

# Design and Implementation of Sparse Global Analyses for C-like Languages

Hakjoo Oh   Kihong Heo   Wonchan Lee   Woosuk Lee   Kwangkeun Yi

Seoul National University

{pronto,khheo,wclee,wslee,kwang}@ropas.snu.ac.kr

## Abstract

In this article we present a general method for achieving global static analyzers that are precise, sound, yet also scalable. Our method generalizes the sparse analysis techniques on top of the abstract interpretation framework to support relational as well as non-relational semantics properties for C-like languages. We first use the abstract interpretation framework to have a global static analyzer whose scalability is unattended. Upon this underlying sound static analyzer, we add our generalized sparse analysis techniques to improve its scalability while preserving the precision of the underlying analysis. Our framework determines what to prove to guarantee that the resulting sparse version should preserve the precision of the underlying analyzer.

We formally present our framework; we present that existing sparse analyses are all restricted instances of our framework; we show more semantically elaborate design examples of sparse non-relational and relational static analyses; we present their implementation results that scale to analyze up to one million lines of C programs. We also show a set of implementation techniques that turn out to be critical to economically support the sparse analysis process.

**Categories and Subject Descriptors** F.3.2 [Semantics of Programming Languages]: Program Analysis

**Keywords** Static analysis, abstract interpretation, sparse analysis

## 1. Introduction

Precise, sound, scalable yet global static analyzers have been unachievable in general. Other than almost syntactic properties, once the target property becomes slightly deep in semantics it's been a daunting challenge to achieve the four goals in a single static analyzer. This situation explains why, for example, in the static error detection tools for full C, there exists a clear dichotomy: either "bug-finders" that risk being unsound yet scalable or "verifiers" that risk being unscalable yet sound. No such tools are scalable to globally analyze million lines of C code while being sound and precise enough for practical use.

In this article we present a general method for achieving global static analyzers that are precise, sound, yet also scalable. Our approach generalizes the sparse analysis ideas on top of the abstract

interpretation framework. Since the abstract interpretation framework [8, 10] guides us to design sound yet arbitrarily precise static analyzers for any target language, we first use the framework to have a global static analyzer whose scalability is unattended. Upon this underlying sound static analyzer, we add our generalized sparse analysis techniques to improve its scalability while preserving the precision of the underlying analysis. Our framework determines what to prove to guarantee that the resulting sparse version should preserve the precision of the underlying analyzer.

Our framework bridges the gap between the two existing technologies – abstract interpretation and sparse analysis – towards the design of sound, yet scalable global static analyzers. Note that while abstract interpretation framework provides a theoretical knob to control the analysis precision without violating its correctness, the framework does not provide a knob to control the resulting analyzer's scalability preserving its precision. On the other hand, existing sparse analysis techniques [6, 13, 14, 17, 18, 22, 37, 39, 41] achieve scalability, but they are mostly algorithmic and tightly coupled with particular analyses.<sup>1</sup> The sparse techniques are not general enough to be used for an arbitrarily complicated semantic analysis.

**Contributions** Our contributions are as follows.

- We propose a general framework for designing sparse static analysis. Our framework is semantics-based and precision-preserving. We prove that our framework yields a correct sparse analysis that has the same precision as the original.
- We present a new notion of data dependency, which is a key to the precision-preserving sparse analysis. In contrast to the conventional def-use chains, sparse analysis with our data dependency is fully precise.
- We design sparse non-relational and relational analysis which are still general as themselves. We can instantiate these designs with a particular non-relational and relational abstract domains, respectively.
- We prove the practicality of our framework by experimentally demonstrating the achieved speedup of an industrial-strength static analyzer [21, 24, 26, 32–35]. The sparse analysis can analyze programs up to 1 million lines of C code with interval domain and up to 100K lines of C code with octagon domain.

**Outline** Section 2 explains our sparse analysis framework. Section 3 and 4 design sparse non-relational and relational analyses, respectively, based on our framework. Section 5 discusses several issues involved in the implementations. Section 6 presents the experimental studies. Section 7 discusses related work.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PLDI'12, June 11–16, Beijing, China.

Copyright © 2012 ACM 978-1-4503-1205-9/12/04...\$10.00

<sup>1</sup>A few techniques [7, 36] are in general settings but instead they take coarse-grained approach to sparsity.

## 2. Sparse Analysis Framework

### 2.1 Notation

Given function  $f \in A \rightarrow B$ , we write  $f|_C$  for the restriction of function  $f$  to the domain  $\text{dom}(f) \cap C$ . We write  $f \setminus_C$  for the restriction of  $f$  to the domain  $\text{dom}(f) - C$ . We abuse the notation  $f|_a$  and  $f \setminus_a$  for the domain restrictions on singleton set  $\{a\}$ . We write  $f[a \mapsto b]$  for the function got from function  $f$  by changing the value for  $a$  to  $b$ . We write  $f[a_1 \mapsto b_1, \dots, a_n \mapsto b_n]$  for  $f[a_1 \mapsto b_1] \dots [a_n \mapsto b_n]$ . We write  $f[\{a_1, \dots, a_n\} \overset{w}{\mapsto} b]$  for  $f[a_1 \mapsto f(a_1)] \sqcup b, \dots, a_n \mapsto f(a_n) \sqcup b]$  (weak update).

### 2.2 Program

A program is a tuple  $\langle \mathbb{C}, \hookrightarrow \rangle$  where  $\mathbb{C}$  is a finite set of control points and  $\hookrightarrow \subseteq \mathbb{C} \times \mathbb{C}$  is a relation that denotes control flows of the program;  $c' \hookrightarrow c$  indicates that  $c$  is a next control point of  $c'$ . Each control point is associated with a command, denoted  $\text{cmd}(c)$ . A path  $p = p_0 p_1 \dots p_n$  is a sequence of control points such that  $p_0 \hookrightarrow p_1 \hookrightarrow \dots \hookrightarrow p_n$ . We write  $\text{Paths} = \text{lfp} \lambda P. \{c_0 c_1 \mid c_0 \hookrightarrow c_1\} \cup \{p_0 \dots p_n c \mid p \in P \wedge p_n \hookrightarrow c\}$  for the set of all paths in the program.

**Collecting Semantics** Collecting semantics of program  $P$  is an invariant  $\llbracket P \rrbracket \in \mathbb{C} \rightarrow 2^{\mathbb{S}}$  that represents a set of reachable states at each control point, where the concrete domain of states,  $\mathbb{S} = \mathbb{L} \rightarrow \mathbb{V}$ , maps concrete locations ( $\mathbb{L}$ ) to concrete values ( $\mathbb{V}$ ). The collecting semantics is characterized by the least fixpoint of semantic function  $F \in (\mathbb{C} \rightarrow 2^{\mathbb{S}}) \rightarrow (\mathbb{C} \rightarrow 2^{\mathbb{S}})$  such that,

$$F(X) = \lambda c \in \mathbb{C}. f_c \left( \bigcup_{c' \hookrightarrow c} X(c') \right). \quad (1)$$

where  $f_c \in 2^{\mathbb{S}} \rightarrow 2^{\mathbb{S}}$  is a semantic function at control point  $c$ . Because our framework is independent from target languages, we leave out the definition of the concrete semantic function  $f_c$ .

