



SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Bachelor's Thesis in Informatics

**Combatting the Precision Loss of Partial
Contexts in Abstract Interpretation**

Felix Sebastian Kraye



SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Bachelor's Thesis in Informatics

**Combatting the Precision Loss of Partial
Contexts in Abstract Interpretation**

**Bekämpfung des Präzisionsverlustes durch
partielle Kontexte in Abstrakter
Interpretation**

Author:	Felix Sebastian Kraye
Supervisor:	Prof. Dr. Helmut Seidl
Advisor:	Michael Schwarz
Submission Date:	15th of February 2023

I confirm that this bachelor's thesis is my own work and I have documented all sources and material used.

Munich, 15th of February 2023

Felix Sebastian Kraye

Acknowledgments

Abstract

Contents

Acknowledgments	iii
Abstract	iv
1 Introduction	1
2 Background	2
2.1 Static Analysis	2
2.2 Flow sensitive analysis	3
2.3 Constraint systems	4
2.4 Interprocedural analysis	4
2.5 Context sensitivity	6
2.6 Partial context sensitivity	7
2.7 Precision loss	9
3 Combatting Precision Loss	11
3.1 Formal description	11
3.1.1 Taint analysis	11
3.1.2 Improving the values-of-variables analysis	13
3.2 Implementation	14
3.2.1 Taint analysis	14
3.2.2 Benefiting other analyses	19
4 Evaluation	24
4.1 Testing	24
4.2 Benchmarking	24
5 Conclusion	25
Abbreviations	26
List of Figures	27
List of Tables	28

Bibliography

29

1 Introduction

WIP:

- introduce GOBLINT
- show problem on a small example

Related work

Structure First we will introduce the basics of static analysis. This will go by introducing constraint systems and how these are used to gain information about the program statically. It will be accompanied by an example of a value-of-variables analysis acting on a toy language we will use for examples in this thesis. This will be extended to an interprocedural approach where partial context sensitivity will be introduced. Here the source of the precision loss will be pointed out. We then will propose an approach to combat this precision loss. The approach will first be introduced theoretically, after which we also present the challenges and results of implementing it in the GOBLINT analyzer. To give an evaluation to the proposed approach, a benchmark of the implementation will be performed and inspected. Our conclusions are presented in the last chapter.

2 Background

2.1 Static Analysis

Static analysis is defined by Rival [RY20] as "[...]an automatic technique that approximates in a conservative manner semantic properties of programs before their execution". This means that the program is analyzed just by the given source code without execution. The goal is to prove certain properties about the program in a "sound" manner, i.e., any property that is proven to hold actually does hold. However, from failing to prove a property one cannot conclude that the given property does not hold.

In order to prove properties, e.g. finding that a program does not contain races or identifying dead code, we need to gain information about the program. This is done by performing various kinds of analyses. We will focus on flow sensitive analyses from now on, i.e., analyses which find properties of the program dependent on the location within it. We will introduce a syntax to formalize flow sensitive analyses in the following sections. This formalization approach is heavily based on [ASV12].

// TODO: talk about Abstract Interpretation.

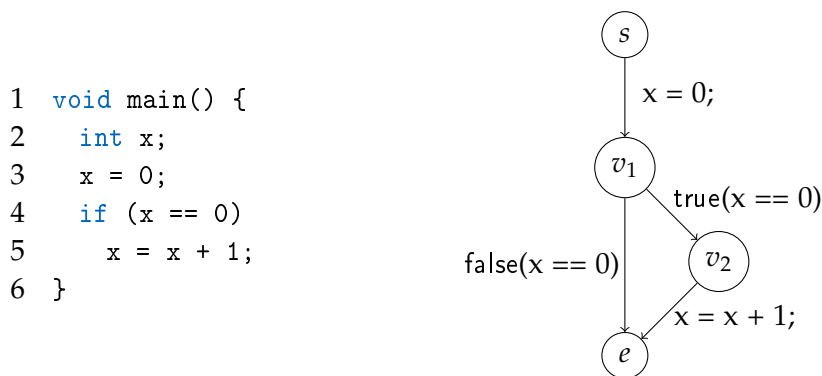


Figure 2.1: Example program (left) and corresponding CFG (right)

2.2 Flow sensitive analysis

As noted above flow sensitive analyses find properties of the program dependent on the point within the program. Expressed differently this means a flow sensitive analysis will find an overapproximation of states the program may be in for any given point within the program, from now on called "program point". This state can describe many things dependent on the analysis performed.

First let us define what a program point is: Consider a Control flow Graph (CFG), where nodes represent points between instructions within the program. Edges are labeled with instructions or checks (from now on collectively called "actions") and describe the transitions between these points (see example Figure 2.1). Then any node on this CFG is what we call a program point.

Concretely let N be the set of all program points. Furthermore, let \mathbb{D} be a Domain containing abstract states describing concrete states of the program. This means that some $d \in \mathbb{D}$ can describe many states the program can be in.

Then an analysis is expected to find a mapping $\eta : N \rightarrow \mathbb{D}$ which maps program points to abstract states describing that location within the program, i.e., for $[v] \in N$, $\eta [v]$ should be an abstract state describing all possible states (and possibly more) the program can be in at program point $[v]$.

As an example we will introduce a values-of-variables analysis for integers. This analysis finds a mapping from a set of program variables X to abstractions of their possible values at any given program point. Our toy language will support global variables (globals) as well as local variables (locals). The global variables can be accessed and changed by any procedure, while local ones are only visible to the procedure in which it was declared and can only be accessed and changed by this procedure. Therefore, our set X of variables is the disjoint union of globals G and locals L : $X = G \uplus L$. In the scope of this thesis we will focus on abstracting integer values by sets of integers. Thereby the goal of our values-of-variables analysis is to find a mapping $X \rightarrow 2^{\mathbb{N}}$ for each program point.

Combining this with the considerations from above, we chose the mapping $\mathbb{D}_v = X \rightarrow 2^{\mathbb{N}}$ as the Domain for the values-of-variables analysis. In summary this means that the resulting $\eta_v : N \rightarrow \mathbb{D}_v$ for this analysis describes a mapping $\eta_v [v]$ for some program point $[v] \in N$, where $\eta_v [v] x$ is a set containing all values $x \in X$ may possibly hold at $[v]$. From this we can conclude that x cannot hold any value outside $\eta_v [v] x$ at program point $[v]$.

