

Pointer Analysis, Conditional Soundness, and Proving the Absence of Errors*

Christopher L. Conway¹, Dennis Dams², Kedar S. Namjoshi², and
Clark Barrett¹

¹ New York University, Dept. of Computer Science

{cconway,barrett}@cs.nyu.edu

² Bell Laboratories, Alcatel-Lucent

{dennis,keedar}@research.bell-labs.com

Abstract. It is well known that the use of points-to information can substantially improve the accuracy of a static program analysis. Commonly used algorithms for computing points-to information are known to be sound only for memory-safe programs. Thus, it appears problematic to utilize points-to information to verify the memory safety property without giving up soundness. We show that a sound combination is possible, even if the points-to information is computed separately and only conditionally sound. This result is based on a refined statement of the soundness conditions of points-to analyses and a general mechanism for composing conditionally sound analyses.

1 Introduction

It is well known that information about pointer relationships is essential for effective analysis and optimization of C programs [2, 18]. Such information can be provided by a variety of algorithms that compute an approximation of the *points-to* relations of a program (e.g., [3, 27, 12]). For variables x and y , x *points to* y if there is some execution of the program such that the value of x is either the address of y or, if x and y are aggregate objects (such as arrays or structures), the value of an element of x is an address within the extent of y .

A (may) points-to analysis is sound if the relation it computes over-approximates the true points-to relation of the program. Typical analysis algorithms are known to be sound only for “well-behaved” programs, i.e., programs with behavior that is well-defined by the C standard [22]. For instance, typical points-to analysis algorithms consider the points-to sets of pointer values x and $x+1$ to be the same. This is justified if $x+1$ does not “overflow” the bounds of the object pointed to by x . However, if the expression does overflow (i.e., the program is not “well-behaved”), the object pointed to by $x+1$ is undefined.

In extending the Orion static analyzer [11] to verify memory safety, we found that performing the analysis without access to points-to information resulted in an overwhelming number of false alarms. In principle, a single, combined analysis can be defined that computes memory safety and points-to information simultaneously. Since the memory safety information being computed in one “half” is available to the points-to “half,” the points-to information is kept sound even for ill-behaved executions. Conversely, the memory safety component has access to up-to-date points-to information, enabling a higher precision analysis.

However, such a fine-grained combination of analyses may not be scalable. We would like to treat existing scalable points-to analyses as plug-in components. Moreover, one may wish to perform the memory safety analysis in a “bottom-up” fashion, computing a general summary for each function on any possible input—in this case, separately computed points-to information helps limit the possible values of pointer parameters and global variables.

These considerations lead to the central questions addressed in this paper: Is it possible to obtain a sound combination of independent points-to and memory safety analyses, especially as the first obtains sound results assuming memory safety? More generally, under what conditions can conditionally sound analyses be combined? What guarantees can be made for the combination?

This paper makes several contributions:

* This is an expanded version of a paper presented at SAS 2008 [5].

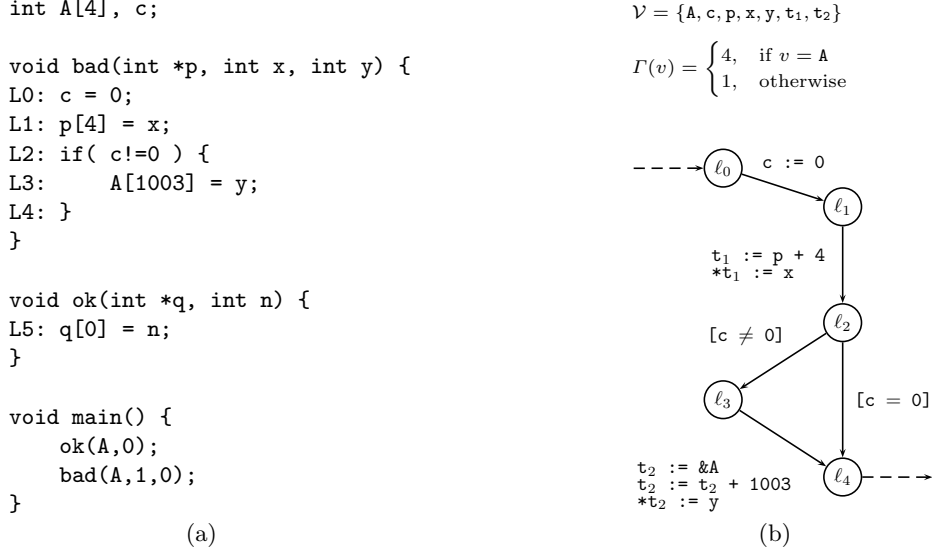


Fig. 1. An unsafe C program

- We formalize the notion of conditional soundness, show how to compose conditionally sound analyses, and derive the conditional soundness guarantee of the composition. Although we describe conditional soundness in the specific context of points-to analysis and a particular kind of memory safety, we believe our framework can be used to refine the soundness results of a variety of static analyses, e.g., analyses that are sound assuming sequential consistency or numerical analyses that are sound assuming the absence of integer overflows. Conditional soundness can be formulated in terms of the (unconditional) preservation of a class of temporal safety properties using Cousot and Cousot’s power construction [7–9]. Our formulation, while more specialized, is simpler (e.g., it is state-based rather than path-based) and captures the behavior of several interesting analyses more directly.
- We show that a set of points-to analyses similar to and sharing the soundness properties of commonly-used flow-sensitive and insensitive analyses—such as those of Emami, Ghiya, and Hendren [16]; Wilson and Lam [28]; Andersen [3]; Steensgaard [27]; and Das [12]—provide results that are sound for any *memory-safe execution* of a program. This statement is both stronger and more precise than the traditional statement that such analyses are sound for “well-behaved” programs.
- This more precise characterization of a points-to analysis, along with the combination theorem for conditional analyses, shows that the combination of an independent points-to analysis with a memory safety analysis is conditionally sound. The soundness result guarantees that the *absence* of errors can be proved. Conversely, for a program with memory errors, at least one representative error—but not necessarily all errors—along any unsafe execution will be detected.

Motivating Example. Figure 1(a) is an example of a program which is *not* well-behaved: there is an off-by-one error at label L1 and an off-by-one-thousand error at label L3 when c is not 0. Assume that the functions `ok` and `bad` are analyzed in a bottom-up fashion, without reference to the actual parameters supplied in `main`. Without any points-to information regarding q , the only safe assumption is that the expression $q[0]$ at label L5 can alias any location in memory—a conservative memory safety analysis would be forced to assume that the behavior of the program is undefined from this point on. But this is not the case: the function `ok` is memory-safe so long as q points to an array of at least one element. Points-to information is necessary to obtain a precise memory safety analysis.

A typical points-to analysis (e.g., Andersen’s [3]) will determine that p and q both point to A and not to c , n , x , or y . Using points-to information, a memory safety analysis can (correctly) infer that $q[0]$ is an

in-bounds location at L5 and, thus, the execution of `ok` is well-defined. Further, it can (correctly) detect that `p[4]` is out-of-bounds at L1 and emit a useful error report.

In many implementations `p[4]` will alias `c` at L1—so that `c` is set to 1, making L3 reachable—but `p` will not point to `c` according to the points-to relation. Since `c` is initialized to 0 in `bad` and—according to the points-to relation—no expression aliasing `c` is subsequently assigned to, the analysis is likely to (incorrectly) infer that the error at L3 is unreachable. Thus, it may seem that relying on the points-to relation will lead a static analyzer to miss real errors. Note, though, that the reachability of the error at L3 is due solely to the unsafe assignment at L1: the points-to relation can be relied upon up to the first occurrence of a memory safety error.

This line of reasoning is not specific to the example: as we will show, it applies to any conditionally sound analysis, enabling such an analysis to detect at least one error along any erroneous execution.

2 Program Analysis and Conditional Soundness

To present program analysis in a formal setting, we use the framework of abstract interpretation [6]. A full syntax of program statements is given in the next section. For now, we are concerned only with the relationship between concrete and abstract interpretations. We omit any discussion of techniques (such as widening and extrapolation) which serve to make program analyses finite and computable—we are concerned solely with issues of soundness.

Let C be a distinguished set of *concrete states*. A *domain* (D, γ) is a pair, where D is a set of *abstract states* and $\gamma : D \rightarrow 2^C$ is a *concretization function*. When the meaning is clear, we overload D to refer both to a domain and to its underlying set of states and use γ_D to refer to the concretization function. We lift γ_D to sets of states: $\gamma_D(D') = \bigcup_{d \in D'} \gamma_D(d)$, where $D' \subseteq D$. We say a set $D' \subseteq D$ *over-approximates* $C' \subseteq C$ iff $\gamma_D(D') \supseteq C'$. Similarly, a function $f_D : D \rightarrow 2^D$ *over-approximates* $f_C : C \rightarrow 2^C$ iff $f_D(d)$ over-approximates $f_C(c)$ whenever c is in $\gamma_D(d)$.

We define the soundness of a program interpretation in terms of a *collecting semantics*. Given a (concrete or abstract) domain D , we will define a *semantic operator* $\llbracket \cdot \rrbracket$ which maps a program \mathcal{P} to a set $\llbracket \mathcal{P} \rrbracket \subseteq D$ of *reachable states*. The semantics $\llbracket \mathcal{P} \rrbracket$ is defined inductively in terms of *semantic interpretations over D* : a set $\mathcal{I}[\mathcal{P}] \subseteq D$ of *initial states* and a *transfer function* $\mathcal{F}[\mathcal{P}] : D \rightarrow 2^D$. We lift $\mathcal{F}[\mathcal{P}]$ to sets of states: $\mathcal{F}[\mathcal{P}](D') = \bigcup_{d \in D'} \mathcal{F}[\mathcal{P}](d)$, where $D' \subseteq D$.

An *analysis* \mathcal{A} is represented as a tuple $(D, \mathcal{I}, \mathcal{F})$, where D is a domain and \mathcal{I} and \mathcal{F} are semantic interpretations over D . We use $D_{\mathcal{A}}$, $\mathcal{I}_{\mathcal{A}}$, and $\mathcal{F}_{\mathcal{A}}$ to denote the constituents of an analysis \mathcal{A} and $\gamma_{\mathcal{A}}$ to denote the concretization function of the domain $D_{\mathcal{A}}$.