### 2.3 Baseline Abstraction

We abstract the collecting semantics of program  $P$  by the following Galois connection

$$\mathbb{C} \rightarrow 2^{\mathbb{S}} \xleftarrow{\gamma} \mathbb{C} \rightarrow \hat{\mathbb{S}} \quad (2)$$

where  $\alpha$  and  $\gamma$  are pointwise liftings of abstract and concretization function  $\alpha_{\mathbb{S}}$  and  $\gamma_{\mathbb{S}}$  (such that  $2^{\mathbb{S}} \xleftarrow{\gamma_{\mathbb{S}}} \hat{\mathbb{S}}$ ), respectively.

We consider a particular, but general and practical, family of abstract domains where abstract state  $\hat{\mathbb{S}}$  is map  $\hat{\mathbb{L}} \rightarrow \hat{\mathbb{V}}$  where  $\hat{\mathbb{L}}$  is a finite set of abstract locations, and  $\hat{\mathbb{V}}$  is a (potentially infinite) set of abstract values. All non-relational abstract domains, such as intervals [8], are members of this family. Furthermore, the family covers some numerical, relational domains. Practical relational analyses exploit *packed* relationality [4, 12, 31, 40]; the abstract domain is of form  $\text{Packs} \rightarrow \hat{\mathbb{R}}$  where  $\text{Packs}$  is a set of variable groups selected to be related together.  $\hat{\mathbb{R}}$  denotes numerical constraints among variables in those groups. In such *packed* relational analysis, each variable pack is treated as an abstract location ( $\hat{\mathbb{L}}$ ) and numerical constraints amount to abstract values ( $\hat{\mathbb{V}}$ ). Examples of the numerical constraints are domain of octagons [31] and polyhedrons [11]. In practice, relational analyses are necessarily packed relational [4, 12] because of otherwise unacceptable costs.

Abstract semantics is characterized as a least fixpoint of abstract semantic function  $\hat{F} \in (\mathbb{C} \rightarrow \hat{\mathbb{S}}) \rightarrow (\mathbb{C} \rightarrow \hat{\mathbb{S}})$  defined as,

$$\hat{F}(\hat{X}) = \lambda c \in \mathbb{C}. \hat{f}_c \left( \bigcup_{c' \hookrightarrow c} \hat{X}(c') \right). \quad (3)$$

where  $\hat{f}_c \in \hat{\mathbb{S}} \rightarrow \hat{\mathbb{S}}$  is a monotone abstract semantic function for control point  $c$ . We assume that  $\hat{F}$  is sound with respect to  $F$ , i.e.,

$\alpha \circ F \sqsubseteq \hat{F} \circ \alpha$ , then the soundness of abstract semantics is followed by fixpoint transfer theorem [10].

### 2.4 Sparse Analysis by Eliminating Unnecessary Propagation

The abstract semantic function given in (3) propagates some abstract values unnecessarily. For example, suppose that we analyze statement  $x := y$  using a non-relational domain, like interval domain [8]. We know for sure that the abstract semantic function for the statement *defines* a new abstract value only at variable  $x$  and *uses* only the abstract value of variable  $y$ . Thus, it is unnecessary to propagate the whole abstract states. However, the function given in (3) blindly propagates the whole abstract states of all predecessors  $c'$  to control point  $c$ .

To make the analysis sparse, we need to eliminate this unnecessary propagation by making the semantic function propagate abstract values along data dependency, not control flows; that is, we make the semantic function propagate only the abstract values newly computed at one control point to the other where they are actually used. In the rest of this section, we explain how to make abstract semantic function (3) sparse while preserving its precision and soundness.

### 2.5 Definition and Use Set

We first need to precisely define what are “definitions” and “uses”. They are defined in terms of abstract semantics, i.e., abstract semantic function  $\hat{f}_c$ , not concrete semantics.

**Definition 1** (Definition set). *Let  $S$  be the least fixpoint  $\text{lfp} \hat{F}$  of the original semantic function  $\hat{F}$ . Definition set  $D(c)$  at control point  $c$  is a set of abstract locations that are assigned a new abstract value by abstract semantic function  $\hat{f}_c$ , i.e.*

$$D(c) \triangleq \{l \in \hat{\mathbb{L}} \mid \exists \hat{s} \sqsubseteq \bigcup_{c' \hookrightarrow c} S(c'). \hat{f}_c(\hat{s})(l) \neq \hat{s}(l)\}.$$

**Definition 2** (Use set). *Let  $S$  be the least fixpoint  $\text{lfp} \hat{F}$  of the original semantic function  $\hat{F}$ . Use set  $U(c)$  at control point  $c$  is a set of abstract locations that are used by abstract semantic function  $\hat{f}_c$ , i.e.*

$$U(c) \triangleq \{l \in \hat{\mathbb{L}} \mid \exists \hat{s} \sqsubseteq \bigcup_{c' \hookrightarrow c} S(c'). \hat{f}_c(\hat{s})|_{D(c)} \neq \hat{f}_c(\hat{s} \setminus l)|_{D(c)}\}.$$

**Example 1.** *Consider the following simple subset of  $C$ :*

$$\begin{array}{l} \text{cmd} \rightarrow x := e \mid *x := e \\ e \rightarrow x \mid \&x \mid *x. \end{array}$$

*The meaning of each statement and each expression is fairly standard. We design a pointer analysis for this as follows:*

$$\begin{aligned} \hat{s} \in \hat{\mathbb{S}} &= \text{Var} \rightarrow \mathcal{P}^{\text{Var}} \\ \hat{f}_c(\hat{s}) &= \begin{cases} \hat{s}[x \mapsto \hat{E}(e)(\hat{s})] & \text{cmd}(c) = x := e \\ \hat{s}[y \mapsto \hat{E}(e)(\hat{s})] & \text{cmd}(c) = *x := e \\ & \text{and } \hat{s}(x) = \{y\} \\ \hat{s}[\hat{s}(x) \overset{w}{\mapsto} \hat{E}(e)(\hat{s})] & \text{cmd}(c) = *x := e \end{cases} \\ \hat{E}(e)(\hat{s}) &= \begin{cases} \hat{s}(x) & e = x \\ \{x\} & e = \&x \\ \bigcup_{y \in \hat{s}(x)} \hat{s}(y) & e = *x \end{cases} \end{aligned}$$