2.3 Constraint systems

We now formulate a way in which we can describe an analysis in the form of constraints. For this we need a partial ordering \sqsubseteq on the domain \mathbb{D} .

Then we create a system of constraints which can be solved for a solution. Consider the edges (u, A, v) of the CFG, where each edge denotes a transition from program point $[u]$ to program point $[v]$ via the action A . Now let each of these edges give rise to a constraint

$$\eta [v] \sqsupseteq \llbracket A \rrbracket^\# (\eta [u])$$

where $\llbracket A \rrbracket^\#$ denotes the abstract effect of the action A defining our analysis. In addition, we need a start state. This is given by $\text{init}^\# : \mathbb{D}$ which is defined depending on the analysis. This gives rise to the start constraint $\eta [s] \sqsupseteq \text{init}^\#$ for the starting point of the program $[s] \in N$.

We will show these ideas with our example of the values-of-variables analysis: Let us define the partial ordering \sqsubseteq_v that is necessary for building the constraints. We will do this by stating that a mapping $M_1 \in \mathbb{D}_v$ is ordered below or equal to another mapping M_2 , if and only if for every variable $x \in X$, the set x is mapped to in M_1 is a subset of or equal to the one x is mapped to in M_2 . Formulated formally this is:

$$M_1, M_2 \in \mathbb{D}_v : M_1 \sqsubseteq_v M_2 \iff \forall x \in X : M_1 x \subseteq M_2 x$$

Next we define the start state $\text{init}^\# = M_\top$ for this domain as the mapping that maps every variable to the full set of integers \mathbb{N} , i.e., $\forall x \in X : M_\top x = \mathbb{N}$. This is because we assume variables to be randomly initialized in our toy language.

It remains to define the abstract effect of actions $\llbracket A \rrbracket_v^\#$ for our values-of-variables analysis. We will just show the effect of a simple variable assignment:

$$\llbracket x = y; \rrbracket_v^\# M = M \oplus \{x \mapsto (M y)\}$$

where $M \oplus \{x \mapsto s\}$ denotes that the mapping M is updated such that x will be mapped to the set s . A full definition of abstract effects of a values-of-variables analysis can be found at `</ / TODO >`.

2.4 Interprocedural analysis

So far we only have defined how a program without procedure calls is analyzed. Now we want to introduce procedure calls of the form $f()$. For simplicity, we will only

consider argumentless procedure calls without a return value. Arguments and return values can be simulated by using global variables.

Since a call has its own set of local variables to work with and a call stack can contain multiple of the same procedure (e.g. for recursion), we will analyze procedures in their own environment. However, we need to consider global variables and how the procedure affects these.

The idea is to give procedures their own starting states and analyze them similarly as we have done before. The final state of the called procedure is then used to be combined back into the state of the caller before the call. Formalized for an edge $(u, f();, v)$ this looks as follows:

$$\begin{aligned}\eta [s_f] &\sqsupseteq \text{enter}^\# (\eta [u]) \\ \eta [v] &\sqsupseteq \text{combine}^\# ((\eta [u]), (\eta [e_f]))\end{aligned}$$

where $[s_f]$ and $[e_f]$ are the start and end node of the CFG for procedure $f()$. The functions $\text{combine}^\# : \mathbb{D} \times \mathbb{D} \rightarrow \mathbb{D}$ and $\text{enter}^\# : \mathbb{D} \rightarrow \mathbb{D}$ are defined by the analysis. $\text{enter}^\#$ handles computing the start state for the procedure $f()$, while $\text{combine}^\#$ describes in what way the caller state and the end state of the callee are merged after the call.

It is worth mentioning at this point that even though a procedure can be called from multiple points within the program we still only analyze the procedure once. For n procedure calls $(u_n, f();, v_n)$ we get n constraints for $[s_f]$: $\eta [s_f] \sqsupseteq \text{enter}^\# (\eta [u_n])$. We can express this differently in a single constraint as follows:

$$\eta [s_f] \sqsupseteq \bigsqcup \{d \mid \exists (u_n, f();, v_n) \in \text{Edges} : \text{enter}^\# (\eta [u_n]) = d\}$$

where \bigsqcup is the least upper bound, i.e., the least $d \in \mathbb{D}$ according to the ordering \sqsubseteq that is ordered above all of its argument elements.

For our values-of-variables analysis we will show how $\text{enter}_v^\#$ and $\text{combine}_v^\#$ are defined. We need to take global variables into account when computing the start state and combining the caller state with the returned callee state after the call. Therefore, we define the two functions as follows:

$$\begin{aligned}\text{enter}_v^\# M &= M|_{\text{Globals}} \oplus \{x \mapsto \mathbb{N} \mid \forall x \in X\}|_{\text{Locals}_{ce}} \\ \text{combine}_v^\# (M_{cr}, M_{ce}) &= M_{cr}|_{\text{Locals}_{cr}} \oplus M_{ce}|_{\text{Globals}}\end{aligned}$$

where $M|_{\text{Locals}}$ and $M|_{\text{Globals}}$ refers to the mapping M restricted to only the local or global variables respectively. Note that Locals_{ce} refers to the locals of the callee while Locals_{cr} refers to the locals of the caller.

To explain these two functions let us first look at $\text{enter}_v^\#$. This function takes the part of the mapping from the caller that contains information about global variables and

adds the information of uninitialized local variables used in the procedure to the state. For $\text{combine}_v^\#$ the local part of the caller is kept, but it is updated with the global part of the callee return state, because the later contains the updated information about global variables after the procedure call.

2.5 Context sensitivity

In the previous chapter we approached the analysis of procedures by analyzing them only once with an abstract start state describing all possible concrete states the procedure could start with. We call this behavior "context insensitive" as the procedure is analyzed without differentiating between different states with which it is called.

