dx.doi.org/10.1016/0164-1212(90)90094-3).
- Krinke, J., 2003a. Advanced Slicing of Sequential and Concurrent Programs (Ph.D. thesis). Universität Passau.
- Krinke, J., 2003b. Context-sensitive slicing of concurrent programs. In: Proceedings of the 9th European Software Engineering Conference Held Jointly with 11th ACM SIGSOFT International Symposium on Foundations of Software Engineering. pp. 178–187. <http://dx.doi.org/10.1145/940071.940096>.
- Kumar, S., Horwitz, S., 2002. Better slicing of programs with jumps and switches. In: Proceedings of the 5th International Conference on Fundamental Approaches to Software Engineering. FASE 2002, In: Lecture Notes in Computer Science (LNCS), vol. 2306, Springer, pp. 96–112.
- Kumar K., P., Sanyal, A., Karkare, A., Padhi, S., 2019. A static slicing method for functional programs and its incremental version. In: Proceedings of the 28th International Conference on Compiler Construction. In: CC 2019, Association for Computing Machinery, New York, NY, USA, pp. 53–64. <http://dx.doi.org/10.1145/3302516.3307345>.
- Landi, W., Ryder, B.G., 1992. A safe approximate algorithm for interprocedural aliasing. *ACM SIGPLAN Not.* 27 (7), 235–248.
- Liang, D., Harrold, M.J., 1998. Slicing objects using system dependence graphs. In: Proceedings of the International Conference on Software Maintenance. ICSM '98, IEEE Computer Society, Washington, DC, USA, pp. 358–367. URL <http://dl.acm.org/citation.cfm?id=850947.853342>.
- Litvak, S., Dor, N., Bodik, R., Rinetzy, N., Sagiv, M., 2010. Field-sensitive program dependence analysis. In: Proceedings of the Eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering. FSE '10, Association for Computing Machinery, New York, NY, USA, pp. 287–296. <http://dx.doi.org/10.1145/1882291.1882334>.
- Lyle, J.R., 1984. Evaluating Variations on Program Slicing for Debugging (Data-Flow, Ada) (Ph.D. thesis). University of Maryland at College Park, USA.
- Majumdar, A., Drape, S.J., Thomborson, C.D., 2007. Slicing obfuscations: Design, correctness, and evaluation. In: Proceedings of the 2007 ACM Workshop on Digital Rights Management. DRM '07, ACM, New York, NY, USA, pp. 70–81. <http://dx.doi.org/10.1145/1314276.1314290>, URL <http://doi.acm.org/10.1145/1314276.1314290>.
- Maydan, D.E., Hennessy, J.L., Lam, M.S., 1991. Efficient and exact data dependence analysis. In: Proceedings of the ACM SIGPLAN 1991 Conference on Programming Language Design and Implementation. pp. 1–14. <http://dx.doi.org/10.1145/113445.113447>.
- Muchnick, S.S., 1997a. Advanced Compiler Design and Implementation. Morgan Kaufmann (Chapter 8.12).
- Muchnick, S.S., 1997b. Advanced Compiler Design and Implementation. Morgan Kaufmann (Chapter 12.2).
- Ochoa, C., Silva, J., Vidal, G., 2005. Lightweight program specialization via dynamic slicing. In: Proceedings of the 2005 ACM SIGPLAN Workshop on Curry and Functional Logic Programming. WFLP '05, ACM, New York, NY, USA, pp. 1–7. <http://dx.doi.org/10.1145/1085099.1085101>, URL <http://doi.acm.org/10.1145/1085099.1085101>.
- Ottenstein, K.J., Ottenstein, L.M., 1984. The program dependence graph in a software development environment. *SIGSOFT Softw. Eng. Notes* 9 (3), 177–184. <http://dx.doi.org/10.1145/390010.808263>, URL <http://doi.acm.org/10.1145/390010.808263>.
- Pugh, W., Wonnacott, D., 1992. Eliminating false data dependences using the omega test. In: Proceedings of the ACM SIGPLAN 1992 Conference on Programming Language Design and Implementation. PLDI '92, Association for Computing Machinery, New York, NY, USA, pp. 140–151. <http://dx.doi.org/10.1145/143095.143129>.
- Ramalingam, G., Field, J., Tip, F., 1999. Aggregate structure identification and its application to program analysis. In: Proceedings of the 26th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. POPL '99, Association for Computing Machinery, New York, NY, USA, pp. 119–132. <http://dx.doi.org/10.1145/292540.292553>.
- Reps, T., Horwitz, S., Sagiv, M., 1995. Precise interprocedural dataflow analysis via graph reachability. In: Proceedings of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. POPL '95, Association for Computing Machinery, New York, NY, USA, pp. 49–61. <http://dx.doi.org/10.1145/199448.199462>.
- Reps, T., Horwitz, S., Sagiv, M., Rosay, G., 1994. Speeding up slicing. *SIGSOFT Softw. Eng. Notes* 19 (5), 11–20.
- Rodrigues, N.F., Barbosa, L.S., 2005. Component identification through program slicing. In: In Proc. of Formal Aspects of Component Software, (FACS 2005). Elsevier ENTCS. Elsevier, pp. 291–304.
- Silva, J., 2012. A vocabulary of program slicing-based techniques. *ACM Comput. Surv.* 44 (3).
- Silva, J., Tamarit, S., Tomás, C., 2012. System dependence graphs in sequential erlang. In: Proceedings of the 15th International Conference on Fundamental Approaches to Software Engineering. FASE 2012, In: Lecture Notes in Computer Science (LNCS), vol. 7212, Springer, pp. 486–500.
- Snelting, G., Robschink, T., Krinke, J., 2006. Efficient path conditions in dependence graphs for software safety analysis. *ACM Trans. Softw. Eng. Methodol.* 15 (4), 410–457. <http://dx.doi.org/10.1145/1178625.1178628>.
- Späth, J., Ali, K., Bodden, E., 2019. Context-, flow-, and field-sensitive data-flow analysis using synchronized pushdown systems. *Proc. ACM Program. Lang.* 3 (POPL), <http://dx.doi.org/10.1145/3290361>.
- Sridharan, M., Fink, S.J., Bodik, R., 2007. Thin slicing. In: Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation. PLDI '07, Association for Computing Machinery, New York, NY, USA, pp. 112–122. <http://dx.doi.org/10.1145/1250734.1250748>.
- Tip, F., 1995. A survey of program slicing techniques. *J. Program. Lang.* 3 (3), 121–189.
- Tóth, M., Bozó, I., Horváth, Z., Lövei, L., Tejfel, M., Kozsik, T., 2010. Impact analysis of erlang programs using behaviour dependency graphs. In: Proceedings of the Third Summer School Conference on Central European Functional Programming School. CEFP '09, Springer-Verlag, Berlin, Heidelberg, pp. 372–390. URL <http://dl.acm.org/citation.cfm?id=1939128.1939139>.
- Walkinshaw, N., Roper, M., Wood, M., 2003. The Java system dependence graph. In: Proceedings Third IEEE International Workshop on Source Code Analysis and Manipulation. pp. 55–64.
- Weiser, M., 1981. Program slicing. In: Proceedings of the 5th International Conference on Software Engineering. ICSE '81, IEEE Press, Piscataway, NJ, USA, pp. 439–449.