*Now suppose that we analyze program  $\textcircled{10}x := \&y; \textcircled{11}*p := \&z; \textcircled{12}y := x$ ; (superscripts are control points). Suppose that points-to set of pointer  $p$  is  $\{x, y\}$  at control point  $\textcircled{11}$  according to the fixpoint. Definition set and use set at each control point are as follows.*

$$\begin{array}{ll} D(\textcircled{10}) &= \{x\} & U(\textcircled{10}) &= \emptyset \\ D(\textcircled{11}) &= \{x, y\} & U(\textcircled{11}) &= \{p, x, y\} \\ D(\textcircled{12}) &= \{y\} & U(\textcircled{12}) &= \{x\} \end{array}$$

Note that  $U(\textcircled{1})$  contains  $D(\textcircled{1})$  because of the weak update ( $\overset{w}{\mapsto}$ ): the semantics of weak update  $\hat{s}[l \overset{w}{\mapsto} v] = \hat{s}[l \mapsto \hat{s}(l) \sqcup v]$  is defined to use the target location  $l$ . This implicit use information, which does not explicitly appear in the program text, is naturally captured because our approach is semantics-based. Previous sparse analyses, e.g., [18], are mostly algorithmic and hence rely on extra techniques for correctness.<sup>2</sup>

## 2.6 Data Dependencies

Once identifying definition and use sets at each control point, we can discover data dependencies of abstract semantic function  $\hat{F}$  between two control points. Intuitively, if the abstract value of abstract location  $l$  defined at control point  $c_0$  is used at control point  $c_n$ , there is a data dependency between  $c_0$  and  $c_n$  on  $l$ . Formal definition of data dependency is given below:

**Definition 3** (Data dependency). *Data dependency is ternary relation  $\rightsquigarrow \subseteq \mathbb{C} \times \hat{\mathbb{L}} \times \mathbb{C}$  defined as follows:*

$$c_0 \rightsquigarrow c_n \triangleq \exists c_0 \dots c_n \in \text{Paths}, l \in \hat{\mathbb{L}}. \\ l \in D(c_0) \cap U(c_n) \wedge \forall i \in (0, n). l \notin D(c_i)$$

The definition means that if there exists a path from control point  $c_0$  to  $c_n$ , a value of abstract location  $l$  can be defined at  $c_0$  and used at  $c_n$ , and there is no intermediate control point  $c_i$  that may change the value of  $l$ , then a data dependency exists between control points  $c_0$  and  $c_n$  on location  $l$ .

**Example 2.** *In the program presented in Example 1, we can find two data dependencies,  $\textcircled{10} \rightsquigarrow \textcircled{11}$  and  $\textcircled{11} \rightsquigarrow \textcircled{12}$ .*

**Comparison with Def-use Chains** Our notion of data dependency is different from the conventional notion of def-use chains. If we want to conservatively collect all the possible def-use chains, we should exclude only the paths from definition points to use points when there exists a point that always kills the definition. However, data dependency in Definition 3 excludes a path even when there exists a point that might, but not always, kill the definition. We can slightly modify Definition 3 to express def-use chain relation  $\rightsquigarrow_{\text{du}}$  as follows:

$$c_0 \rightsquigarrow_{\text{du}} c_n \triangleq \exists c_0 \dots c_n \in \text{Paths}, l \in \hat{\mathbb{L}}. \\ l \in D(c_0) \cap U(c_n) \wedge \forall i \in (0, n). l \notin D_{\text{must}}(c_i)$$

where  $D_{\text{must}}(c) \triangleq \{l \in \hat{\mathbb{L}} \mid \forall \hat{s} \sqsubseteq \bigsqcup_{c' \leftarrow c} (\text{Ifp} \hat{F})(c'). \hat{f}_c(\hat{s})(l) \neq \hat{s}(l)\}$ .

**Example 3.** *We can find three def-use chains,  $\textcircled{10} \rightsquigarrow_{\text{du}} \textcircled{11}$ ,  $\textcircled{10} \rightsquigarrow_{\text{du}} \textcircled{12}$ , and  $\textcircled{11} \rightsquigarrow_{\text{du}} \textcircled{12}$  in Example 1. Note that, in data dependency (Example 2), the value flow  $\textcircled{10} \rightsquigarrow_{\text{du}} \textcircled{12}$  is captured by a composition of two data dependencies  $\textcircled{10} \rightsquigarrow \textcircled{11}$  and  $\textcircled{11} \rightsquigarrow \textcircled{12}$ .*

The reason why we use our notion of data dependencies instead of def-use chains, even though both preserve the soundness of the original analysis, becomes evident in Section 2.8. That is, while sparse analysis with approximated def-use chains may lose precision, data dependencies still preserve the precision of the original analysis even when approximations are involved.

## 2.7 Sparse Abstract Semantic Function

Using data dependency, we can make abstract semantic function sparse, which propagates between control points only the abstract values that participate in the fixpoint computation. Sparse abstract function  $\hat{F}_s$ , whose definition is given below, is the same as the original except that it propagates abstract values along to the data

dependency, not to control dependency:

$$\hat{F}_s(\hat{X}) = \lambda c \in \mathbb{C}. \hat{f}_c(\bigsqcup_{c_d \overset{l}{\rightsquigarrow} c} \hat{X}(c_d) | l).$$

As this definition is only different in that it is defined over data dependency ( $\rightsquigarrow$ ), we can reuse abstract semantic function  $\hat{f}_c$ , and its soundness result, from the original analysis design.

The following lemma states that the analysis result with sparse abstract semantic function is the same as the one of original analysis.

**Lemma 1** (Correctness). *Let  $S$  and  $S_s$  be  $\text{Ifp} \hat{F}$  and  $\text{Ifp} \hat{F}_s$ . Then,*

$$\forall c \in \mathbb{C}. \forall l \in D(c). S(c)(l) = S_s(c)(l).$$

The lemma guarantees that the sparse analysis result is identical to the original result only up to the entries that are defined in every control point. This is fair since the sparse analysis result does not contain the entries unnecessary for its computation.

## 2.8 Sparse Analysis with Approximated Data Dependency

Sparse analysis designed until Section 2.7 might not be practical since we can decide definition set  $D$  and use set  $U$  only with the original fixpoint  $\text{Ifp} \hat{F}$  computed.

To design a practical sparse analysis, we can approximate data dependency using an approximated definition set  $\hat{D}$  and use set  $\hat{U}$ .