Definition 1. Let $\mathcal{A} = (D, \mathcal{I}, \mathcal{F})$ be an analysis. The k -reachability predicate $R_{\mathcal{A}}^k[\mathcal{P}]$ for a program \mathcal{P} w.r.t. \mathcal{A} holds if a state is reachable in \mathcal{A} in exactly k steps. We define $R_{\mathcal{A}}^k[\mathcal{P}]$ inductively as a subset of D :

$$R_{\mathcal{A}}^0[\mathcal{P}] = \mathcal{I}[\mathcal{P}] \qquad R_{\mathcal{A}}^k[\mathcal{P}] = \mathcal{F}[\mathcal{P}](R_{\mathcal{A}}^{k-1}[\mathcal{P}]), \quad k > 0$$

The semantics $\llbracket \cdot \rrbracket_{\mathcal{A}}$ w.r.t. \mathcal{A} maps a program \mathcal{P} to a subset of D , the reachable states in \mathcal{P} w.r.t. \mathcal{A} :

$$\llbracket \mathcal{P} \rrbracket_{\mathcal{A}} = \bigcup_{k \geq 0} R_{\mathcal{A}}^k[\mathcal{P}]$$

To judge the soundness of an analysis, we need a concrete semantics against which it can be compared. The *concrete domain* D_C is given by the pair (C, γ_C) , where γ_C is the trivial concretization function: $\gamma_C(c) = \{c\}$. We assume that a *concrete analysis* $\mathcal{C} = (D_C, \mathcal{I}_C, \mathcal{F}_C)$ is given. The concrete analysis uniquely defines a *concrete semantics* $\llbracket \cdot \rrbracket_{\mathcal{C}}$.

Definition 2. An analysis \mathcal{A} is *sound* iff for every program \mathcal{P} , $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$ over-approximates $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}}$ (i.e., $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}) \supseteq \llbracket \mathcal{P} \rrbracket_{\mathcal{C}}$).

It is sufficient for \mathcal{A} to have sound semantic interpretations.

Definition 3. Let \mathcal{I}_D be a semantic interpretation over a domain D . \mathcal{I}_D is sound iff for every program \mathcal{P} , $\mathcal{I}_D[\mathcal{P}]$ over-approximates $\mathcal{I}_C[\mathcal{P}]$. \mathcal{F}_D is sound iff for every program \mathcal{P} , $\mathcal{F}_D[\mathcal{P}]$ over-approximates $\mathcal{F}_C[\mathcal{P}]$.

Theorem 1. Let \mathcal{A} be an analysis. If $\mathcal{I}_\mathcal{A}$ and $\mathcal{F}_\mathcal{A}$ are sound, then \mathcal{A} is sound.

Proof. See Theorem 3 and Note 1, below. □

Conditional Soundness. So far, we have defined a style of analysis which is *unconditionally sound*, mirroring the traditional approach to abstract interpretation. However, as we have noted, points-to analysis is sound only under certain assumptions about the behavior of the program analyzed. To address this, we introduce the notion of *conditional soundness with respect to a predicate θ* . An analysis will be *θ -sound* if it over-approximates the concrete states of a program that are reachable via *only θ -states*. We first define a semantics restricted to θ .

Definition 4. Let $\mathcal{A} = (D, \mathcal{I}, \mathcal{F})$ be an analysis and θ a predicate on D (we view the predicate θ , equivalently, as a subset of D). The θ -restricted k -reachability predicate $R_{\mathcal{A}\downarrow\theta}^k[\mathcal{P}]$ for program \mathcal{P} w.r.t. \mathcal{A} holds if a state is reachable in \mathcal{A} in exactly k steps via only θ states. $R_{\mathcal{A}\downarrow\theta}^k[\mathcal{P}]$ is defined inductively:

$$R_{\mathcal{A}\downarrow\theta}^0[\mathcal{P}] = \mathcal{I}[\mathcal{P}] \qquad R_{\mathcal{A}\downarrow\theta}^k[\mathcal{P}] = \mathcal{F}[\mathcal{P}](\theta \cap R_{\mathcal{A}\downarrow\theta}^{k-1}[\mathcal{P}]), \quad k > 0$$

The θ -restricted semantics $\llbracket \cdot \rrbracket_{\mathcal{A}\downarrow\theta}$ w.r.t. \mathcal{A} maps a program \mathcal{P} to a subset of D , the θ -reachable states in \mathcal{P} w.r.t. \mathcal{A} :

$$\llbracket \mathcal{P} \rrbracket_{\mathcal{A}\downarrow\theta} = \bigcup_{k \geq 0} R_{\mathcal{A}\downarrow\theta}^k[\mathcal{P}]$$

Note that $R_{\mathcal{A}\downarrow\theta}^k[\mathcal{P}]$ may include non- θ states—neither $\mathcal{I}[\mathcal{P}]$ nor the range of $\mathcal{F}[\mathcal{P}]$ are restricted to θ —but those states will not yield successors in $R_{\mathcal{A}\downarrow\theta}^{k+1}[\mathcal{P}]$. The θ -restricted semantics give us a lower bound for the approximation of a θ -sound analysis.

Definition 5. Let \mathcal{A} be an analysis and θ a predicate on C . \mathcal{A} is θ -sound iff for every program \mathcal{P} , $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$ over-approximates $\llbracket \mathcal{P} \rrbracket_{C\downarrow\theta}$.

Note that an unconditionally sound analysis is also θ -sound for *any* θ . More generally, any θ -sound analysis is also φ -sound, for any φ stronger than θ .

This notion of conditional soundness does not just give us a more precise statement of the behavior of certain analyses—it provides us with a sufficient condition to show an analysis *proves the absence of error states*.

Theorem 2. Let \mathcal{P} be a program and \mathcal{A} a θ -sound analysis. If there are no reachable non- θ states in \mathcal{P} w.r.t. \mathcal{A} , then there are no reachable concrete non- θ states in \mathcal{P} (i.e., if $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}) \subseteq \theta$, then $\llbracket \mathcal{P} \rrbracket_C \subseteq \theta$).

Proof. Since \mathcal{A} is θ -sound, we know that $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$ over-approximates $\llbracket \mathcal{P} \rrbracket_{C\downarrow\theta}$ (i.e., $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}) \supseteq \llbracket \mathcal{P} \rrbracket_{C\downarrow\theta}$). By assumption, $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}) \subseteq \theta$. Hence, $\llbracket \mathcal{P} \rrbracket_{C\downarrow\theta} \subseteq \theta$. It suffices to show: if $\llbracket \mathcal{P} \rrbracket_{C\downarrow\theta} \subseteq \theta$, then $R_C^k[\mathcal{P}]$ is equal to $R_{C\downarrow\theta}^k[\mathcal{P}]$ for all $k \geq 0$ (and, thus, $\llbracket \mathcal{P} \rrbracket_{C\downarrow\theta} = \llbracket \mathcal{P} \rrbracket_C$). We proceed by induction on k .

Assume $k = 0$. By definition, $R_C^0[\mathcal{P}]$ and $R_{C\downarrow\theta}^0[\mathcal{P}]$ both equal $\mathcal{I}_C[\mathcal{P}]$. Thus, $R_C^0[\mathcal{P}]$ equals $R_{C\downarrow\theta}^0[\mathcal{P}]$.

Now, assume $k > 0$ and $R_C^{k-1}[\mathcal{P}]$ equals $R_{C\downarrow\theta}^{k-1}[\mathcal{P}]$. Since $\llbracket \mathcal{P} \rrbracket_{C\downarrow\theta} \subseteq \theta$, we must have also $R_C^{k-1}[\mathcal{P}] \subseteq \theta$. Hence,

$$\begin{aligned} R_C^k[\mathcal{P}] &= \mathcal{F}_C[\mathcal{P}](\theta \cap R_C^{k-1}[\mathcal{P}]) \\ &= \mathcal{F}_C[\mathcal{P}](R_{C\downarrow\theta}^{k-1}[\mathcal{P}]) \\ &= \mathcal{F}_C[\mathcal{P}](R_{C\downarrow\theta}^{k-1}[\mathcal{P}]) \\ &= R_{C\downarrow\theta}^k[\mathcal{P}] \end{aligned}$$

□

It is sufficient for \mathcal{A} to have a θ -sound transfer function.

Definition 6. Let $\mathcal{F}_{\mathcal{D}}$ be a transfer function over domain D . $\mathcal{F}_{\mathcal{D}}$ is θ -sound iff for any program \mathcal{P} , $\mathcal{F}_{\mathcal{D}}[\mathcal{P}](D')$ over-approximates $\mathcal{F}_{\mathcal{C}}[\mathcal{P}](C')$ whenever D' over-approximates C' and $C' \subseteq \theta$.

Theorem 3. Let \mathcal{A} be an analysis. If $\mathcal{I}_{\mathcal{A}}$ is sound and $\mathcal{F}_{\mathcal{A}}$ is θ -sound, then \mathcal{A} is θ -sound.

Proof. Take a program \mathcal{P} . We must show that, if c is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\theta}$, it is also in $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})$. By definition, if c is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\theta}$, then it is in some $R_{\mathcal{C}}^k \downarrow_{\theta}[\mathcal{P}]$. Similarly, if c is in $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})$, then c is in some $\gamma_{\mathcal{A}}(R_{\mathcal{A}}^k[\mathcal{P}])$. Thus, it suffices to show that $R_{\mathcal{A}}^k[\mathcal{P}]$ over-approximates $R_{\mathcal{C}}^k \downarrow_{\theta}[\mathcal{P}]$ for all $k \geq 0$. We proceed by induction on k .

Assume $k = 0$. By definition, $R_{\mathcal{A}}^0[\mathcal{P}]$ is $\mathcal{I}_{\mathcal{A}}[\mathcal{P}]$ and $R_{\mathcal{C}}^0 \downarrow_{\theta}[\mathcal{P}]$ is $\mathcal{I}_{\mathcal{C}}[\mathcal{P}]$. Since $\mathcal{I}_{\mathcal{A}}$ is sound, $\mathcal{I}_{\mathcal{A}}[\mathcal{P}]$ over-approximates $\mathcal{I}_{\mathcal{C}}[\mathcal{P}]$. Hence, $R_{\mathcal{A}}^0[\mathcal{P}]$ over-approximates $R_{\mathcal{C}}^0 \downarrow_{\theta}[\mathcal{P}]$.