This is not very precise as we will exemplify by applying the values-of-variables analysis to the program in Figure 2.2. We ignore the marked lines of the program for now. The procedure $\text{incr}()$ is called twice: Once with $a = 1$ in Line 10 and once with $a = -3$ in Line 13. This leads to two constraints for node $[s_{\text{incr}}]$:

$$\eta_v [s_{\text{incr}}] \sqsupseteq_v \text{enter}_v^\# \eta_v [v_2] = \{a \rightarrow \{1\}\}$$

$$\eta_v [s_{\text{incr}}] \sqsupseteq_v \text{enter}_v^\# \eta_v [v_5] = \{a \rightarrow \{-3\}\}$$

leading to $\eta_v [s_{\text{incr}}] = \{a \rightarrow \{-3, 1\}\}$. At the end point of the call the state will be $\eta_v [e_{\text{incr}}] = \{a \rightarrow \{-2, 2\}\}$, which is then combined with the states of nodes in the main procedure. Therefore, the state at Node $[v_6]$ will be $\{a \rightarrow \{-2, 2\}\}$, which is used to check the $\text{assert}(a < 0)$; in Line 14. The result of this assertion cannot be determined by the analysis even though it is easy for humans to see that it should hold.

This could have been avoided, if the procedure was analyzed twice, once with each starting state. To achieve this we will need to perform some modifications on our current approach: Instead of searching a mapping $\eta : N \rightarrow \mathbb{ID}$ we now seek $\eta : (N \times \mathbb{ID}) \rightarrow \mathbb{ID}$. This allows us to have different states for the same program point. We call the second part of $N \times \mathbb{ID}$ "context". For now this context will be the same as the starting state of the current procedure. Therefore, we need to adjust the constraints for $\text{enter}^\#$ and $\text{combine}^\#$:

$$\eta [s_f, \text{enter}^\# (\eta [u, d])] \sqsupseteq \text{enter}^\# (\eta [u, d])$$

$$\eta [v, d] \sqsupseteq \text{combine}^\# ((\eta [u, d]), (\eta [e_f, \text{enter}^\# (\eta [u, d])]))$$

The main procedure will always be analyzed just once, as in our toy language it is only called initially. The context for its nodes can be chosen arbitrarily.

There are no changes we need to perform on the values-of-variables analysis to make it context-sensitive. Solely the changes to the general analysis framework above suffice. Applying this changed analysis to the example in Figure 2.2 would lead to the procedure `incr()` being analyzed twice with different contexts, assuming we still ignore the marked lines. This leads to the following two entry constraints for different unknowns of the constraint system:

$$\begin{aligned}\eta_v [s_{incr}, \{a \rightarrow \{1\}\}] &\sqsubseteq_v \{a \rightarrow \{1\}\} \\ \eta_v [s_{incr}, \{a \rightarrow \{-3\}\}] &\sqsubseteq_v \{a \rightarrow \{-3\}\}\end{aligned}$$

For node v_6 only the state $\eta_v [e_{incr}, \{a \rightarrow \{-3\}\}] = \{a \rightarrow \{-2\}\}$ is combined with the caller state from before the call. With this information we can safely say that the assertion in the following Line 14 will hold.

2.6 Partial context sensitivity

While the context-sensitive approach from the previous section might be very precise, it can be quite costly in terms of computation time. To reach a middle ground between a context insensitive and a fully context-sensitive analysis, we change the approach so that contexts are different from the entry state of a call. With this we can group entry states by contexts to analyze a procedure multiple times. This time we group not once per individual entry state, but once per group of entry states.

Let \mathbb{C} be our context-domain. Instead of adding a state from the value-domain \mathbb{D} to the program points N we add an element from the context-domain. Thereby we change the definition of the mapping we want to compute as follows: We now have $\eta : (N \times \mathbb{C}) \rightarrow \mathbb{D}$.

We also define a new function $\text{context}^\# : \mathbb{D} \rightarrow \mathbb{C}$ here. This function will calculate the context when entering a procedure. Additionally, the constraints for $\text{enter}^\#$ and $\text{combine}^\#$ are changed as follows:

$$\begin{aligned}\eta [s_f, \text{context}^\# (\eta [u, c])] &\sqsubseteq \text{enter}^\# (\eta [u, c]) \\ \eta [v, c] &\sqsubseteq \text{combine}^\# ((\eta [u, c]), (\eta [e_f, \text{context}^\# (\eta [u, c])]))\end{aligned}$$

for an edge $(u, f(), v)$.

This formalization results in multiple constraints for a single contextualized starting variable $[s_f, c']$. We can alternatively formulate this as

$$\eta [s_f, c'] \sqsubseteq \bigsqcup \{ \text{enter}^\# (\eta [u_n, c_n]) \mid \exists (u_n, f(), v_n) \in \text{Edges} : \text{context}^\# (\eta [u_n, c_n]) = c' \}$$