**Definition 4** (Approximated Data Dependency). *Approximated data dependency is ternary relation  $\rightsquigarrow_a \subseteq \mathbb{C} \times \hat{\mathbb{L}} \times \mathbb{C}$  defined as follows:*

$$c_0 \rightsquigarrow_a c_n \triangleq \exists c_0 \dots c_n \in \text{Paths}, l \in \hat{\mathbb{L}}. \\ l \in \hat{D}(c_0) \cap \hat{U}(c_n) \wedge \forall i \in (0, n). l \notin \hat{D}(c_i)$$

The definition is the same except that it is defined using  $\hat{D}$  and  $\hat{U}$ . The derived sparse analysis is to compute the fixpoint of the following abstract semantic function:

$$\hat{F}_a(\hat{X}) = \lambda c \in \mathbb{C}. \hat{f}_c(\bigsqcup_{c_d \overset{l}{\rightsquigarrow}_a c} \hat{X}(c_d) | l).$$

One thing to note is that not all  $\hat{D}$  and  $\hat{U}$  make the derived sparse analysis compute the same result as the original. First, both  $\hat{D}(c)$  and  $\hat{U}(c)$  at each control point should be an over-approximation of  $D(c)$  and  $U(c)$ , respectively (we can easily show that the analysis computes different result if one of them is an under-approximation). Next, all spurious definitions that are included in  $\hat{D}$  but not in  $D$  should be also included in  $\hat{U}$ . The following example illustrates what happens when there exists an abstract location which is a spurious definition but is not included in the approximated use set.

**Example 4.** *Consider the same program presented in Example 1. except that we now suppose the points-to set of pointer  $p$  being  $\{y\}$ . Then, definition set and use set at each control point are as follows:*

$$\begin{array}{ll} D(\textcircled{10}) = \{x\} & U(\textcircled{10}) = \emptyset \\ D(\textcircled{11}) = \{y\} & U(\textcircled{11}) = \{p\} \\ D(\textcircled{12}) = \{y\} & U(\textcircled{12}) = \{x\}. \end{array}$$

*Note that  $U(\textcircled{11})$  does not contain  $D(\textcircled{11})$  because of strong update. The following is one example of unsafe approximation.*

$$\begin{array}{ll} \hat{D}(\textcircled{10}) = \{x\} & \hat{U}(\textcircled{10}) = \emptyset \\ \hat{D}(\textcircled{11}) = \{x, y\} & \hat{U}(\textcircled{11}) = \{p\} \\ \hat{D}(\textcircled{12}) = \{y\} & \hat{U}(\textcircled{12}) = \{x\}. \end{array}$$

*This approximation is unsafe because spurious definition  $\{x\}$  at control point  $\textcircled{11}$  is not included in approximated use set  $\hat{U}(\textcircled{11})$ . With this approximation, abstract value of  $x$  at  $\textcircled{10}$  is not propagated*

<sup>2</sup>Sparse pointer analysis [18] uses  $\chi$  and  $\mu$  functions for this purpose. .

to ⑫, while it is propagated in the original analysis (⑩  $\rightsquigarrow_a$  ⑫, but ⑩  $\rightsquigarrow_a$  ⑫). However, if  $\{x\} \subseteq \hat{U}(\textcircled{1})$ , then the abstract value will be propagated through two data dependency, ⑩  $\rightsquigarrow_a$  ⑪ and ⑪  $\rightsquigarrow_a$  ⑫. Note that  $x$  is not defined at ⑪, thus the propagated abstract value for  $x$  is not modified at ⑪.

We can formally define safe approximation of definition set and use set as follows:

**Definition 5.** Set  $\hat{D}(c)$  and  $\hat{U}(c)$  are a safe approximation of definition set  $D(c)$  and use set  $U(c)$ , respectively, if and only if

- (1)  $\hat{D}(c) \supseteq D(c) \wedge \hat{U}(c) \supseteq U(c)$ ; and
- (2)  $\hat{D}(c) - D(c) \subseteq \hat{U}(c)$ .

The remaining things is to prove that the safe approximation  $\hat{D}$  and  $\hat{U}$  yields the correct sparse analysis, which the following lemma states:

**Lemma 2** (Correctness of Safe Approximation). *Suppose sparse abstract semantic function  $\hat{F}_a$  is derived by the safe approximation  $\hat{D}$  and  $\hat{U}$ . Let  $S$  and  $S_a$  be  $\text{lfp}\hat{F}$  and  $\text{lfp}\hat{F}_a$ . Then,*

$$\forall c \in \mathbb{C}. \forall l \in \hat{D}(c). S_a(c)(l) = S(c)(l).$$

**Precision Loss with Conservative Def-use Chains** While approximated data dependency does not degrade the precision of an analysis, conservative def-use chains from approximated definition set and use set make the analysis less precise even if the approximation is safe. The following example illustrates the case of imprecision.

**Example 5.** Consider the same setting of Example 4, assuming the points-to set of  $\mathfrak{p}$  being  $\{x\}$ . Approximated definition set and use set establish the following three def-use chains: ⑩  $\rightsquigarrow_{\text{du}}$  ⑪, ⑪  $\rightsquigarrow_{\text{du}}$  ⑫, and ⑩  $\rightsquigarrow_{\text{du}}$  ⑫ (we assume here that relation  $\rightsquigarrow_{\text{du}}$  is similarly modified as in Definition 5). With these conservative def-use chains, the points-to set of  $x$  propagated to control point ⑫ is  $\{y\} \cup \{z\}$ , which is bigger set than  $\{y\}$ , the one that appears in the original analysis.

## 2.9 Designing Sparse Analysis Steps in the Framework

In summary, the design of sparse analysis within our framework is done in the following two steps:

- (1) Design a static analysis based on abstract interpretation framework [8]. Note that the abstract domain should be a member of the family explained in Section 2.3.
- (2) Design a method to find a safe approximation  $\hat{D}$  and  $\hat{U}$  of definition set  $D$  and use set  $U$  (Definition 5).

## 3. Designing Sparse Non-Relational Analysis

As a concrete example, we show how to design sparse non-relational analyses within our framework. Following Section 2.9, we proceed in two steps: (1) We design a conventional non-relational analysis based on abstract interpretation. Relying on the abstract interpretation framework [8, 9], we can flexibly design a static analysis of our interest with soundness guaranteed. (2) We design a method to find  $\hat{D}$  and  $\hat{U}$  and prove that they are safe approximations (Definition 5). The sparse analysis designed in this section is the core of our interval domain-based static analyzer,  $\text{Intervals}_{\text{sparse}}$ , which will be evaluated in Section 6.

For brevity, we restrict our presentation to the following simple subset of  $\mathbb{C}$ , where a variable has either an integer value or a pointer (i.e.  $\mathbb{V} = \mathbb{Z} + \mathbb{L}$ ):

$$\begin{aligned} \text{cmd} &\rightarrow x := e \mid *x := e \mid \{x < n\} \\ \text{where } e &\rightarrow n \mid x \mid \&x \mid *x \mid e + e \end{aligned}$$

Assignment  $x := e$  corresponds to assigning the value of expression  $e$  to variable  $x$ . Store  $*x := e$  performs indirect assignments; the value of  $e$  is assigned to the location that  $x$  points to. An assume command  $\{x < n\}$  makes the program continue only when the condition evaluates to true.