Now, assume $k > 0$ and $R_{\mathcal{A}}^{k-1}[\mathcal{P}]$ over-approximates $R_{\mathcal{C}}^{k-1} \downarrow_{\theta}[\mathcal{P}]$. By definition,

$$\begin{aligned} R_{\mathcal{A}}^k[\mathcal{P}] &= \mathcal{F}_{\mathcal{A}}[\mathcal{P}](R_{\mathcal{A}}^{k-1}[\mathcal{P}]) \\ R_{\mathcal{C}}^k \downarrow_{\theta}[\mathcal{P}] &= \mathcal{F}_{\mathcal{C}}[\mathcal{P}](\theta \cap R_{\mathcal{C}}^{k-1} \downarrow_{\theta}[\mathcal{P}]) \end{aligned}$$

Since $\mathcal{F}_{\mathcal{A}}$ is θ -sound, $\mathcal{F}_{\mathcal{A}}[\mathcal{P}](A')$ over-approximates $\mathcal{F}_{\mathcal{C}}[\mathcal{P}](C')$ whenever A' over-approximates C' and $C' \subseteq \theta$. Take $R_{\mathcal{A}}^{k-1}[\mathcal{P}]$ for A' and $\theta \cap R_{\mathcal{C}}^{k-1} \downarrow_{\theta}[\mathcal{P}]$ for C' . By assumption, $R_{\mathcal{A}}^{k-1}[\mathcal{P}]$ over-approximates $R_{\mathcal{C}}^{k-1} \downarrow_{\theta}[\mathcal{P}]$. Hence, $R_{\mathcal{A}}^k[\mathcal{P}]$ over-approximates $\theta \cap R_{\mathcal{C}}^{k-1} \downarrow_{\theta}[\mathcal{P}]$. By definition, $\theta \cap R_{\mathcal{C}}^{k-1} \downarrow_{\theta}[\mathcal{P}]$ is a subset of θ . Thus, $R_{\mathcal{A}}^k[\mathcal{P}]$ over-approximates $R_{\mathcal{C}}^k \downarrow_{\theta}[\mathcal{P}]$. \square

Note 1. Theorem 3 generalizes Theorem 1: If $\mathcal{F}_{\mathcal{A}}$ is unconditionally sound, then it is θ -sound for any θ . Hence, if $\mathcal{I}_{\mathcal{A}}$ and $\mathcal{F}_{\mathcal{A}}$ are unconditionally sound, \mathcal{A} is unconditionally sound.

In Section 4, we will show that points-to analysis is SAFEDEREF-sound, where SAFEDEREF is a predicate that captures memory safety.

Parameterized Analysis. Having defined a precise notion of conditional soundness, we now consider how the results of a θ -sound analysis can be used to refine a second analysis. Suppose that \mathcal{A} is an analysis and we have already computed the set of reachable states $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$. We may wish to use the information present in $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$ to refine a second analysis over a different domain B . For example, we could use the reduced product construction [7] to form a new domain over a subset of $D_{\mathcal{A}} \times B$ including only those states (a, b) where a is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$ and the states a and b “agree” (e.g., $\gamma_{\mathcal{A}}(a) \cap \gamma_B(b) \neq \emptyset$).

Traditional methods for combining analyses take a “white box” approach—e.g., Cousot and Cousot [7] assume that the state transformers are available and can be combined in a mechanical way; Lerner et al. [24] assume that analyses can be run in parallel, one step at a time. In contrast, we will assume that any prior analysis is a black box: we have access to its result (in the form of a set of reachable abstract states), its domain (which allows us to interpret the result), and some (possibly conditional) soundness guarantee. This naturally models the use of off-the-shelf program analyses to provide refinement advice.

We will define such a refinement in terms of a parameterized analysis which produces a new, refined analysis from the results of a prior analysis. An *analysis generator* $\tilde{\mathcal{G}}$ is a tuple $(D, E, \tilde{\mathcal{I}}, \tilde{\mathcal{F}})$ where: D and E are domains (the *input* and *output* domains, respectively) and $\tilde{\mathcal{I}}$ and $\tilde{\mathcal{F}}$ are *parameterized interpretations* mapping a set of states $D' \subseteq D$ to semantic interpretations $\tilde{\mathcal{I}}\langle D' \rangle$ and $\tilde{\mathcal{F}}\langle D' \rangle$ over E . We denote by $\tilde{\mathcal{G}}\langle D' \rangle$ the analysis over E defined by the parameterized interpretations on input D' : $\tilde{\mathcal{G}}\langle D' \rangle = (E, \tilde{\mathcal{I}}\langle D' \rangle, \tilde{\mathcal{F}}\langle D' \rangle)$. As one might expect, the soundness of $\tilde{\mathcal{G}}\langle D' \rangle$ depends on the input D' .

Definition 7. An analysis generator $\tilde{\mathcal{G}}$ with input domain D is sound iff for every set of states $D' \subseteq D$, $\tilde{\mathcal{G}}\langle D' \rangle$ is θ -sound with $\theta = \gamma_D(D')$.

Given an analysis generator, it is natural to consider the analysis formed by composing the generator with an analysis over its input domain. If \mathcal{A} is an analysis and $\tilde{\mathcal{G}}$ is an analysis generator with input domain $D_{\mathcal{A}}$ (i.e., the input domain of $\tilde{\mathcal{G}}$ is the underlying domain of \mathcal{A}), the *composed analysis* $\tilde{\mathcal{G}} \circ \mathcal{A}$ is defined by providing the result of \mathcal{A} as a parameter to $\tilde{\mathcal{G}}$ (i.e., $\tilde{\mathcal{G}} \circ \mathcal{A} = \tilde{\mathcal{G}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})$). An important property of the composed analysis is preservation of soundness.

$$n \in \mathbb{Z} \quad \mathbf{x}, \mathbf{y} \in \text{Vars}$$

$$\begin{aligned} L \in \text{Lvals} &::= \mathbf{x} \mid * \mathbf{x} \\ E \in \text{Exprs} &::= L \mid n \mid \mathbf{x} \oplus \mathbf{y} \mid \mathbf{x} \leq \mathbf{y} \mid \&\mathbf{x} \\ S \in \text{Stmts} &::= L := E \mid [E] \end{aligned}$$

Fig. 2. Grammar for a minimal C-like language.

Theorem 4. *If $\tilde{\mathcal{G}}$ is sound and \mathcal{A} is θ -sound, then the composed analysis $\tilde{\mathcal{G}} \circ \mathcal{A}$ is θ -sound.*

Proof. Let B be the output domain of $\tilde{\mathcal{G}}$ and \mathcal{P} a program. We must show that, if c is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\theta}$, it is also in $\gamma_B(\llbracket \mathcal{P} \rrbracket_{\tilde{\mathcal{G}} \circ \mathcal{A}})$, which is equal to $\gamma_B(\llbracket \mathcal{P} \rrbracket_{\tilde{\mathcal{G}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})})$.

By definition, if c is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\theta}$, then it is in some $R_{\mathcal{C}}^k \downarrow_{\theta}[\mathcal{P}]$. Hence, by the definitions of $R_{\mathcal{C}}^k \downarrow_{\theta}$ and $\mathcal{F}_{\mathcal{C}}$, there is some sequence c_0, c_1, \dots, c_k of concrete states from $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\theta}$ such that:

- c_k is equal to c ;
- c_i is in $\mathcal{F}_{\mathcal{C}}[\mathcal{P}](c_{i-1})$, for each $0 < i \leq k$;
- c_0 is in $\mathcal{I}_{\mathcal{C}}[\mathcal{P}]$; and
- each state c_0, c_1, \dots, c_{k-1} is in θ .

Since \mathcal{A} is θ -sound, we know that $\llbracket \mathcal{P} \rrbracket_{\mathcal{A}}$ over-approximates $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\theta}$. Hence, each of the states c_0, c_1, \dots, c_k is in $\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})$ and, thus, each c_i is in $R_{\mathcal{C}}^i \downarrow_{\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})}[\mathcal{P}]$, for $0 \leq i \leq k$. Hence, each of the states c_0, c_1, \dots, c_k is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})}$. In particular, the state c (i.e., c_k) is in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})}$.

Since $\tilde{\mathcal{G}}$ is sound, we know that $\llbracket \mathcal{P} \rrbracket_{\tilde{\mathcal{G}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})}$ over-approximates $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}} \downarrow_{\gamma_{\mathcal{A}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})}$. Hence, c is in $\gamma_B(\llbracket \mathcal{P} \rrbracket_{\tilde{\mathcal{G}}(\llbracket \mathcal{P} \rrbracket_{\mathcal{A}})})$. This shows θ -soundness for $\tilde{\mathcal{G}} \circ \mathcal{A}$. \square

In the remainder of this paper, we will apply Theorem 4 to the problem of verifying memory safety using points-to information. In Section 3, we define the concrete semantics for a little language that captures the pointer semantics of C. In Section 4, we define a memory safety property SAFEDEREF and a set of SAFEDEREF-sound points-to analyses similar to the points-to analyses found in the literature. In Section 5, we define a memory safety analysis parameterized by points-to information. These, together with Theorem 4, allow us to prove the absence of memory safety errors.

Note that since the points-to analysis is SAFEDEREF-sound, we could in theory use it to prove the absence of memory safety errors. However, the points-to domain does not track memory safety information with adequate precision to detect errors without an impractical number of false alarms. The value of Theorem 4 is that it allows for the definition of a specialized memory safety analysis which uses the points-to analysis to increase its precision while remaining SAFEDEREF-sound.

3 Concrete Semantics

To make precise statements about program analyses requires a concrete program semantics. We will define the semantics of the little language presented in Fig. 2. The language eliminates all but those features of C that are essential to the question at hand. The semantics of the language is chosen to model the requirements of ANSI/ISO C [22] without making implementation-specific assumptions. Undefined or implementation-defined behaviors are modeled with explicit nondeterminism. Note that an ANSI/ISO-compliant C compiler is free to implement undefined behaviors in a specific, deterministic manner. By modeling undefined behaviors using non-determinism, the soundness statements made about each analysis apply to *any* standard-compliant compilation strategy.

The most important features of C that we exclude here are fixed-size integer types, narrowing casts, dynamic memory allocation, and functions.³ We also ignore the “strict aliasing” rule [22, §6.5]. Each of these can be handled, at the cost of a higher degree of complexity in our definitions.