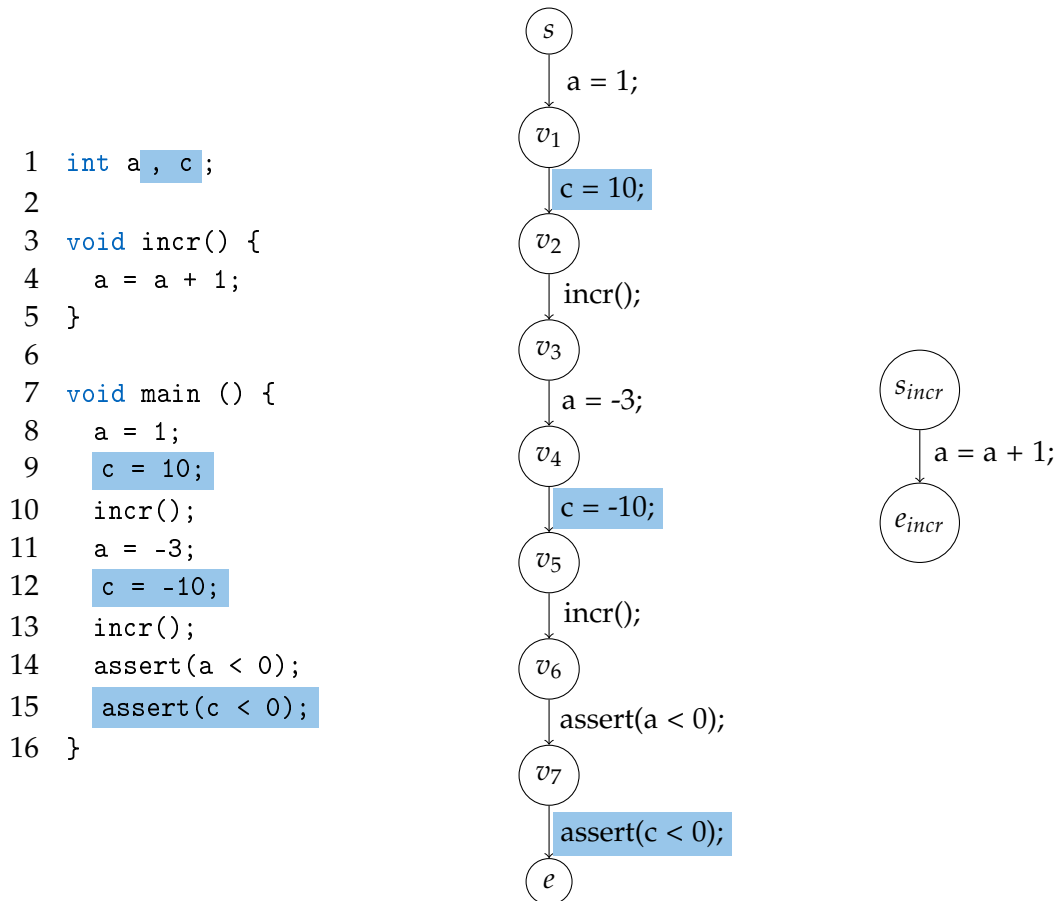


Figure 2.2: Example program (left) and corresponding CFGs for main (middle) and incr (right)

i.e., the constraint for the variable $[s_f, c']$ is the least upper bound of all entry states for some call of f , which have the same context c' as the constraint variable. Or expressed differently, all states computed by $\text{enter}^\# d$ for f are grouped by the context $c' = \text{context}^\# d$, where each group is joined by \sqcup to produce a constraint for a starting variable $[s_f, c']$ with the respective context.

With this formal model we have the option to perform an analysis completely context sensitively ($\mathbb{C} = \mathbb{ID}$ and $\text{context}^\# = \text{enter}^\#$), completely context insensitively ($\mathbb{C} = \{\bullet\}$) or anything in between. Note that we define $\{\bullet\}$ as the "unit domain" which contains exactly one element with the trivial ordering $\bullet \sqsubseteq \bullet$.

We have to note here that there are some severe issues with the approach for (partially) context-sensitive analyses described in this thesis: The resulting system of constraints may not be finite and some variables in the constraint system may depend on an infinite number of other variables. This can result in a very hard or even impossible-to-solve constraint system. While we will stick with this formal approach for this thesis for the sake of simplicity, an escape route to this issue is described in [ASV12].

2.7 Precision loss

The main source of the precision loss in context-insensitive or partially context-sensitive analyses is the join over all states with the same context, i.e., when we take the least upper bound of a group of entry states. Consider a procedure f that has no effect, i.e., $s_f = e_f$. Even for this procedure, the $\text{combine}^\#$ function receives the less precise result of the join \sqcup to combine it with the caller state. In this simple case, the result would be more precise if the $\text{combine}^\#$ function could directly use the result from the corresponding $\text{enter}^\#$ as the callee state for combining.

Even for procedures that do change the state, there might be some parts of the state which are untouched by the call. If we can identify these untouched parts, we could reduce the precision loss experienced by using partial contexts.

Let us clarify the source of the precision loss mentioned above with an example: For this we once again consider the example program Figure 2.2. This time we take the marked lines into account. When the program is analyzed context insensitively, not only does the state at the start node for $\text{incr}()$ s_{incr} represent two possible values for the global variable a , but also for c . Therefore, the state for this node is

$$\eta_v [s_{\text{incr}}, \bullet] = \{a \rightarrow \{-3, 1\}, c \rightarrow \{-10, 10\}\}$$

Even though the variable c is never changed within $\text{incr}()$, the mapping $c \rightarrow \{-10, 10\}$ is still copied into the caller state when combining the states for node v_6 . Thus, the

information gained by the context insensitive values-of-variables analysis does not suffice to determine the assertion in Line 15 to hold. This loss of precision could easily be avoided if we had some idea which global variables are definitely not changed by a procedure call.

3 Combatting Precision Loss

In this chapter we describe our approach to reduce the precision loss described in Section 2.7. First we use the syntax for flow sensitive analyses from Chapter 2 to formally define the idea. After that we explain the concrete implementation of the approach into the GOBLINT analyzer.

3.1 Formal description

3.1.1 Taint analysis

The basic idea to combat the precision loss is to track for each procedure which variables have been written or have possibly been altered in some other way. This information is then used in the values-of-variables analysis when combining the abstract state from the caller with the abstract return state given by the callee at the end of the procedure. In the following we call a variable that has been written or altered in the current procedure context "tainted". Therefore, we introduce a new taint analysis tracking which variables have been tainted within the context of the current procedure. It is worth mentioning that our notion of taintedness is related but different from other uses of the term "taint analysis".

Let us now formulate the syntax for the taint analysis we use in this thesis: Since we want to find a collection of tainted variables per program point, a suitable domain for this analysis is the powerset of the set of variables X ordered by the subset relation:

$$\mathbb{D}_t = 2^X \text{ with } \sqsubseteq_t = \subseteq$$

To be compatible with the notion of partially context-sensitive analyses from Section 2.6 we need to also specify a context domain \mathbb{C}_t which we define later. Note that we seek to compute a mapping from program points (with context) to sets of variables, i.e., $\eta_t : (N \times \mathbb{C}_t) \rightarrow \mathbb{D}_t$. To interpret this with the goal of our taint analysis in mind, we note that $\eta_t[n, \bullet] = T$ denotes that T is the set of possibly tainted variables at program point n . Expressed differently this means that for any variable $x \in T$ we cannot exclude that this variable was altered between the start of the current procedure up until the program point n . Note here that the tainted set T not only includes variables which

have been tainted by statements of the current procedure, but also variables which have been tainted within procedures called by the current one.