### 3.1 Step 1: Designing Non-sparse Analysis

**Abstract Domain** From the baseline abstraction (in Section 2.3), we consider a family of state abstractions  $2^{\mathbb{S}} \xrightarrow[\alpha_{\mathbb{S}}]{\gamma_{\mathbb{S}}} \hat{\mathbb{S}}$  such that, (Because it is standard, we omit the definition of  $\alpha_{\mathbb{S}}$ .)

$$\hat{\mathbb{S}} = \hat{\mathbb{L}} \rightarrow \hat{\mathbb{V}} \quad \hat{\mathbb{L}} = \text{Var} \quad \hat{\mathbb{V}} = \hat{\mathbb{Z}} \times \hat{\mathbb{P}} \quad \hat{\mathbb{P}} = 2^{\hat{\mathbb{L}}}$$

An abstract location is a program variable. An abstract value is a pair of an abstract integer  $\hat{\mathbb{Z}}$  and an abstract pointer  $\hat{\mathbb{P}}$ . A set of integers is abstracted to an abstract integer ( $2^{\mathbb{Z}} \xrightarrow[\alpha_{\mathbb{Z}}]{\gamma_{\mathbb{Z}}} \hat{\mathbb{Z}}$ ). Note that the abstraction is generic so we can choose any non-relational numeric domains of our interest, such as intervals ( $\hat{\mathbb{Z}} = \{[l, u] \mid l, u \in \mathbb{Z} \cup \{-\infty, +\infty\} \wedge l \leq u\} \cup \{\perp\}$ ). For simplicity, we do not abstract pointers (because they are finite): pointer values are kept by a points-to set ( $\hat{\mathbb{P}} = 2^{\hat{\mathbb{L}}}$ ). Other pointer abstractions are also orthogonally applicable.

**Abstract Semantics** The abstract semantics is defined by the least fixpoint of semantic function (3),  $\hat{F}$ , where the abstract semantic function  $\hat{f}_c \in \hat{\mathbb{S}} \rightarrow \hat{\mathbb{S}}$  is defined as follows:

$$\hat{f}_c(\hat{s}) = \begin{cases} \hat{s}[x \mapsto \hat{\mathcal{E}}(e)(\hat{s})] & \text{cmd}(c) = x := e \\ \hat{s}[\hat{s}(x).\hat{\mathbb{P}} \xrightarrow{w} \hat{\mathcal{E}}(e)(\hat{s})] & \text{cmd}(c) = *x := e \\ \hat{s}[x \mapsto \langle \hat{s}(x).\hat{\mathbb{Z}} \cap_{\mathbb{Z}} \alpha_{\mathbb{Z}}(\{z \in \mathbb{Z} \mid z < n\}), \hat{s}(x).\hat{\mathbb{P}} \rangle] & \text{cmd}(c) = \{x < n\} \end{cases}$$

Auxiliary function  $\hat{\mathcal{E}}(e)(\hat{s})$  computes abstract value of  $e$  under  $\hat{s}$ . Assignment  $x := e$  updates the value of  $x$ . Store  $*x := e$  weakly<sup>3</sup> updates the value of abstract locations that  $*x$  denotes.  $\{x < n\}$  confines the interval value of  $x$  according to the condition.  $\hat{\mathcal{E}} \in e \rightarrow \hat{\mathbb{S}} \rightarrow \hat{\mathbb{V}}$  is defined as follows:

$$\begin{aligned} \hat{\mathcal{E}}(n)(\hat{s}) &= \langle \alpha_{\mathbb{Z}}(\{n\}), \perp \rangle \\ \hat{\mathcal{E}}(x)(\hat{s}) &= \hat{s}(x) \\ \hat{\mathcal{E}}(\&x)(\hat{s}) &= \langle \perp, \{x\} \rangle \\ \hat{\mathcal{E}}(*x)(\hat{s}) &= \bigsqcup \{ \hat{s}(a) \mid a \in \hat{s}(x).\hat{\mathbb{P}} \} \\ \hat{\mathcal{E}}(e_1 + e_2)(\hat{s}) &= \langle v_1.\hat{\mathbb{Z}} \dot{+}_{\mathbb{Z}} v_2.\hat{\mathbb{Z}}, v_1.\hat{\mathbb{P}} \cup v_2.\hat{\mathbb{P}} \rangle \\ &\quad \text{where } v_1 = \hat{\mathcal{E}}(e_1)(\hat{s}), v_2 = \hat{\mathcal{E}}(e_2)(\hat{s}) \end{aligned}$$

Note that the above analysis is parameterized by an abstract numeric domain  $\hat{\mathbb{Z}}$  and sound operators  $\dot{+}_{\mathbb{Z}}$  and  $\cap_{\mathbb{Z}}$ .

### 3.2 Step 2: Finding Definitions and Uses

The second step is to find safe approximations of definitions and uses. The framework provides a mathematical definitions regarding correctness but does not provide how to find safe  $\hat{D}$  and  $\hat{U}$ . In the rest part of this section, we present a semantics-based, systematic way to find them.

We propose to find  $\hat{D}$  and  $\hat{U}$  from a conservative approximation of  $\hat{F}$ . We call the approximated analysis by pre-analysis. Let  $\hat{\mathbb{D}}_{\text{pre}}$  and  $\hat{F}_{\text{pre}}$  be the domain and semantic function of such a pre-analysis, which satisfies the following two conditions.

$$\mathbb{C} \rightarrow \hat{\mathbb{S}} \xrightarrow[\alpha_{\text{pre}}]{\gamma_{\text{pre}}} \hat{\mathbb{D}}_{\text{pre}} \quad \alpha_{\text{pre}} \circ \hat{F} \sqsubseteq \hat{F}_{\text{pre}} \circ \alpha_{\text{pre}}$$

By abstract interpretation framework [8, 9], such a pre-analysis is guaranteed to be conservative, i.e.,  $\alpha_{\text{pre}}(\text{lfp}\hat{F}) \sqsubseteq \text{lfp}\hat{F}_{\text{pre}}$ . As an

<sup>3</sup>For brevity, we consider only weak updates. Applying strong update is orthogonal to our sparse analysis design.

example, in experiments (Section 6), we use a simple abstraction as follows:

$$\mathbb{C} \rightarrow \hat{\mathbb{S}} \xleftarrow[\alpha_{pre}]{\gamma_{pre}} \hat{\mathbb{S}} \quad \begin{array}{l} \alpha_{pre} = \lambda \hat{X}. \sqcup \{ \hat{X}(c) \mid c \in \text{dom}(\hat{X}) \} \\ \hat{F}_{pre} = \lambda \hat{s}. \sqcup_{c \in \mathbb{C}} \hat{f}_c(\hat{s}) \end{array}$$

The abstraction ignores the control flows of programs and computes a single global invariant (a.k.a., flow-insensitivity).