The syntactic classes of variables, lvalues, expressions, and statements, are defined in Fig. 2. We use n to represent an integer constant and \mathbf{x} and \mathbf{y} to represent arbitrary variables. We use \oplus to represent an arbitrary binary arithmetic operator and \trianglelefteq to represent a relational operator. Pointer operations include arithmetic, indirection ($*$), and address-of ($\&$). Statements include assignments and tests ($[E]$, where E is an expression).

Variables in our language are viewed as arrays of memory cells. Each cell may hold either an unbounded integer or a pointer value. The only type information present is the allocated size of each variable—the “type system” merely maps variables to their sizes and provides no safety guarantees.

A *program* \mathcal{P} is a tuple $(\mathcal{V}, \Gamma, \mathcal{L}, \mathcal{S}, \tau, en)$, where: $\mathcal{V} \subseteq Vars$ is a finite set of *program variables*; $\Gamma : \mathcal{V} \rightarrow \mathbb{N}$ is a *typing environment* mapping variables to their allocated sizes; \mathcal{L} is a finite set of *program points*; $\mathcal{S} \subseteq Stmts$ is a finite set of *program statements* whose variables are from \mathcal{V} ; $\tau \subseteq \mathcal{L} \times \mathcal{S} \times \mathcal{L}$ is a *transition relation*; and $en \in \mathcal{L}$ is a distinguished *entry point*. In the following, we assume a fixed program $\mathcal{P} = (\mathcal{V}, \Gamma, \mathcal{L}, \mathcal{S}, \tau, en)$.

Example 1. Figure 1(b) gives a fragment of the program representation for the code in Fig. 1(a), corresponding to the function `bad`. We have introduced temporaries \mathbf{t}_1 and \mathbf{t}_2 in order to simplify expression evaluation and compressed multiple statements onto a single transition when they represent a single statement in the source program.

In order to reason about points-to and memory safety analyses, we need a memory model on which to base the concrete semantics. The unit of memory allocation is a *home* in the set \mathbb{H} . Each home h represents a contiguous block of memory cells, e.g., a statically declared array. A *location* $h[i]$ represents the cell at integer *offset* i in home h . The set of locations with homes from \mathbb{H} is denoted \mathbb{L} . The function **size** : $\mathbb{H} \rightarrow \mathbb{N}$ maps a home to its allocated size. When $0 \leq i < \mathbf{size}(h)$, location $h[i]$ is *in bounds*; otherwise it is *out of bounds*. Memory locations contain values from the set $Vals = \mathbb{Z} \cup \mathbb{L}$. A *memory state* is a partial function $m : \mathbb{L} \rightarrow Vals$. The set of all memory states is denoted \mathbb{M} . The set of concrete states C is the set of pairs (p, m) where $p \in \mathcal{L}$ represents the program position and m is a memory state.

An *allocation for* \mathcal{V} is an injective function **home** : $\mathcal{V} \rightarrow \mathbb{H}$ such that **size**(**home**(\mathbf{x})) = $\Gamma(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{V}$. Given such an allocation, the *lvalue* of $\mathbf{x} \in \mathcal{V}$ is **lval**(\mathbf{x}) = **home**(\mathbf{x})[0]. We write $m(\mathbf{x})$ for $m(\mathbf{lval}(\mathbf{x}))$ and $m[\mathbf{x} \mapsto v]$ for $m[\mathbf{lval}(\mathbf{x}) \mapsto v]$, where m is a memory state. We say a location $h[i]$ is *within* a variable \mathbf{x} if $h = \mathbf{home}(\mathbf{x})$ and $h[i]$ is in bounds.

Figure 3 defines the concrete interpretations \mathcal{E} and **post** of, respectively, expressions and statements. Note that both \mathcal{E} and **post** result in *sets* of, respectively, values and concrete states—the set-based semantics is needed as undefined operations may have a nondeterministic result. \mathcal{E} returns the distinguished value \perp in the case where an expression is not just ill-defined, but erroneous (e.g., reading an out-of-bounds memory location)—in this case the next state can have *any* memory state at *any* program point.

The operator $\tilde{\oplus}$ is used to denote the semantic counterpart to the syntactic operator \oplus . The definition of $\tilde{\oplus}$ is as usual for integer values. If an integer j is added to (alt. subtracted from) a location $h[i]$, where both $0 \leq i \leq \mathbf{size}(h)$ and $0 \leq i + j \leq \mathbf{size}(h)$ (alt. $0 \leq i - j \leq \mathbf{size}(h)$), then the result is $h[i + j]$ (alt. $h[i - j]$). If a location $h[j]$ is subtracted from a location $h'[i]$, where $h = h'$ and $0 \leq i, j \leq \mathbf{size}(h)$, the result is $i - j$. In all other cases, the result is undefined.

Note that arithmetic on pointer values is only defined for locations within (or one location beyond) a single home. E.g., adding an offset to a location sufficient to create an out-of-bounds location does not make the value point to a new home; subtracting locations from two different homes does not indicate the “distance” between the homes.

The operator $\tilde{\trianglelefteq}$ is used to denote the semantic counterpart of a relational operator \trianglelefteq . The definition of $\tilde{\trianglelefteq}$ is as usual for integer values. If two locations $h[i]$ and $h[j]$ have the same home, then $\tilde{\trianglelefteq}$ is equal to the

³ The omission of dynamic allocation in the discussion of points-to analysis and memory safety may seem an oversimplification. However, it is not essential to our purpose here. Points-to analyses typically handle dynamic allocation by treating each allocation site as if it were the static declaration of a global array of unknown size.

$$\begin{aligned}
\mathcal{E}(m, n) &= \{n\} \\
\mathcal{E}(m, \mathbf{x}) &= \begin{cases} \mathbb{Z}, & \text{if } m(\mathbf{x}) \text{ is undefined} \\ \{m(\mathbf{x})\}, & \text{otherwise} \end{cases} \\
\mathcal{E}(m, * \mathbf{x}) &= \begin{cases} \perp, & \text{if } m(\mathbf{x}) \text{ is undefined, not a location, or out of bounds} \\ \mathbb{Z} & \text{if } m(m(\mathbf{x})) \text{ is undefined} \\ \{m(m(\mathbf{x}))\}, & \text{otherwise} \end{cases} \\
\mathcal{E}(m, \mathbf{x} \oplus \mathbf{y}) &= \begin{cases} \mathbb{Z}, & \text{if } m(\mathbf{x}) \tilde{\oplus} m(\mathbf{y}) \text{ is undefined} \\ \{m(\mathbf{x}) \tilde{\oplus} m(\mathbf{y})\}, & \text{otherwise} \end{cases} \\
\mathcal{E}(m, \&\mathbf{x}) &= \{\mathbf{lval}(\mathbf{x})\} \\
\mathcal{E}(m, \mathbf{x} \trianglelefteq \mathbf{y}) &= \begin{cases} \{0, 1\}, & \text{if } m(\mathbf{x}) \tilde{\trianglelefteq} m(\mathbf{y}) \text{ is undefined} \\ \{m(\mathbf{x}) \tilde{\trianglelefteq} m(\mathbf{y})\}, & \text{otherwise} \end{cases} \\
\mathbf{post}(m, p, \mathbf{x} := E) &= \begin{cases} \mathcal{L} \times \mathbb{M}, & \text{if } \mathcal{E}(m, E) = \perp \\ \{(p, m[\mathbf{x} \mapsto v]) \mid v \in \mathcal{E}(m, E)\}, & \text{otherwise} \end{cases} \\
\mathbf{post}(m, p, * \mathbf{x} := E) &= \begin{cases} \mathcal{L} \times \mathbb{M}, & \text{if } m(\mathbf{x}) \text{ is undefined, not a location, or out of bounds; or if } \mathcal{E}(m, E) = \perp \\ \{(p, m[m(\mathbf{x}) \mapsto v]) \mid v \in \mathcal{E}(m, E)\}, & \text{otherwise.} \end{cases} \\
\mathbf{post}(m, p, [E]) &= \begin{cases} \emptyset, & \text{if } \mathcal{E}(m, E) = \{0\} \\ \{(p, m)\}, & \text{otherwise.} \end{cases}
\end{aligned}$$

Fig. 3. The concrete interpretation.

integer comparison $i \tilde{\trianglelefteq} j$. The value of $\tilde{\trianglelefteq}$ is otherwise undefined, with the following exception: equality (resp. disequality) on two in-bounds locations with different homes or between an in-bounds location and 0 (the *null pointer constant*) evaluates to FALSE (resp. TRUE).

We now define the concrete interpretation of a program.

Definition 8. The concrete semantics $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}}$ of a program $\mathcal{P} = (\mathcal{V}, \Gamma, \mathcal{L}, \mathcal{S}, \tau, en)$ is defined by the analysis $\mathcal{C} = (D_{\mathcal{C}}, \mathcal{I}_{\mathcal{C}}, \mathcal{F}_{\mathcal{C}})$, where

$$\begin{aligned}
\mathcal{I}_{\mathcal{C}}[\mathcal{P}] &= \{(en, m) \mid \forall l \in \mathbb{L}. m(l) \text{ is not a location}\} \\
\mathcal{F}_{\mathcal{C}}[\mathcal{P}](p, m) &= \bigcup_{(p, S, p') \in \tau} \mathbf{post}(m, p', S)
\end{aligned}$$

If (p, S, p') is in τ and $c' \in \mathbf{post}(m, p', S)$, we say c' is an S -successor of (p, m) .

Example 2. Figure 4(a) gives a subset of the reachable concrete states of the program in Fig. 1(b). At ℓ_0 , \mathbf{p} is $\mathbf{A}[0]$ (the base address of the array \mathbf{A}), \mathbf{x} is 1, and \mathbf{y} is 0. At ℓ_1 , due to the assignment to out-of-bounds location $\mathbf{A}[4]$, the next state is undefined: *every* program point is reachable with *any* memory state.

Note 2. The following modifications would allow us to model other aspects of the C programming language. To model dynamic memory allocation, we could add a component to the memory state to map locations to their allocation status. To model fixed-size integers and narrowing casts, we would need to define scalar values that span multiple locations (i.e., bytes), which require the modeling of alignment, padding, and byte order. To model functions, we could add an explicit representation of the call stack and variable scope. Handling strict aliasing would require more detailed type information, and would introduce more undefined behaviors in the case where a location is accessed through an illegal “type pun.” This would correspondingly weaken the soundness results below.