It remains to define \mathbb{C}_t , $\text{init}_t^\#$, $\text{enter}_t^\#$ and $\text{combine}_t^\#$ as well as the abstract effects of actions $\llbracket A \rrbracket^\#$. Recall that the notion of a "tainted" variable is defined in relation to the current procedure. This means we want to start without any variable being initially tainted when entering a procedure. It is worth pointing out that the entry to a procedure call does not depend on the state where it is called. Therefore, we design our analysis to be context-insensitive, i.e.,

$$\mathbb{C}_t = \{\bullet\} \text{ and trivially } \text{context}_t^\# T = \bullet$$

With these considerations we can also define $\text{enter}_t^\#$ and $\text{init}_t^\#$ as follows:

$$\text{enter}_t^\# T = \text{init}_t^\# = \emptyset$$

It is worth pointing out here that the function $\text{enter}_t^\# T$ is always equal to the empty set irregardless of its argument T . Therefore, it computes the same entry state for each call of a certain procedure.

When combining the caller state with the returned callee state, we note that we need to keep the tainted set from before the call, as a tainted variable can never get "untainted" again, no matter what the procedure does. Additionally, we add the tainted set returned by the callee, since anything tainted in the call needs to be considered tainted after the call as well. This is because we want to know which variables have been altered in a procedure call, no matter if the tainting happened within the procedure itself or within a further procedure call. This leaves us with the following equation for the $\text{combine}_t^\#$ function:

$$\text{combine}_t^\# (T_{cr}, T_{ce}) = T_{cr} \cup (T_{ce} \setminus \text{Locals}_{ce})$$

Note that we removed the callee local variables Locals_{ce} because these are not accessible by the caller and all of its callers anyway, so it is not useful to keep track of them.

Lastly we define the abstract effects of actions. Most of these (including checks) do not do anything besides propagating through the state from before. The only major exception are variable assignments. For these we note that the specific variable, which the value is assigned to is added to the tainted set. This is independent of the expression that evaluates to the assigned value, as we are only interested in the fact that the variable on the left of the assignment is altered. This leaves us with the following abstract effects of actions:

$$\llbracket A \rrbracket^\# T = \begin{cases} T \cup \{x\} & \text{if } A \equiv (x = e;) \\ T & \text{else} \end{cases}$$

where e is any arbitrary expression.

This concludes our definition of the taint analysis. In the following section we see how this information helps us to improve the values-of-variables analysis.

3.1.2 Improving the values-of-variables analysis

Recall the source of the precision loss we want to reduce. This happened when a global variable was updated with a less precise value after a procedure call even though this specific variable was not changed by the call.

Thanks to the taint analysis we defined in the previous section, we now do have the information which variables can be altered by a procedure $f()$ and which surely stay untouched. These are exactly those variables which are not in the tainted set of the end node $[e_f]$ for that procedure.

With this insight we can now update the $\text{combine}_v^\#$ function of our values-of-variables analysis as follows:

$$\text{combine}_v^\#(M_{cr}, M_{ce}) = M_{cr}|_{\text{Locals}_{cr} \cup (\text{Globals} \setminus T_{ce})} \oplus M_{ce}|_{\text{Globals} \cap T_{ce}}$$

where for an edge $(u, f();, v)$ we have $T_{ce} = \eta_t[e_f, c]$ for the respective context c to be combined.

Similar to before the $\text{combine}_v^\#$ function takes the caller mapping, restricts it to a subset of caller reachable variables and updates this mapping with the callee mapping restricted to the rest of caller reachable variables. In other words, the caller reachable variables are partitioned into two sets such that one subset is taken from the caller state while the other one is taken from the callee state. Before this change the partitioning was done strictly in such a way that the local variables were taken from the caller state and all global variables from the callee state. After this change, the global variables that are not tainted by the callee are also taken from the caller state and not from the callee anymore. Thereby the precision loss for untainted variables is eliminated.

One might wonder if this change could lead to a case, where the callee state has a more precise value for a variable that is discarded because this variable is not in the tainted set. Concretely this situation would be described by

$$\exists \text{Edge}(u, f();, v), x \in \text{Globals} : x \notin \eta_t[e_f] \wedge (\eta_v[e_f] x \subset \eta_v[u] x)$$

From $x \notin \eta_t[e_f]$ we know that x has not been altered in the procedure $f()$ since the node $[s_f]$, and therefore it holds that

$$\eta_v[e_f] x = \eta_v[s_f] x$$

By the definitions of \sqsubseteq_v and $\text{enter}_v^\#$ we get:

$$\eta_v[s_f] x \supseteq (\text{enter}_v^\#(\eta_v[u])) x = \eta_v[u] x$$

Therefore, $\eta_v[e_f] x \supseteq \eta_v[u] x$ which is a contradiction to the proposed case that we can therefore exclude.

3.2 Implementation

Before explaining the process of implementing the proposed taint analysis and its usage to improve other analyses, we introduce the GOBLINT analyzer and its structure in this paragraph. The core functionality of GOBLINT is to statically analyze C programs using an approach similar to the one described in Chapter 2. This generally works as follows: After the C input file is preprocessed, a CFG is generated. This graph is then used together with the specifications of various analyses to generate a constraint system. GOBLINT solves this constraint system and produces different kinds of outputs to the user according to the solution (e.g. notifications, warnings or a visualization of the full solution).

It is worth mentioning that GOBLINT can perform multiple analyses on a program at the same time. For this a compound domain is built (for the value domain as well as for the context domain), that is a tuple of all the domains of the analyses to be performed. To generate constraints, all activated analyses are taken into account where the specification of each analysis acts on its corresponding part of the compound domain. Information can be transferred between the different analyses via a system called "queries".

Figure 3.1 shows the inner structure of the analyzer. We can see that GOBLINT provides parametrized domains which can be used in the specifications of the analyses. It is also shown that multiple analyses are then combined into one MCP that is then used with the CFG to generate constraints which are solved.

For a deeper insight into the inner workings of GOBLINT refer to [Api14].