We now define  $\hat{D}$  and  $\hat{U}$  by using pre-analysis. Let  $\hat{T}_{pre} \in \mathbb{C} \rightarrow \hat{\mathbb{S}}$  be the pre-analysis result in terms of original analysis, i.e.,  $\hat{T}_{pre} = \gamma_{pre}(\text{Ifp} \hat{F}_{pre})$ . The definitions of  $\hat{D}$  and  $\hat{U}$  are naturally derived from the semantic definition of  $\hat{f}_c$ .

$$\hat{D}(c) = \begin{cases} \{x\} & \text{cmd}(c) = x := e \\ \hat{T}_{pre}(c)(x).\hat{\mathbb{P}} & \text{cmd}(c) = *x := e \\ \{x\} & \text{cmd}(c) = \{x < n\} \end{cases}$$

$\hat{D}$  is defined to include locations whose values are potentially defined (changed). In the definition of  $\hat{f}_c$  for  $x := e$  and  $\{x < n\}$ , we notice that abstract location  $x$  may be defined. In  $*x := e$ , we see that  $\hat{f}_c$  may define locations  $\hat{s}(x).\hat{\mathbb{P}}$  for a given input state  $\hat{s}$  at program point  $c$ . Here, we use the pre-analysis: because we cannot have the input state  $\hat{s}$  prior to the analysis, we instead use its conservative abstraction  $\hat{T}_{pre}(c)$ . Such  $\hat{D}$  satisfies the safe approximation condition (Definition 5), because we collect all potentially defined locations, pre-analysis is conservative, and  $\hat{f}$  is monotone.

Before defining  $\hat{U}$ , we define an auxiliary function  $\mathcal{U} \in e \rightarrow \hat{\mathbb{S}} \rightarrow 2^{\hat{\mathbb{L}}}$ . Given expression  $e$  and state  $\hat{s}$ ,  $\mathcal{U}(e)(\hat{s})$  finds the set of abstract locations that are referenced during the evaluation of  $\hat{\mathcal{E}}(e)(\hat{s})$ . Thus,  $\mathcal{U}$  is naturally derived from the definition of  $\hat{\mathcal{E}}$ .

$$\begin{array}{l} \mathcal{U}(n)(\hat{s}) = \emptyset \\ \mathcal{U}(x)(\hat{s}) = \{x\} \\ \mathcal{U}(\&x)(\hat{s}) = \emptyset \\ \mathcal{U}(*x)(\hat{s}) = \{x\} \cup \hat{s}(x).\hat{\mathbb{P}} \\ \mathcal{U}(e_1 + e_2)(\hat{s}) = \mathcal{U}(e_1)(\hat{s}) \cup \mathcal{U}(e_2)(\hat{s}) \end{array}$$

When  $e$  is either  $n$  or  $\&x$ ,  $\hat{\mathcal{E}}$  does not refer any abstract location. Because  $\hat{\mathcal{E}}(x)(\hat{s})$  references abstract location  $x$ ,  $\mathcal{U}(x)(\hat{s})$  is defined by  $\{x\}$ .  $\hat{\mathcal{E}}(*x)(\hat{s})$  references location  $x$  and each location  $a \in \hat{s}(x)$ , thus the set of referenced locations is  $\{x\} \cup \hat{s}(x).\hat{\mathbb{P}}$ .  $\hat{U}$  is defined as follows: (For brevity, let  $\hat{s}_c = \hat{T}_{pre}(c)$ )

$$\hat{U}(c) = \begin{cases} \mathcal{U}(e)(\hat{s}_c) & \text{cmd}(c) = x := e \\ \{x\} \cup \hat{s}_c(x).\hat{\mathbb{P}} \cup \mathcal{U}(e)(\hat{s}_c) & \text{cmd}(c) = *x := e \\ \{x\} & \text{cmd}(c) = \{x < n\} \end{cases}$$

Using  $\hat{T}_{pre}$  and  $\mathcal{U}$ , we collect abstract locations that are potentially used during the evaluation of  $e$ . Because  $\hat{f}_c$  is defined to refer to abstract location  $x$  in  $*x := e$  and  $\{x < n\}$ ,  $\mathcal{U}$  additionally includes  $x$ . Note that, in  $*x := e$ ,  $\hat{U}(c)$  includes  $\hat{s}_c(x).\hat{\mathbb{P}}$  because  $\hat{f}_c$  performs weak updates.

**Lemma 3.**  $\hat{D}$  and  $\hat{U}$  are safe approximations.

**Sparse Pointer Analysis as Instances** We can instantiate this design of non-relational analysis to the recent two successful scalable sparse analysis presented in [17, 18].

Semi-sparse analysis [17] applies sparse analysis only for top-level variables whose addresses are never taken. We do the same thing by designing pre-analysis which computes a fixpoint  $\hat{T}_{pre}$  such that  $\hat{T}_{pre}(c)(x).\hat{\mathbb{P}} = \hat{\mathbb{L}}$  for all  $x$ s that are not top-level variables.

Staged Flow-Sensitive Analysis [18] uses auxiliary flow-insensitive pointer analysis to get an over-approximation of def-use information on pointer variables. By coincidence, our sparse non-relational analysis already does the same analysis for pointer variables except it also tracks numeric constraints of variables. We can design pre-analysis whose precision is incomparable to the

original one, as in [18], although the framework cannot guarantee the correctness anymore.

## 4. Designing Sparse Relational Analysis

As another example, we show how to design sparse relational analyses. In Section 4.1, we define the family of relational analyses that our framework considers. In Section 4.2, we define safe  $\hat{D}$  and  $\hat{U}$  for the analysis.

We consider *packed* relational analysis [4, 31]. A pack is a set of variables selected to be related together. In the rest of this section, we assume a set of variable packs,  $Packs \subseteq 2^{Var}$  such that  $\bigcup Packs = Var$ , are given by users or a pre-analysis [12, 31]. In a packed relational analysis, abstract states ( $\hat{\mathbb{S}}$ ) map variable packs ( $Packs$ ) to a relational domain ( $\hat{\mathbb{R}}$ ), i.e.,  $\hat{\mathbb{S}} = Packs \rightarrow \hat{\mathbb{R}}$ , which is a member of the abstract domains that our framework deals with (Section 2.3). Note that the packed relational domain ( $\hat{\mathbb{S}}$ ) generalizes  $\hat{\mathbb{R}}$ :  $\hat{\mathbb{R}}$  is a special case of  $\hat{\mathbb{S}}$  in which  $Packs$  is the entire variables set ( $Packs = \{Var\}$ ).