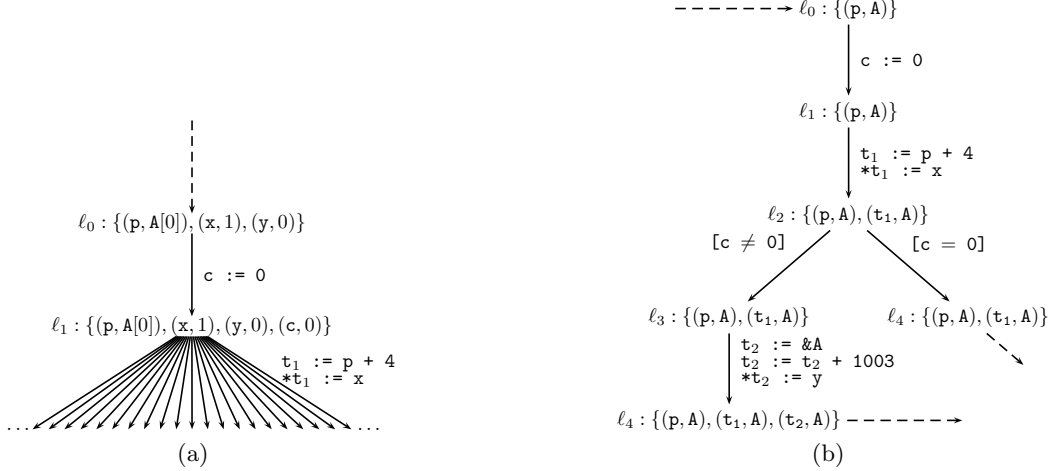


Fig. 4. Concrete and points-to semantics for the program in Fig. 1(b)

4 Pointer Analysis

The goal of pointer analysis is to compute an over-approximate points-to set for each variable in the program, i.e., the set of homes “into” which a variable may point in some reachable state.

A *points-to state* is a relation between variables. We denote the set of points-to states by Pts . When it is convenient, we treat a points-to state also as a relation between variables and memory locations: for points-to state pts , variables x, y , and location $h[i]$, we say $(x, h[i])$ is in pts when (x, y) is in pts and $h[i]$ is within y (i.e., $h[i]$ is in bounds and $h = \mathbf{home}(y)$). We write $pts(x)$ for the *points-to set* of the variable x in pts , i.e., the set of variables y (alt. locations l) such that (x, y) (alt. (x, l)) is in pts .

The concretization function γ_{Pts} takes a points-to state to the set of concrete states where at *most* its points-to relationships hold. Say that variable x *points to* y in memory state m if there exist locations l_1, l_2 such that l_1 is within x , l_2 is within y , and $m(l_1) = l_2$. Then m is in $\gamma_{Pts}(pts)$ iff for all x, y such that x points to y in m , the pair (x, y) is in pts . Note that there may be other pairs in pts as well—the points-to relation is over-approximate. Note also that *only in-bounds location values must agree with the points-to state*; out-of-bounds locations are unconstrained.

Note 3. Since the set of variables is known *a priori* and fixed, we will assume for convenience that the function \mathbf{home} is a bijection. To extend our results to programs with dynamic allocation, we would have to introduce an inverse map from homes to abstract variables, such that every allocated block of memory has some representative in the domain of the points-to relation. For example, we may introduce a set of variables \mathbf{malloc}_p , for p in \mathcal{L} , such that any memory location allocated by a call to \mathbf{malloc} at program point p is within \mathbf{malloc}_p . In any case, we would like to maintain the property: if $m(x)$ is an in-bounds location, then $m(x)$ is within some (perhaps abstract) variable y .

Figure 5 defines the interpretations \mathcal{E}_{Pts} and \mathbf{post}_{Pts} for a selection of, respectively, expressions and statements in the points-to domain. The interpretations are chosen to match those used by common points-to analyses. A key feature is the treatment of the indirection operator $*$, which assumes that its argument is within bounds. Without this assumption, the interpretation would have to use the “top” points-to state (i.e., all pairs of variables) for the result of any indirect assignment.

We lift Pts to the set $\mathcal{L} \times Pts$ in the natural way.

$$\begin{aligned}
\mathcal{E}_{Pts}(pts, n) &= \emptyset \\
\mathcal{E}_{Pts}(pts, \mathbf{x}) &= pts(\mathbf{x}) \\
\mathcal{E}_{Pts}(pts, *x) &= \{z \in \mathcal{V} \mid \exists y \in \mathcal{V} : pts(\mathbf{x}, y) \wedge pts(y, z)\} \\
\mathcal{E}_{Pts}(pts, \mathbf{x} \oplus \mathbf{y}) &= pts(\mathbf{x}) \cup pts(\mathbf{y}) \\
\mathcal{E}_{Pts}(pts, \&\mathbf{x}) &= \{\mathbf{x}\} \\
\mathcal{E}_{Pts}(pts, \mathbf{x} \trianglelefteq \mathbf{y}) &= \emptyset
\end{aligned}$$

$$\begin{aligned}
\mathbf{post}_{Pts}(pts, \mathbf{x} := E) &= pts \cup \{(\mathbf{x}, \mathbf{y}) \mid \mathbf{y} \in \mathcal{E}_{Pts}(pts, E)\} \\
\mathbf{post}_{Pts}(pts, *x := E) &= \bigcup_{(x, y) \in pts} \mathbf{post}_{Pts}(pts, y := E) \\
\mathbf{post}_{Pts}(pts, [\mathbf{x} \trianglelefteq \mathbf{y}]) &= pts
\end{aligned}$$

Fig. 5. Abstract interpretation over points-to states.

Definition 9. A flow- and path-sensitive points-to analysis **Pts** is given by the tuple $(Pts, \mathcal{I}_{Pts}, \mathcal{F}_{Pts})$, where

$$\begin{aligned}
\mathcal{I}_{Pts}[\mathcal{P}] &= \{(en, \emptyset)\} \\
\mathcal{F}_{Pts}[\mathcal{P}](p, pts) &= \bigcup_{(p, S, p') \in \tau} (p', \mathbf{post}_{Pts}(pts, S))
\end{aligned}$$

Example 3. Figure 4(b) shows a subset of the reachable points-to states for the program in Fig. 1(b). At ℓ_0 , \mathbf{p} points to **A**. The transition from ℓ_1 to ℓ_2 causes \mathbf{t}_1 to point to **A** as well. The presence of an out-of-bounds array access has no effect on the points-to state: the analysis assumes that evaluating $*\mathbf{t}_1$ is safe.

Definition 10. Let SAFEDEREF be the predicate that holds in a concrete state (p, m) if, for every transition (p, S, p') in τ where S includes an expression of the form $*x$, $m(x)$ is an in-bounds location.

To show that **Pts** is SAFEDEREF-sound, we must first show that the function \mathcal{E}_{Pts} over-approximates the locations given by the function \mathcal{E} on concrete states in SAFEDEREF. We formalize this with the following Lemma.

Lemma 1. For concrete state $c = (p, m)$, points-to state pts , variable \mathbf{x} , location l , and expression E : if (1) SAFEDEREF(c) holds, (2) m is in $\gamma_{Pts}(pts)$, (3) there exists an edge (p, S, p') in τ such that E appears in S , (4) l is in $\mathcal{E}(m, E)$, and (5) l is within \mathbf{x} , then \mathbf{x} is in $\mathcal{E}_{Pts}(pts, E)$.

Proof. We proceed by cases over E .

- $E = n$. No l can satisfy (4). The property holds trivially.
- $E = \mathbf{y}$. If $m(\mathbf{y})$ is undefined, then $\mathcal{E}(m, E) = \mathbb{Z}$ and no l can satisfy (4). Assume $m(\mathbf{y})$ is defined. Then, by definition, the only value in $\mathcal{E}(m, E)$ is $m(\mathbf{y})$. Assume $m(\mathbf{y})$ is a location within \mathbf{x} . Since m is in $\gamma_{Pts}(pts)$, (\mathbf{y}, \mathbf{x}) must be in pts . Hence, \mathbf{x} is in $\mathcal{E}(pts, E)$.
- $E = *\mathbf{y}$. Since c is a SAFEDEREF state, $m(\mathbf{y})$ must be an in-bounds location l' . If $m(l')$ is undefined, then $\mathcal{E}(m, E) = \mathbb{Z}$ and no l can satisfy (4). Assume $m(l')$ is defined. Then, by definition, the only value in $\mathcal{E}(m, E)$ is $m(l')$. Assume $m(l')$ is a location within \mathbf{x} . Take \mathbf{z} such that l' is within \mathbf{z} : such a \mathbf{z} is guaranteed to exist if **home** is a bijection (see Note 3). Since m is in $\gamma_{Pts}(pts)$, (\mathbf{y}, \mathbf{z}) and (\mathbf{z}, \mathbf{x}) must be in pts . Hence, \mathbf{x} is in $\mathcal{E}_{Pts}(pts, E)$.
- $E = \mathbf{y} \oplus \mathbf{z}$. The only case where $\mathcal{E}(m, E)$ contains a location value is when \oplus is $+$, one of the operands is location, the other is an integer, and the result is well-defined. Assume, wlog, that $m(\mathbf{y}) = h[i]$ (with $0 \leq i \leq \mathbf{size}(h)$) and $m(\mathbf{z}) = j$ (with $0 \leq i + j \leq \mathbf{size}(h)$). By definition, the only value in $\mathcal{E}(m, E)$ is $h[i + j]$. Assume $h[i + j]$ is within \mathbf{x} . Then $h[i]$ is within \mathbf{x} . Since m is in $\gamma_{Pts}(pts)$, \mathbf{x} must be in $pts(\mathbf{y})$. Hence, \mathbf{x} is in $\mathcal{E}_{Pts}(pts, E)$.

- $E = \&x$. By definition, the only value in $\mathcal{E}(m, E)$ is $\mathbf{lval}(x)$, which is within \mathbf{x} . By definition, \mathbf{x} is in $\mathcal{E}_{Pts}(pts, E)$.
- $E = \mathbf{y} \trianglelefteq \mathbf{z}$. No l can satisfy (4). The property holds trivially.

□

Theorem 5. *The points-to analysis \mathbf{Pts} is SAFEDEREF-sound.*

Proof. By Theorem 3, it suffices to show that \mathcal{I}_{Pts} is sound and \mathcal{F}_{Pts} is SAFEDEREF-sound.