3.2.1 Taint analysis

To define an analysis the GOBLINT analyzer provides an interface, where the relevant parts can be seen in Figure 3.2. This interface requires two modules `D` and `C` which define the domain and the context-sensitive part of the domain. After that some functions are required:

- `name` to uniquely refer to an analysis.
- `startstate` to define the state used when entering the analysis (similar to `init#`).
- `query` to implement the query system of GOBLINT. This allows an analysis to broadcast information to be used in analyses.
- *Transfer functions* which define the abstract effects of actions (similar to $\llbracket A \rrbracket^{\#}$).
- *Functions for interprocedural analysis*

- *Function for analysis of multithreaded programs*

For our taint analysis we create a new module implementing this interface. As a name for GOBLINT internally we chose `taintPartialContexts` because `taint` was already used, and the name needs to be unique.

Domain

The next step is to choose D and C . According to the concept of our analysis described in Subsection 3.1.1 the domain should be a set of variables. However, we are now analyzing C instead of our toy language. In C not every left-hand side of an assignment is just a simple variable, but can be one of many more complex things, e.g., the memory location `*xptr` pointed to by the pointer `xptr`, the fourth place `a[3]` in an array `a`, the member `frac.n` of a struct `frac` and many more. This concept is called Left Value (of an assignment) (`lval`) and there is an implementation of this type provided by GOBLINT in the `Lval.CilLval` module. To be as precise as possible we use a set of `lvals` instead of a set of variables for the implementation of the taint analysis.

Another point worth mentioning is that we sometimes need the notion of "all variables" (or rather "all `lvals`") when we want to express that everything is tainted. While conceptually using the full set X poses no issue, in a concrete implementation this is extremely impractical and not even realizable if the set is infinitely large. For this case GOBLINT provides a parametrized domain `ToppedSet(Base)`. This domain is either a set of elements of the `Base` type or alternatively a `Top` element which can be interpreted as the "full set of all `Base` elements". Therefore, we finally have $D = \text{ToppedSet}(\text{Lval.CilLval})$ for our domain. Note that this also defines the ordering on the domain to be the regular subset ordering.

It remains to define the module C : We noted in Subsection 3.1.1 that our analysis by itself is context insensitive. Therefore, the context domain of our analysis C is empty, which is expressed with the `Unit` domain provided by GOBLINT. Note here, that this does not mean that the taint analysis is always performed context insensitively, i.e., only ran once per function. Since GOBLINT uses a compound domain, there may be other context-sensitive analyses, forcing the whole compound analysis to analyze a function multiple times. Our taint analysis however never contributes differing sub contexts to the compound context.

The `startstate` function

This function computes the initial state for our analysis similar to the `init#` function we introduced in Chapter 2. As discussed in Subsection 3.1.1, we implement this function so that it returns the empty set.

We note here however, that in practice we do not use the way, `startstate` is implemented in the scope of our thesis. In `GOBLINT` this function is called before the main function is even entered. Thus, `enter` (which we define later) is still used to compute the entry state for the main function. We chose to implement `startstate` in this way for consistency.

Transfer functions

These implement the effect of actions on the state, similar to the abstract effects of actions $\llbracket A \rrbracket^\#$ in Chapter 2. Variable declarations are handled by `vdecl` while `branch` handles checks for if-statements and loops. For these two actions our analysis just propagates the state from before, so the two mentioned functions use the default implementation from the `Analysis.IdentitySpec` of `GOBLINT`.

Much more interesting is the case of the `assign` function which handles the effect of an assignment to an `lval`. For this case we want to add the `lval` to our tainted set. The parameters for the `assign` function are: `ctx` which amongst other things contains the state from before, the `lval` to which a value is assigned and an expression that evaluates to the value that is assigned. We are only interested in `ctx` and the `lval`, as to us only the fact that a value is assigned is relevant and not its concrete value.

Tainting `lvals` is not as straightforward as it might seem at first. Just adding it to the state from before, i.e., the tainted set, only suffices if the `lval` is a specific location in the memory, e.g., a specific (local or global) variable. The `lval` could however also be a reference to a location in the memory, e.g., a pointer. For these it is not helpful to just taint the reference because we need to know the specific memory locations that are or could be tainted. To solve this issue we make use of `GOBLINT`'s `MayPointTo` query. This takes a reference to the memory and asks all other activated analyses if they have any information about where this reference may point to. Just like everything else in the static analyzer `GOBLINT`, the answer is an overapproximation, so we can be sure not to miss any location that could be referenced.

In conclusion, tainting an `lval` goes as follows: If the `lval` is a specific memory location, this `lval` is added to the tainted set. If it is a reference to the memory described by some expression, send a `MayPointTo` query to ask other analyses which memory locations this expression may point to and add the returned set of `lvals` to the tainted set. We implemented this functionality in a helper function `taint_lval`. Therefore, calling this function is the only thing the `assign` function needs to do as seen in Figure 3.3.

Functions for interprocedural analysis

Here we define the functions `context`, `enter` and `combine`. These functions work similar to their abstract counterparts as described in Chapter 2. In addition to these known functions, the interface also requires two additional functions: `return` which handles return statements right before a function is left and `special` which handles calls to library functions or other functions, for which we do not have the source code we can analyze.

Our implementation does not differ a lot from the proposed formal description in Subsection 3.1.1. The only major difference is that we need to handle return values, as these taint the lval they are assigned to. Therefore, we implement them as follows:

context: Our context domain is the `Unit` domain, and we do not want to generate different contexts for this analysis. Thus, this function returns the unit element.

enter: This function just returns the empty set as the entry state for the called function like we discussed in Subsection 3.1.1.

combine: The `combine` function first checks if there is an lval the return value is assigned to. If so it taints this respective lval in the caller state using the helper function `taint_lval` introduced in the "Transfer Functions" section. After that it computes the union of the resulting state with the returned callee state and returns it.

Summarized, the result of this function computes the union of both states it receives for combining and additionally adds lvals which are possibly tainted by the return value.

return: In our formal description in Subsection 3.1.1, the `combinet#` removed variables unreachable by the caller. In the concrete implementation, we gave this functionality to the `return` function, so the removal happens right before the `combine`. Since in C we also have function arguments, we also remove the ones of the function we are returning from.