Carlos Galindo received the bachelor's degree in computer science from the Universitat Politècnica de València (UPV), Valencia, Spain, in 2018, completing his bachelor's thesis abroad while taking part in a student exchange program with the Eidgenössische Technische Hochschule Zürich (ETHz), Switzerland, and the master's degree with specialization in software engineering, in 2019. He is currently working toward the Ph.D. degree in computer science, in the area of program slicing and static analysis at the VRAIN Research Institute, Valencia, Spain.

Jens Krinke is associate professor with the Software Systems Engineering Group, University College London, where he is director of the Centre for Research on Evolution, Search, and Testing (CREST). His main focus is software analysis for software engineering purposes. His current research interests include software similarity, modern code review, program analysis, and software testing. He is well known for his work on program slicing and clone detection.

Sergio Pérez received the Bachelor's degree in computer science engineering in 2012 from Universitat Politècnica de València (UPV), Valencia, Spain. He has been a member of the MIST Research Group since 2015. He completed his Ph.D. degree in the same university in 2023. He has mainly centered his research on debugging techniques, especially on program slicing and static analysis. His research interests also include regression testing, property-based testing, object-oriented programming, and functional programming.

Josep Silva received the B.S. and Ph.D. degrees in computer science from the Universitat Politècnica de València, Valencia, Spain, in 2001 and 2006, respectively, and the M.S. degree specialization in software engineering, in 2007. Since 2001, he has been a professor in three universities. Currently, he is an associate professor with the Universitat Politècnica de València (UPV), and an aggregate professor with UNED. In 2019, he joined the VRAIN Research Institute as a member of the MiST Research Group. His current research is centered on debugging, program slicing, and testing.