Take a program $\mathcal{P} = (\mathcal{V}, \Gamma, \mathcal{L}, \mathcal{S}, \tau, en)$. We must show: (1) $\mathcal{I}_{Pts}[\mathcal{P}]$ over-approximates $\mathcal{I}_C[\mathcal{P}]$ and (2) $\mathcal{F}_{Pts}[\mathcal{P}](p, pts)$ over-approximates $\mathcal{F}_C[\mathcal{P}](p, m)$ whenever (p, m) is a SAFEDEREF state and m is in $\gamma_{Pts}(pts)$.

- (1) Let (en, m) be an initial concrete state. By the definition of \mathcal{I}_C , $m(l)$ is not a location for any $l \in \mathbb{L}$. By the definition of \mathcal{I}_{Pts} , the only initial points-to state is (en, \emptyset) . By the definition of γ_{Pts} , m is in $\gamma_{Pts}(\emptyset)$. Hence, (en, m) is in $\gamma_{Pts}(\mathcal{I}_{Pts}[\mathcal{P}])$. This shows soundness for \mathcal{I}_{Pts} .
- (2) Take concrete states $c = (p, m)$, c' , points-to state pts , and statement S such that m is in $\gamma_{Pts}(pts)$, SAFEDEREF(c) holds, and c' is an S -successor of c . It suffices to show that c' is in $\gamma_{Pts}(\mathcal{F}_{Pts}[\mathcal{P}](p, pts))$. We proceed by cases on S :

- $S = \mathbf{x} := E$. Follows from Lemma 1.
- $S = *x := E$. It suffices to show that, if $m(x)$ is within \mathbf{y} , then $pts(x, \mathbf{y})$. This follows immediately from the definition of γ_{Pts} .
- $S = [E]$. By definition of \mathbf{post} and \mathbf{post}_{Pts} , $c = c'$ and $\mathbf{post}_{Pts}(pts, S) = pts$. Thus c' is in $\gamma_{Pts}(\mathbf{post}_{Pts}(pts, S))$.

This shows SAFEDEREF-soundness for \mathcal{F}_{Pts} .

Hence, \mathbf{Pts} is SAFEDEREF-sound. □

We can extract more traditional flow-sensitive, global, and flow-insensitive pointer analyses from $\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}}$ as follows.

- A *flow-sensitive, program-point-sensitive (path-insensitive) analysis* is derived by assigning to each program point p the least points-to state (by subset inclusion) pts^\sharp such that, if (p, pts) is in $\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}}$, then $pts \subseteq pts^\sharp$.
- A *flow-sensitive, global (program-point-insensitive) analysis* is derived by assigning to every program point the least points-to state (by subset inclusion) pts^\sharp such that, if (p, pts) is in $\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}}$ for *any* program point p , then $pts \subseteq pts^\sharp$.
- A *flow-insensitive analysis* is derived by replacing τ in Definition 9 with the relation τ^\sharp , where the edge (p, S, q) is in τ^\sharp whenever some edge (t, S, u) is in τ , for *any* program points t and u . Intuitively, if a statement occurs anywhere in the program, then it may occur between any two program points—the interpretation ignores the control-flow structure of the program.
- *Flow-insensitive, program-point-sensitive and flow-insensitive, global* combinations can be defined as above, substituting the flow-insensitive semantics for $\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}}$.

Theorem 6. *Each of the flow-, path-, and program-point-sensitive and insensitive variations of the points-to analysis is SAFEDEREF-sound.*

Proof. Take a program \mathcal{P} and a points-to analysis \mathcal{Q} from among the above. It is clear that, for any state (p, pts) in $\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}}$, there is a state (p, pts') in $\llbracket \mathcal{P} \rrbracket_{\mathcal{Q}}$ such that $pts \subseteq pts'$. Since γ_{Pts} is monotonic, any concrete state in $\gamma_{Pts}(\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}})$ is also in $\gamma_{Pts}(\llbracket \mathcal{P} \rrbracket_{\mathcal{Q}})$. By Theorem 5, $\gamma_{Pts}(\llbracket \mathcal{P} \rrbracket_{\mathbf{Pts}})$ over-approximates $\llbracket \mathcal{P} \rrbracket_C \downarrow_{\text{SAFEDEREF}}$. Hence, the $\llbracket \mathcal{P} \rrbracket_{\mathcal{Q}}$ over-approximates $\llbracket \mathcal{P} \rrbracket_C \downarrow_{\text{SAFEDEREF}}$: \mathcal{Q} is SAFEDEREF-sound. □

Note 4. The flow-sensitive, program-point-sensitive analysis yields a points-to relation similar to that of Emami et al. [16]. The flow-insensitive, global analysis procedure yields a points-to relation similar to that of Andersen [3]. The Steensgaard [27] and Das [12] relations add additional approximation to the global relation. We claim (but do not prove formally here) that these procedures approximate $\llbracket \mathcal{P} \rrbracket_{\text{pts}}$ and, thus, are *at least* SAFEDEREF-sound.

By the definition of conditional soundness, it is possible some condition θ weaker than SAFEDEREF exists such that some or all of the above analyses are θ -sound. It is our belief that this is not the case: no realistic points-to analysis is θ -sound for any θ weaker than SAFEDEREF. A proof of this proposition is beyond the scope of this paper.

In summary, we have shown that a set of points-to analyses which share the assumptions of widely used analyses from the literature are sound for all memory-safe executions. This claim is both stronger and more precise than any correctness claims the authors have encountered: our points-to analyses (and, by extension, those cited above) compute a relation which is conservative not only for “well-behaved” (i.e., memory-safe) programs, but for all well-behaved *executions*, even the well-behaved executions of ill-behaved programs

We have shown that, if we can prove the absence of non-SAFEDEREF states in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}}$, the points-to analyses we have defined above will be sound. It remains to describe an analysis parameterized by points-to information which can perform a precise memory safety analysis.

5 Checking Memory Safety

We wish to define an analysis procedure that will soundly prove the absence of non-SAFEDEREF states in the concrete program. Note that the only attributes of a location value that are relevant to the property SAFEDEREF are its offset and the size of its home; if we can precisely track these attributes, we can ignore the home component of a location (i.e., which variable it is within) so long as we have access to over-approximate points-to information.

Note 5. In our description of the analysis, we will omit the merging, widening, and covering operations necessary to make the reachability computation tractable. In our implementation of a memory safety analysis in Orion, we constrain integer values and pointer offsets using a relational abstract domain (e.g., convex polyhedra [10]) and use merging and widening to efficiently over-approximate the semantics given below.

Our analysis will track *abstract values* from the set \widehat{Vals} . An abstract value is either an integer or an *abstract location*, a pair (i, n) representing a location at offset i in a home of size n . Each abstract value \hat{v} represents a set of concrete values, according to the abstraction function $\alpha : Vals \rightarrow \widehat{Vals}$. For integer values, α is the identity (i.e., $\alpha(n) = n$); for concrete location values, α preserves the offset and size (i.e., $\alpha(h[i]) = (i, \text{size}(h))$). An abstract location (i, n) is *in bounds* if it represents only in bounds concrete locations (i.e., $0 \leq i < n$); otherwise it is *out of bounds*. An *abstract memory state* is a partial function $b : \mathbb{L} \rightarrow \widehat{Vals}$. We denote by B the set of abstract memory states.

The concretization function $\gamma_B : B \rightarrow 2^{\mathcal{C}}$ takes an abstract memory state b to the set of concrete memories abstracted by b . A concrete memory m is in $\gamma_B(b)$ iff for all l either $m(l)$ and $b(l)$ are both undefined or $\alpha(m(l)) = b(l)$.

Figure 6 defines the interpretations \mathcal{E}_B and post_B for, respectively, expressions and statements with respect to B . Note that the interpretations rely on points-to information. In the limiting case, where no points-to information is available (i.e., the points-to relation includes all pairs), the expression $\star \mathbf{x}$ can take the value of any location abstracted by $b(\mathbf{x})$. As in the concrete interpretation, \mathcal{E}_B returns the value \perp in the case where expression evaluation is (potentially) erroneous.

The operator $\hat{\oplus}$ is used to denote the abstract counterpart to the syntactic operator \oplus . The definition of $\hat{\oplus}$ is as usual for integer values. If an integer j is added to (resp. subtracted from) an abstract location (i, m) , where both $0 \leq i \leq m$ and $0 \leq i + j \leq m$ (resp. $0 \leq i - j \leq m$), the result is $(i + j, m)$ (resp. $(i - j, m)$). In all other cases, the result is undefined.

$$\begin{aligned}
\mathcal{E}_B(b, pts, n) &= \{n\} \\
\mathcal{E}_B(b, pts, \mathbf{x}) &= \begin{cases} \mathbb{Z}, & \text{if } b(\mathbf{x}) \text{ is undefined} \\ \{b(\mathbf{x})\}, & \text{otherwise} \end{cases} \\
\mathcal{E}_B(b, pts, \star \mathbf{x}) &= \begin{cases} \perp, & \text{if } b(\mathbf{x}) \text{ is undefined, not a location, or out of bounds} \\ \widehat{Vals}, & \text{if } b(l) \text{ is undefined for some } l \text{ in } pts(\mathbf{x}), \text{ where } \alpha(l) = b(\mathbf{x}) \\ \{b(l) \mid pts(\mathbf{x}, l), \alpha(l) = b(\mathbf{x})\}, & \text{otherwise} \end{cases} \\
\mathcal{E}_B(b, pts, \mathbf{x} \oplus \mathbf{y}) &= \begin{cases} \mathbb{Z}, & \text{if } b(\mathbf{x}) \widehat{\oplus} b(\mathbf{y}) \text{ is undefined} \\ \{b(\mathbf{x}) \widehat{\oplus} b(\mathbf{y})\}, & \text{otherwise} \end{cases} \\
\mathcal{E}_B(b, pts, \&\mathbf{x}) &= \{(0, \mathbf{size}(\mathbf{home}(\mathbf{x})))\} \\
\mathcal{E}_B(b, pts, \mathbf{x} \leq \mathbf{y}) &= \begin{cases} \{0, 1\}, & \text{if } b(\mathbf{x}) \widehat{\leq}_{pts(\mathbf{x}) \cap pts(\mathbf{y})} b(\mathbf{y}) \text{ is undefined} \\ b(\mathbf{x}) \widehat{\leq}_{pts(\mathbf{x}) \cap pts(\mathbf{y})} b(\mathbf{y}), & \text{otherwise} \end{cases} \\
\mathbf{post}_B(b, pts, p, \mathbf{x} := E) &= \begin{cases} \mathcal{L} \times B, & \text{if } \mathcal{E}_B(b, pts, E) = \perp \\ \{(p, b[\mathbf{x} \mapsto \hat{v}]) \mid \hat{v} \in \mathcal{E}_B(b, pts, E)\}, & \text{otherwise} \end{cases} \\
\mathbf{post}_B(b, pts, p, \star \mathbf{x} := E) &= \begin{cases} \mathcal{L} \times B, & \text{if } b(\mathbf{x}) \text{ is undefined, not a location, or out of bounds; or if} \\ & \mathcal{E}_B(b, pts, E) = \perp \\ \{(p, b[l \mapsto \hat{v}]) \mid pts(\mathbf{x}, l), \alpha(l) = b(\mathbf{x}), \hat{v} \in \mathcal{E}_B(b, pts, E)\}, & \\ & \text{otherwise} \end{cases} \\
\mathbf{post}_B(b, pts, p, [E]) &= \begin{cases} \emptyset, & \text{if } \mathcal{E}_B(b, pts, E) = \{0\} \\ \{(p, b)\}, & \text{otherwise.} \end{cases}
\end{aligned}$$

Fig. 6. Abstract interpretation over B .