It is worth pointing out that we do not just remove all lvals corresponding to local variables or arguments. A function might exist multiple times in the current call stack, resulting in the existence of multiple versions of the same local variable. GOBLINT treats these as being the same variable. Thus, when we remove a local variable we risk also removing a different version of it lower in the call stack, for which we still need the taintedness information. To address this issue, `return` sends an `IsMultiple` query for each variable to be removed and only removes those, that surely not have multiple versions. This query is already provided by GOBLINT.

special: This function addresses library functions or other functions, for which we do not have the source code to analyze. The simple way to handle these, is to just return Top, i.e., saying "everything could be tainted", after a special call.

This is how we handle unknown functions, however GOBLINT provides "Library Descriptors", which contain information about some known C library functions, e.g., `printf`, `malloc`, `cos`, etc. With the respective Library Descriptor of a function, we can gain information about which addresses are "shallowly" written and which are "deeply" written by the call. Shallowly written addresses point to lvals which might be directly written. Deeply written addresses however point to lvals where not only the lval itself, but possibly anything it might recursively point to could be written. Therefore, the special function makes use of GOBLINT's `MayPointTo` and `ReachableFrom` queries in the following way:

First the function checks if a Library Descriptor is available. If not, Top is returned. Otherwise, the shallowly and deeply written addresses are obtained from the Descriptor. Consequently, the union of

- the state before the call
- anything that is possibly tainted by the return value (using `taint_lval` like in `combine`)
- the set of lvals returned by the `MayPointTo` query for any shallowly written address
- the set of lvals returned by the `ReachableFrom` query for any deeply written address

is returned by the special function.

Function for analysis of multithreaded programs

To be able to analyze multithreaded programs, GOBLINT's analysis interface requires the following functions: `threadenter` to compute the startstate for the newly created thread and `threadspawn` which computes the effect of a thread creating instruction to the state of the creating thread.

We implement the former of these two functions similarly to our `startstate` and `enter`. Thus, `threadenter` returns the empty set. In practice, we do not use the fact that we enter threads with an in the scope of this Thesis. We chose to implement `threadenter` in this way for consistency.

To implement the other function, `threadspawn`, we consider how a thread creation effects the state of the creator. We note that for our notion of taintedness the only

relevant effect is, that the thread creating function may write thread ID variables to which it receives a reference as an argument. Thus, this function uses the helper function `taint_lval` defined in the "Transfer Functions" section to add possibly tainted lvals to the state from before and returns the result.

The query function

We want to enable our taint analysis to tell other analyses which lvals are tainted at a specific program point. Therefore, we add a new query `MaybeTainted` to the query system of GOBLINT. The result of this query should be a set containing lvals which may be tainted, i.e., any lval not in the returned set is definitely untainted.

After this addition we are able to make our `taintPartialContexts` analysis answer to this query. Therefore, our analysis implements the query function in such a way that it answers only to `MaybeTainted` queries with the current state but does not answer other queries.

3.2.2 Benefiting other analyses

In this section we discuss how we improved other existing analyses in GOBLINT using the taint analysis we implemented in Subsection 3.2.1. The main analysis that benefited from these changes is the base analysis of the analyzer. This analysis implements a very much extended approach of the basic values-of-variables analysis we formally defined in Chapter 2. The base analysis is however still based on the main goal and basic concept of finding a mapping from program variables to possible values at each program point. Therefore, this analysis uses a mapping from variables to their possible values as part of its domain. However, here the `ValueDomain` of the mapping is much more complex than just a set of possible integers. It provides abstractions for virtually any type in C, including arrays, structs and pointers. Even more though, the `ValueDomain` is highly configurable. Amongst other options it allows choosing between different ways of abstracting integer values or arrays. One interesting option related to the topic of this thesis is the possibility to choose between different degrees of context sensitivity: the analysis can be fully context-sensitive, insensitive with respect to integer variables (abstracted by intervals or in general), only sensitive with respect to pointers or completely insensitive. When choosing anything but the completely context-sensitive option, this analysis experiences the (avoidable) loss of precision described in Section 2.7.

To reduce this loss we need to change the `combine` function of the base analysis so that it uses the results of our `taintPartialContexts` analysis. Let us first describe how the `combine` function was implemented before our changes:

1. The return value is saved. Its value is removed from the callee state.
2. All globals are removed from the caller state.
3. Everything from the callee state is added to the caller, possibly overwriting caller values. This excludes the return value which is handled separately.
4. Some further adjustments according to the configuration are performed to the resulting state.
5. The saved return value is added to the state before it is returned.

To implement our changes we will focus on the steps 2 and 3, where the caller mapping is updated. The other steps will remain the same.

The core idea to implement the concept proposed in Subsection 3.1.2 is as follows: First we get the set of possibly tainted lvals from the callee. We then iterate over its elements one by one, where for each tainted lval we update the caller mapping with the corresponding value from the callee mapping, i.e., we get the value corresponding to that lval from the callee mapping and set the lval to map to this value in the caller mapping. This functionality of updating the caller mapping with the callee mapping using the tainted set is implemented in a helper function `combine_st`.

Before we explain how this helper function is implemented, we first show how it is embedded in the current implementation of the `combine` function. As discussed, we alter steps 2 and 3: First we send a `MayPointTo` query to the return state of the callee. We then check if the query returned the Top set, i.e., the notion that everything is tainted. In this case we perform the unchanged steps 2 and 3 just like before. Amongst other cases, this can happen, when the base analysis is run without our taint analysis being activated.