The distinguishing feature of sparse relational analysis is that definition sets and use sets are defined in terms of variable packs. Consider analyzing statement  $x := 1$ . In non-relational analysis, variable  $x$  may be defined. In packed relational analysis, because an abstract location is a pack, all the variable packs that contain  $x$  may be defined. For example, when  $Packs$  is  $\{\langle\langle x, y \rangle\rangle, \langle\langle x, z \rangle\rangle, \langle\langle y, z \rangle\rangle\}$  (we write  $\langle\langle x_1, \dots, x_n \rangle\rangle$  for the pack of variables  $x_1, \dots, x_n$ ), variable packs that are defined in the statement are  $\langle\langle x, y \rangle\rangle$  and  $\langle\langle x, z \rangle\rangle$ . Naturally, data dependencies are defined in terms of variable packs, i.e.,  $\rightsquigarrow \subseteq \mathbb{C} \times Packs \times \mathbb{C}$

### 4.1 Step 1: Designing Non-sparse Analysis

For brevity, we consider the following pointer-free language: ( $\mathbb{S} = Var \rightarrow \mathbb{Z}$ ).

$$cmd \rightarrow x := e \mid \{x < n\} \quad \text{where } e \rightarrow n \mid x \mid e + e$$

Including pointers in the language does not require novelty but verbosity. We focus only on the key differences between non-relational and relational sparse analysis designs.

**Abstract Domain** From the baseline abstraction (in Section 2.3), we consider a family of state abstractions  $2^{\mathbb{S}} \xleftarrow[\alpha_{\mathbb{S}}]{\gamma_{\mathbb{S}}} \hat{\mathbb{S}}$  such that, ( $\alpha_{\mathbb{S}}$  is defined using  $\alpha_{\hat{\mathbb{R}}}$  such that  $2^{\mathbb{S}} \xleftarrow[\alpha_{\hat{\mathbb{R}}}]{\gamma_{\hat{\mathbb{R}}}} \hat{\mathbb{R}}$ .)

$$\hat{\mathbb{S}} = \hat{\mathbb{L}} \rightarrow \hat{\mathbb{V}} \quad \hat{\mathbb{L}} = Packs \quad \hat{\mathbb{V}} = \hat{\mathbb{R}}$$

Note that the abstraction is generic so we can choose any relational domain for  $\hat{\mathbb{R}}$ .

**Abstract Semantics** In packed relational analysis, we sometimes need to know actual values (such as ranges) of variables. For example, suppose we analyze  $a := b$  with  $Packs = \{\langle\langle a, c \rangle\rangle, \langle\langle b, c \rangle\rangle\}$ . Analyzing the statement amounts to updating the abstract value for pack  $\langle\langle a, c \rangle\rangle$ . However, because variable  $b$  is not contained in the pack, we need to obtain the value of  $b$  from the abstract value associated with  $\langle\langle b, c \rangle\rangle$ . Here, the value for  $b$  is obtained by projecting the relational domain elements for  $\langle\langle b, c \rangle\rangle$  into a non-relational value, such as intervals. To this end, we transform the original program into an internal form that replaces such variables with their actual values: suppose the actual value of  $b$  in terms of intervals is  $\langle 1, 2 \rangle$  then  $a := b$  is transformed into  $a := \langle 1, 2 \rangle$ . Formally, we assume abstract semantic function  $\hat{\mathcal{R}} \in cmd^{rel} \rightarrow \hat{\mathbb{R}} \rightarrow \hat{\mathbb{R}}$  for relational domain  $\hat{\mathbb{R}}$  is defined over the following internal language:

$$cmd^{rel} \rightarrow x := e^{rel} \mid \{x < \hat{\mathbb{Z}}\} \quad \text{where } e^{rel} \rightarrow \hat{\mathbb{Z}} \mid x \mid e^{rel} + e^{rel}$$

where  $\hat{\mathbb{Z}}$  is a (non-relational) abstract integer ( $2^{\mathbb{Z}} \xrightarrow[\alpha_{\mathbb{Z}}]{\gamma_{\mathbb{Z}}} \hat{\mathbb{Z}}$ ). Note that this language is not for program codes, but for our semantics definition.

We now define the semantics of the packed relational analysis. The abstract semantics is defined by the least fixpoint of semantic function (3), where the abstract semantic function  $\hat{f}_c$  is defined as follows:

$$\hat{f}_c(\hat{s}) = \hat{s}[p_1 \mapsto \hat{\mathcal{R}}(cmd_1)(\hat{s}(p_1)), \dots, p_k \mapsto \hat{\mathcal{R}}(cmd_k)(\hat{s}(p_k))]$$

where

$$\begin{aligned} \{p_1, \dots, p_k\} &= \text{pack}(x) \\ cmd_i &= \mathcal{T}(p_i)(\hat{s})(cmd(c)) \end{aligned}$$

For variable  $x$ ,  $\text{pack}(x)$  returns the set of packs that contains  $x$ , i.e.,  $\text{pack}(x) = \{p \in \text{Packs} \mid x \in p\}$ . For both  $x := e$  and  $\{\{x < n\}\}$ , we update only the packs that include  $x$ .  $\mathcal{T}$  is the function that transforms  $cmd$  into  $cmd^{rel}$ . Given a variable pack  $p$ , state  $\hat{s}$ , and command  $cmd$ ,  $\mathcal{T}(p)(\hat{s}) \in cmd \rightarrow cmd^{rel}$  returns transformed command for a given command.

$$\begin{aligned} \mathcal{T}(p)(\hat{s})(x := e) &= x := \mathcal{T}_e(p)(\hat{s})(e) \\ \mathcal{T}(p)(\hat{s})(\{\{x < n\}\}) &= \{\{x < \mathcal{T}_e(p)(\hat{s})(n)\}\} \end{aligned}$$

where  $\mathcal{T}_e(p)(\hat{s}) \in e \rightarrow e^{rel}$  transforms expressions:

$$\begin{aligned} \mathcal{T}_e(p)(\hat{s})(n) &= \alpha_{\mathbb{Z}}(\{n\}) \\ \mathcal{T}_e(p)(\hat{s})(x) &= \begin{cases} x & \text{if } x \in p \\ \pi_x(\hat{s}) & \text{otherwise} \end{cases} \\ \mathcal{T}_e(p)(\hat{s})(e_1 + e_2) &= \mathcal{T}_e(p)(\hat{s})(e_1) + \mathcal{T}_e(p)(\hat{s})(e_2) \end{aligned}$$

where  $\pi_x \in \hat{\mathbb{S}} \rightarrow \hat{\mathbb{Z}}$  is a function that projects a relational domain element onto variable  $x$  to obtain its abstract integer value. To be safe,  $\pi_x$  should satisfies the following condition:

$$\forall \hat{s} \in \hat{\mathbb{S}}. \pi_x(\hat{s}) \sqsupseteq \alpha_{\hat{\mathbb{Z}}}(\{s(x) \mid s \in \gamma_{\hat{\mathbb{R}}}(\prod_{p \in \text{pack}(x)} \hat{s}(p))\})$$

## 4.2 Step 2: Finding Definitions and Uses

We now approximate  $\hat{D}$  and  $\hat{U}$ . In the previous section, we already presented a general, semantics-based method to safely approximate  $\hat{D}$  and  $\hat{U}$  for a given abstract semantics. Because our language in this section is pointer-free, simple syntactic method is enough for our purpose.