The operator $\widehat{\leq}_A$ is used to denote the abstract counterpart to the syntactic operator \leq , parameterized by a points-to set $A \subseteq \mathcal{V}$. The definition of $\widehat{\leq}_A$ is as usual for integer values. If two in-bounds abstract location values (i, m) and (j, n) are compared for equality (resp. disequality) and either $i \neq j$, $m \neq n$, or $A = \emptyset$, then the result is 0 (resp. 1). In all other cases, $\widehat{\leq}_A$ is undefined.

We lift B to the domain $\mathcal{L} \times B$ in the natural way.

Definition 11. *The analysis generator $\widetilde{\mathcal{B}}$ maps a set of points-to states Q to the memory safety analysis $\widetilde{\mathcal{B}}\langle Q \rangle$ defined by the parameterized interpretations*

$$\begin{aligned}
\widetilde{\mathcal{I}}_B\langle Q \rangle[\mathcal{P}] &= \{(en, b) \mid \forall l \in \mathbb{L} : b(l) \text{ is undefined}\} \\
\widetilde{\mathcal{F}}_B\langle Q \rangle[\mathcal{P}](p, b) &= \bigcup_{(p, S, p') \in \tau} \bigcup_{(p, pts) \in Q} \mathbf{post}_B(b, pts, p', S)
\end{aligned}$$

Lemma 2. *For concrete memory state m , points-to state pts , abstract memory state b , value v , and expression E : if (1) m is in $\gamma_{pts}(pts) \cap \gamma_B(b)$ and (2) v is in $\mathcal{E}(m, E)$, then $\alpha(v) = \hat{v}$ for some \hat{v} in $\mathcal{E}_B(b, pts, E)$.*

Proof. We proceed by cases on E .

- $E = n$. By definition, the only value in $\mathcal{E}(m, E)$ or $\mathcal{E}_B(b, pts, E)$ is n and $\alpha(n) = n$.
- $E = \mathbf{x}$. Assume $b(\mathbf{x})$ is undefined. Since m is in $\gamma_B(b)$, $m(\mathbf{x})$ must also be undefined. By definition, $\mathcal{E}(m, E) = \mathcal{E}_B(b, pts, E) = \mathbb{Z}$.
Now, assume $b(\mathbf{x})$ is defined. Since m is in $\gamma_B(b)$, $m(\mathbf{x})$ must also be defined and $\alpha(m(\mathbf{x})) = b(\mathbf{x})$. By definition, the only value in $\mathcal{E}(m, E)$ is $m(\mathbf{x})$ and the only value in $\mathcal{E}_{pts}(b, pts, E)$ is $b(\mathbf{x})$.

- $E = *x$. Assume, wlog, that $m(x)$ is an in-bounds location l . Since m is in $\gamma_B(b)$, $b(x) = \alpha(l)$. Since m is in $\gamma_{pts}(pts)$, l must be in $pts(x)$.

If $b(l')$ is undefined for some l' in $pts(x)$ such that $\alpha(l') = \alpha(l)$, then $\mathcal{E}_B(b, pts, E) = \widehat{Vals}$ and the property holds trivially. Assume $b(l')$ is defined for all such l' . In particular, $b(l)$ is defined. Since m is in $\gamma_B(b)$, $m(l)$ must also be defined and $\alpha(m(l)) = b(l)$. By definition, the only value in $\mathcal{E}(m, E)$ is $m(l)$ and, since l is in $pts(x)$, $b(l)$ is in $\mathcal{E}_B(b, pts, E)$.

- $E = x \oplus y$. If $b(x) \hat{\oplus} b(y)$ is undefined, then $\mathcal{E}_{pts}(b, pts, E) = \mathbb{Z}$. It suffices to show that there are no location values in $\mathcal{E}(m, E)$. The only case where $\mathcal{E}(m, E)$ contains a location value is when \oplus is $+$, one of the operands is a location value, the other is an integer, and the result is well-defined. Assume, wlog, that $m(x)$ is an in-bounds location $h[i]$ and $m(y)$ is an integer j , with $0 \leq i + j \leq \mathbf{size}(h)$. By definition, the only value in $\mathcal{E}(m, E)$ is $h[i] \hat{\oplus} j = h[i + j]$. Since m is in $\gamma_B(b)$, $\alpha(m(x)) = b(x)$ and $\alpha(m(y)) = b(y)$. Hence, $b(x) = (i, \mathbf{size}(h))$ and $b(y) = j$. But $(i, \mathbf{size}(h)) \hat{\oplus} j$ is well-defined, which contradicts the assumption that $b(x) \hat{\oplus} b(y)$ is undefined. Thus, there can be no location values in $\mathcal{E}(m, E)$.

Assume $b(x) \hat{\oplus} b(y)$ is well-defined and both $b(x)$ and $b(y)$ are integers, say i and j . Since m is in $\gamma_B(b)$, $m(x)$ must be i and $m(y)$ must be j . By definition, the only values in $\mathcal{E}(m, E)$ and $\mathcal{E}_B(b, pts, E)$, respectively, are $i \tilde{\oplus} j$ and $i \hat{\oplus} j$. Since $\tilde{\oplus}$ and $\hat{\oplus}$ are defined in the same way for integer operands, $\alpha(i \tilde{\oplus} j) = i \hat{\oplus} j$.

Assume $b(x) \hat{\oplus} b(y)$ is well-defined, one of $b(x)$, $b(y)$ is an abstract location (i, n) (with $0 \leq i \leq n$), the other an integer j (with $0 \leq i + j \leq n$), and \oplus is $+$. Since m is in $\gamma_B(b)$, $m(x)$ must be an location $h[i]$ with $\mathbf{size}(h) = n$ and $m(y)$ must be j . By definition, the only values in $\mathcal{E}(m, E)$ and $\mathcal{E}_B(b, pts, E)$, respectively, are $h[i + j]$ and $(i + j, n)$, and $\alpha(h[i + j]) = (i + j, n)$.

The case where \oplus is $-$ and $b(x)$, $b(y)$ are abstract locations is similar.

- $E = \&x$. By definition, the only values in $\mathcal{E}(m, E)$ and $\mathcal{E}_B(b, pts, E)$, respectively, are $\mathbf{lval}(x)$ and $(0, \mathbf{size}(\mathbf{home}(x)))$. By definition, $\alpha(\mathbf{lval}(x)) = \alpha(\mathbf{home}(x)[0]) = (0, \mathbf{size}(\mathbf{home}(x)))$.
- $E = x \leq y$. If $b(x) \hat{\leq}_{pts(x) \cap pts(y)} b(y)$ is undefined, then the property holds trivially.

Assume $b(x) \hat{\leq}_{pts(x) \cap pts(y)} b(y)$ is well-defined, and $b(x)$, $b(y)$ are both integers, say i and j . Since m is in $\gamma_B(b)$, $m(x)$ must be i and $m(y)$ must be j . By definition, the only values in $\mathcal{E}(m, E)$ and $\mathcal{E}_B(b, pts, E)$, respectively, are $i \hat{\leq} j$ and $i \leq j$. Since $\hat{\leq}$ and \leq are defined in the same way for integer operands, $\alpha(i \hat{\leq} j) = i \leq j$.

Assume $b(x) \hat{\leq}_{pts(x) \cap pts(y)} b(y)$ is well-defined, $b(x)$ and $b(y)$ are in-bounds abstract locations, say (i, n) and (j, r) , and \leq is $==$. Since m is in $\gamma_B(b)$, we have $m(x) = h[i]$ (with $\mathbf{size}(h) = n$) and $m(y) = h'[j]$, (with $\mathbf{size}(h') = r$). By definition, the only value in $\mathcal{E}_B(b, pts, E)$ is 0 and: $i \neq j$, $n \neq r$, or $pts(x) \cap pts(y)$ is empty. We must show that the only value in $\mathcal{E}(m, E)$ is 0: this will be the case when $h[i] \neq h'[j]$. If $i \neq j$, this is immediate. If $n \neq r$, then $\mathbf{size}(h) \neq \mathbf{size}(h')$ and, thus, $h \neq h'$. If $pts(x) \cap pts(y)$ is empty, then, since m is in $\gamma_{pts}(pts)$, $h \neq h'$.

The case when $\hat{\leq}$ is $!=$ and $b(x)$, $b(y)$ are abstract locations is similar.

□

Lemma 3. $\tilde{\mathcal{I}}_B(Q)$ is sound for every points-to state Q .

Proof. Let \mathcal{P} be a program. We must show $\tilde{\mathcal{I}}_B(Q)[\mathcal{P}]$ over-approximates $\mathcal{I}_C[\mathcal{P}]$. Take (p, b) in $\tilde{\mathcal{I}}_B(Q)[\mathcal{P}]$. By definition, b is undefined at all locations. Hence, by the definition of γ_B , every initial concrete state is in $\gamma_B(p, b)$. □

Lemma 4. $\tilde{\mathcal{F}}_B(Q)$ is $\gamma_{pts}(Q)$ -sound for every points-to state Q .