Let us now define what happens, when the result of the query is an explicit set. Before calling the helper function `combine_st`, we need to handle two special cases here:

- For a global variable, there is no mapping in the callee state, but there is one in the caller state. This case can occur in multithreaded mode, if this variable was protected by a mutex before the call, but the mutex was released in the called function. In this case, the mapping for this variable would be removed from the state within the callee. In the `combine` function, we do need to take the information from the callee for such a variable, i.e., remove it from the caller mapping. Therefore, we filter over the caller mapping and remove all globals, for which there exists no mapping in the callee mapping.
- For a global variable, there is a mapping in the caller state, but there is none in the callee state. This case can occur if new information is gained within the call, e.g.,

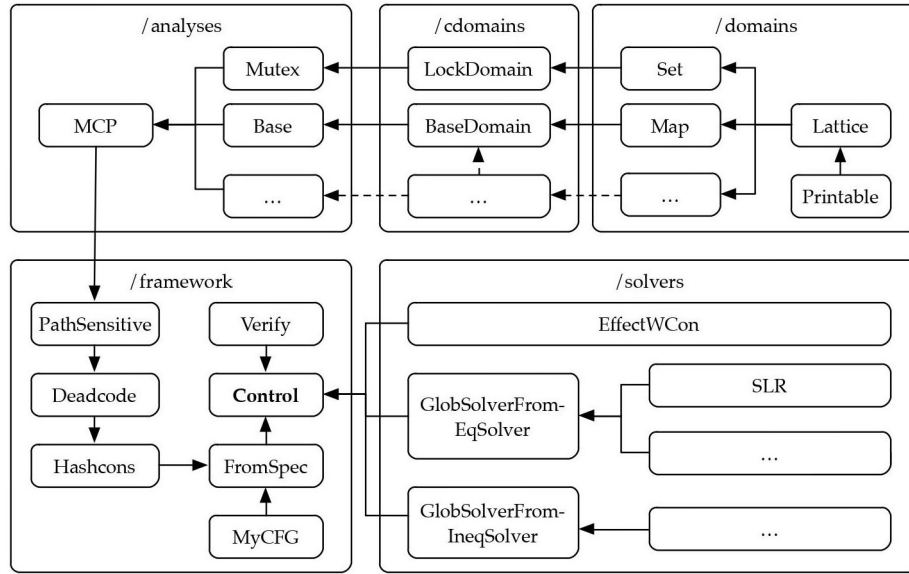


Figure 3.1: Schematic directory structure of GOBLINT. Adapted from [Api14]

some new memory is allocated. This information is not tracked by the tainted set and would therefore not be copied into the caller state. Since we still want to have this new information after the combine, we add all these mappings from the callee to the caller.

These cases have to be handled separately, as for these the corresponding `lval` is not necessarily contained in the tainted set. After the two special cases are handled, we use the `combine_st` helper function to finally update the tainted `lvals` in the caller state. We then proceed with the resulting state to the steps 4 and 5 like before.

We note here that we added a new parameter `f_ask` to the `combine` function. To do this we had to update the analysis interface and consequently all analyses implementing it. This new parameter allows us to send queries to the returned callee state, which was not possible before.

combine_st: This helper function takes the caller state (updated according to the two special cases), the callee state and the set of tainted `lvals`.

- relation analysis (apron) benefited in a similar way
- mention `varEq` and `condVars` for completion??

- Full New Section: ThreadCreate analysis

```
1 module type Spec =
2 sig
3   (* Domain *)
4   module D : Lattice.S
5   module C : Printable.S
6
7   val name : unit -> string
8   val startstate : varinfo -> D.t
9   val query : (D.t, C.t) ctx -> 'a Queries.t -> 'a Queries.result
10
11   (* Transfer functions *)
12   val assign: (D.t, C.t) ctx -> lval -> exp -> D.t
13   val vdecl : (D.t, C.t) ctx -> varinfo -> D.t
14   val branch: (D.t, C.t) ctx -> exp -> bool -> D.t
15
16   (* Functions for interprocedural analysis *)
17   val special : (D.t, C.t) ctx -> lval option -> varinfo -> exp list -> D.t
18   val enter : (D.t, C.t) ctx -> lval option -> fundec -> exp list -> (D.t * D.t) list
19   val return : (D.t, C.t) ctx -> exp option -> fundec -> D.t
20   val combine : (D.t, C.t) ctx -> lval option -> exp -> fundec -> exp list -> C.t option -> D.t
21
22   val context : fundec -> D.t -> C.t
23
24   (* Function for analysis of multithreaded programs *)
25   val threadenter : (D.t, C.t) ctx -> lval option -> varinfo -> exp list -> D.t list
26   val threadspawn : (D.t, C.t) ctx -> lval option -> varinfo -> exp list -> (D.t, C.t) ctx
27 end
```

Figure 3.2: Simplified Interface for implementing analyses in GOBLINT

```
1 let taint_lval ctx (lval:lval) : D.t =
2   let d = ctx.local in
3   (match lval with
4   | (Var v, offs) -> D.add (v, resolve offs) d
5   | (Mem e, _) -> D.union (ctx.ask (Queries.MayPointTo e)) d
6   )
7
8 let assign ctx (lval:lval) (rval:exp) : D.t =
9   taint_lval ctx lval
```

Figure 3.3: Implementation of the helper `taint_lval` and the `assign` function

4 Evaluation

4.1 Testing

- soundness checked with regression tests from GOBLINT

4.2 Benchmarking

- coreutil as benchmarking programs
- various analysis runs performed with goblint: ctx insensitive with and without taint, precision compared, checks passing compared.

5 Conclusion

Abbreviations

CFG Control flow Graph

lval Left Value (of an assignment)

List of Figures

2.1	Example program (left) and corresponding CFG (right)	2
2.2	Example program (left) and corresponding CFGs for main (middle) and incr (right)	8
3.1	Schematic directory structure of GOBLINT. Adapted from [Api14]	21
3.2	Simplified Interface for implementing analyses in GOBLINT	22
3.3	Implementation of the helper <code>taint_lval</code> and the <code>assign</code> function . . .	23

List of Tables

Bibliography

- [Api14] K. Apinis. “Frameworks for analyzing multi-threaded C.” PhD thesis. Technische Universität München, 2014.
- [ASV12] K. Apinis, H. Seidl, and V. Vojdani. “Side-effecting constraint systems: a swiss army knife for program analysis.” In: *Asian Symposium on Programming Languages and Systems*. Springer. 2012, pp. 157–172.
- [RY20] X. Rival and K. Yi. *Introduction to static analysis: an abstract interpretation perspective*. Mit Press, 2020.