Proof. Let \mathcal{P} be a program. We must show $\tilde{\mathcal{F}}_B(Q)[\mathcal{P}](p, b)$ over-approximates $\mathcal{F}_C[\mathcal{P}](p, m)$ whenever (p, m) is in $\gamma_{pts}(Q)$ and m is in $\gamma_B(b)$.

Take concrete states $c = (p, m)$, c' , abstract state $a = (p, b)$, points-to state (p, pts) , and statement S such that c is in both $\gamma_B(b)$ and $\gamma_{pts}(pts)$ and c' is an S -successor of c . It suffices to show that there is some b' in $\mathbf{post}_B(b, pts, S)$ such that c' is in $\gamma_B(b')$. We proceed by cases on S :

- $S = \mathbf{x} := E$. Follows from Lemma 2.
- $S = *x := E$. It suffices to show that, if $m(\mathbf{x})$ is an in-bounds location value l , then $\alpha(l) = b(\mathbf{x})$ and $pts(\mathbf{x}, l)$. This follows immediately from the assumption that m is in both $\gamma_B(b)$ and $\gamma_{Pts}(pts)$.
- $S = [E]$. Follows from Lemma 2.

□

Theorem 7. *The analysis generator $\tilde{\mathcal{B}}$ is sound.*

Proof. By Theorem 3, Lemma 3, and Lemma 4. □

Corollary 1. *If a points-to analysis \mathcal{Q} is SAFEDEREF-sound, the composed memory safety analysis $\tilde{\mathcal{B}} \circ \mathcal{Q}$ is SAFEDEREF-sound.*

Proof. By Theorem 7 and Theorem 4. □

Combining Corollary 1 with Theorems 5 and 6, we can compose $\tilde{\mathcal{B}}$ with any of the points-to analyses described in Section 4 and the resulting analysis will be SAFEDEREF-sound. Recall from Theorem 2 that SAFEDEREF-soundness guarantees the detection of error states. If any non-SAFEDEREF state exists in $\llbracket \mathcal{P} \rrbracket_{\mathcal{C}}$, then a non-SAFEDEREF state is represented by the composed semantics; if only SAFEDEREF states are reachable in the composed analysis then no concrete non-SAFEDEREF state is reachable—the absence of error states can be proved.

6 Related Work

Methods for combining analyses have been described in the abstract interpretation community, starting with Cousot and Cousot [7]. The focus has been on exploiting mutual refinement to achieve the most precise combined analyses, as in Gulwani and Tiwari [19] and Cousot et al. [9]. The power domain of Cousot and Cousot [7, §10.2] provides a general model for analyses with conditional semantics. We believe our notion of conditional soundness provides a simpler model which captures the behavior of a variety of interesting analyses.

Pointer analysis for C programs has been an active area of research for decades [21, 16, 28, 3, 27, 17, 12, 20, 23]. The correctness arguments for points-to algorithms are typically stated informally—each of the analyses has been developed for the purpose of program transformation and understanding, not for use in a sound verification tool. Although Hind [21] proposes the use of pointer analysis in verification, the authors are not aware of any prior work that formally addresses the soundness of verification using points-to information.

Adams et al. [1] explored the use of Das’ algorithm to prune the search space for a tpestate checker and to generate initial predicates for a software model checker. In both cases, the use of the points-to information is essentially heuristic—the correctness of the overall approach does not depend on the points-to analysis being sound.

Dor, Rodeh, and Sagiv [15] describe a variation on traditional points-to analyses intended to improve precision for a sound, inter-procedural memory safety verifier. A proof of soundness is given in Dor’s thesis [14]. However, the proof is not explicit about the obligations of the points-to analysis. We provide a more general framework for reasoning about verification using conditionally sound information.

Bruns and Chandra [4] provide a formal model for reasoning about pointer analysis based on transition systems. The focus of their work is primarily complexity and precision, rather than soundness.

Dhurjati, Kowshik, and Adve [13] define a program transformation which preserves the soundness of a flow-insensitive, equality-based points-to analysis (e.g., those of Steensgaard [27] and Lattner [23]) even for programs with memory safety errors. The use of an equality-based analysis is necessary to achieve an efficient implementation, but it limits the use of the technique in applications where a more precise analysis may be necessary, e.g., in verification. The soundness results we describe here are equally applicable to flow-sensitive, flow-insensitive, equality-based and subset-based pointer analyses.

Our abstraction for memory safety analysis is very similar to the formal models used in CCured [26] and CSSV [15]. Miné [25] describes a combined analysis for embedded control systems which incorporates points-to information. His analysis makes implementation-specific (i.e., unsound in general) assumptions about the layout of memory.

7 Conclusion

This work grew out of a simple, but puzzling question: is it possible to utilize the results of an analysis (points-to) whose soundness is dependent on a property (memory-safety) in a sound analysis for the same property? There seemed to be a circularity that could make a sound combination impossible.

Studying this question, we were led to a more precise statement of the soundness properties of points-to analysis and to the definition of conditional soundness. The final result shows that the combination is sound enough to correctly prove the absence of errors, although it may not be strong enough to point out every possible error.

We have concentrated here on points-to and memory safety analysis, but our conditional soundness framework is by no means restricted to these domains. For example, some static analyses are sound only assuming sequential consistency, that integer overflow does not occur, or that the program is free of floating point exceptions. The soundness claims of such analyses could be refined using the methods we have described in this paper.

Acknowledgments. This material is based upon work supported by the National Science Foundation under Grant No. 0341685. Additional support was provided by NSF Grant No. 0644299.

References

1. S. Adams, T. Ball, M. Das, S. Lerner, S. K. Rajamani, M. Seigle, and W. Weimer. Speeding up dataflow analysis using flow-insensitive pointer analysis. In *Static Analysis Symposium*, pages 230–246, Madrid, Spain, Sept. 2002.
2. A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley, 1988.
3. L. O. Andersen. *Program Analysis and Specialization for the C Programming Language*. PhD thesis, DIKU, University of Copenhagen, May 1994.
4. G. Bruns and S. Chandra. Searching for points-to analysis. In *Foundations of Software Engineering*, pages 61–70, Charleston, South Carolina, Nov. 2002.
5. C. L. Conway, D. Dams, K. S. Namjoshi, and C. Barrett. Points-to analysis, conditional soundness, and proving the absence of errors. In *Static Analysis Symposium*, Valencia, Spain, July 2008.
6. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Principles of Programming Languages*, pages 238–252, Los Angeles, California, 1977.
7. P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *Principles of Programming Languages*, pages 269–282, San Antonio, Texas, 1979.
8. P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, D. Monniaux, and X. Rival. The ASTRÉE analyzer. In *European Symposium on Programming*, pages 21–30, Edinburgh, Scotland, Apr. 2005.
9. P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, D. Monniaux, and X. Rival. Combination of abstractions in the ASTRÉE static analyzer. In *Asian Computing Science Conference (ASIAN)*, Tokyo, Japan, Dec. 2006.
10. P. Cousot and N. Halbwachs. Automatic discovery of linear restraints among variables of a program. In *Principles of Programming Languages*, Tucson, Arizona, Jan. 1978.
11. D. Dams and K. S. Namjoshi. Orion: Building blocks for program analyzers. In *Formal Methods for Components and Objects*, Amsterdam, The Netherlands, Nov. 2005.
12. M. Das. Unification-based pointer analysis with directional assignments. In *Programming Language Design and Implementation*, pages 35–46, Vancouver, British Columbia, 2000.
13. D. Dhurjati, S. Kowshik, and V. Adve. SAFECode: enforcing alias analysis for weakly typed languages. In *Programming Language Design and Implementation*, pages 144–157, Ottawa, Canada, June 2006.
14. N. Dor. *Automatic Verification of Program Cleanness*. PhD thesis, Tel Aviv University, Dec. 2003.

15. N. Dor, M. Rodeh, and M. Sagiv. CSSV: towards a realistic tool for statically detecting all buffer overflows in C. In *Programming Language Design and Implementation*, pages 155–167, San Diego, California, July 2003.
16. M. Emami, R. Ghiya, and L. J. Hendren. Context-sensitive interprocedural points-to analysis in the presence of function pointers. In *Programming Language Design and Implementation*, pages 242–256, June 1994.
17. J. S. Foster, M. Fähndrich, and A. Aiken. Flow-insensitive points-to analysis with term and set constraints. Technical Report UCB/CSD-97-964, University of California, Berkeley, Aug. 1997.
18. R. Ghiya, D. M. Lavery, and D. C. Sehr. On the importance of points-to analysis and other memory disambiguation methods for C programs. In *Programming Language Design and Implementation*, pages 47–58, Snowbird, Utah, June 2001.
19. S. Gulwani and A. Tiwari. Combining abstract interpreters. In *Programming Language Design and Implementation*, Ottawa, Canada, June 2006.
20. N. Heintze and O. Tardieu. Demand-driven pointer analysis. In *Programming Language Design and Implementation*, pages 24–34, Snowbird, Utah, June 2001.
21. M. Hind. Pointer analysis: Haven't we solved this problem yet? In *Program Analysis for Software Tools and Engineering*, Snowbird, Utah, June 2001.
22. ISO Standard - Programming Languages - C, Dec. 1999. ISO/IEC 9899:1999.
23. C. Lattner. *Macroscopic Data Structure Analysis and Optimization*. PhD thesis, University of Illinois at Urbana-Champaign, May 2005.
24. S. Lerner, D. Grove, and C. Chambers. Composing dataflow analyses and transformations. In *Principles of Programming Languages*, pages 270–282, Portland, Oregon, 2002.
25. A. Miné. Field-sensitive value analysis of embedded C programs with union types and pointer arithmetics. In *Languages, Compilers, and Tools for Embedded Systems*, Ottawa, Canada, 2006.
26. G. C. Necula, S. McPeak, and W. Weimer. CCured: type-safe retrofitting of legacy code. In *Principles of Programming Languages*, pages 128–139, Portland, Oregon, Jan. 2002.
27. B. Steensgaard. Points-to analysis in almost linear time. In *Principles of Programming Languages*, pages 32–41, St. Petersburg Beach, Florida, Jan. 1996.
28. R. P. Wilson and M. S. Lam. Efficient context-sensitive pointer analysis for C programs. In *Programming Language Design and Implementation*, pages 1–12, San Diego, California, June 1995.