The distinguishing feature of sparse relational analysis is that the entities that are defined and used are variable packs, not each variable. From the definition of  $\hat{f}_c$ , we notice that packs  $\text{pack}(x)$  are potentially defined both in assignment and assume:

$$\hat{D}(c) = \begin{cases} \text{pack}(x) & cmd(c) = x := e \\ \text{pack}(x) & cmd(c) = \{\{x < n\}\} \end{cases}$$

$\hat{U}$  is defined depending on the definition of  $\pi_x$ . On the assumption that  $\text{Packs}$  contains singleton packs of all program variables, we may define  $\pi_x$  as follows:

$$\pi_x(\hat{s}) = \pi_x^{rel}(\hat{s}(\{x\}))$$

where  $\pi_x^{rel} \in \hat{\mathbb{R}} \rightarrow \hat{\mathbb{Z}}$  which project a relational domain element onto variable  $x$  to obtain its abstract integer value, which is supplied by each relational domain, e.g., see [31]. In section 6, we implement our analyzer based on this definition of  $\pi_x$ .

Now, we define  $\hat{U}$  as follows:

$$\hat{U}(c) = \begin{cases} \text{pack}(x) \cup \{\{l\} \mid l \in \mathcal{U}(e)\} & cmd(c) = x := e \\ \text{pack}(x) & cmd(c) = \{\{x < n\}\} \end{cases}$$

where we  $\mathcal{U}(e)$  denotes the set of variables that appear inside expression  $e$ .

**Lemma 4.**  $\hat{D}$  and  $\hat{U}$  are safe approximations.

## 5. Implementation Techniques

Implementing sparse analysis presents unique challenges regarding construction and management of data dependencies. Because data dependencies for realistic programs are very complex, it is a key to practical sparse analyzers to generate data dependencies efficiently in space and time. We describe the basic algorithm we used for data dependency generation, and discuss two issues that we experienced significant performance variations depending on different implementation choices.

**Generation of Data Dependencies** We use the standard SSA algorithm to generate data dependencies. Because our notion of data dependencies ( $\rightsquigarrow$ , Definition 3) equals to def-use chains with  $\hat{D}$  and  $\hat{U}$  being treated as must-definitions and must-uses, any def-use chain generation algorithms (e.g., reaching definition analysis, SSA algorithm, etc) can be used. We use SSA generation because it is fast and reduces the size of def-use chains [41].

**Interprocedural Extension** With semantics-based approach in mind, interprocedural sparse analysis is no more difficult than its intraprocedural counterpart. Designing a method to find safe definitions and uses for semantic functions regarding procedure calls is all that we need for interprocedural extension.

However, during the implementation, we noticed that this natural extension may not be scalable in practice. The main problem was due to unexpected spurious dependencies among procedures. Consider the following code and suppose we compute data dependencies for global variable  $x$ .

```
int f() { x=0; 1h(); a=x; 2 }
int h() { ... } // does not use variable x
int g() { x=1; 3h(); b=x; 4 }
```

Data dependencies for  $x$  not only include  $1 \rightsquigarrow 2$  and  $3 \rightsquigarrow 4$  but also include spurious dependencies  $1 \rightsquigarrow 4$  and  $3 \rightsquigarrow 2$ , because there are control flow paths from 1 to 4 (as well as 3 to 2) via the common procedure calls to  $h$ . In real C programs, thousands of global variables exist and procedures are called from many different call-sites, which generates overwhelming number of spurious dependencies. In our experiments, such spurious dependencies made the analysis hardly scalable. Staged pointer analysis algorithm [18] takes this approach but no performance problem was reported; we guess that this is because pointer analysis typically ignores non-pointer statements (by sparse evaluation techniques [7, 36]) or locations—for example, number of global pointer variables are just small subset of the entire globals. However, our analyzers trace all semantics of C, i.e., value flows of all types including pointers and numbers.

Thus, we generate data dependencies separately for each procedure. In this approach, we need to specially handle procedure calls: we treat a procedure call as a definition (resp., use) of all abstract locations defined (resp., used) by the callee. In effect, data dependencies are generated in a way dependencies across procedure boundaries always pass through call statements. In addition, we suppose the entry point of each procedure defines all the abstract locations that are used in the body of the procedure (we also handle procedure returns in a similar way). Thus, we no longer need to compute data dependencies over the entire programs but each procedure separately. After generating dependencies inside each procedure, we connect inter-dependences (from call statements to the entry of callee). With this approach, above spurious dependencies reduces, because variable  $x$  is not propagated to procedure  $h$  (because  $h$  does not use  $x$ ). However, there is a problem with this method; data dependencies are not fully sparse. For example, consider a call chain  $f \rightarrow g \rightarrow h$  and suppose  $x$  is defined in procedure  $f$  and used in procedure  $h$ . Even when  $x$  is not used inside  $g$ , value of  $x$  is propagated to  $g$ . In our experiments, this incomplete sparse-

ness of the analysis could not make the resulting sparse analysis scalable enough. We solved the problem by applying the following optimization to the data dependencies ( $\rightsquigarrow$ ) until convergence: suppose  $a \rightsquigarrow b$ ,  $b \rightsquigarrow c$ , and that  $l$  is not defined nor used in  $b$ , then we remove those two dependencies and add  $a \rightsquigarrow c$ . This optimization is clearly sound, and makes the analysis more sparse, leading to significant speed up.

**Using BDDs in Representing Data Dependencies** The second practical issue is memory consumption of data dependencies. Analyzing real C programs must deal with hundreds of thousands of statements and abstract locations. Thus, naive representations for the data dependencies immediately makes memory problems. For example, in analyzing ghostscript-9.00 (the largest benchmark in Table 1), the data dependencies consist of 201 K abstract locations spanning over 2.8M statements. Thus, storing dependency relation  $\rightsquigarrow$  in a naive set-based implementation, which keeps a map ( $\in \mathbb{C} \times \mathbb{C} \rightarrow 2^{\mathbb{L}}$ ), did not work for such large programs (It only worked for programs of moderate sizes less than 150 KLOC). Fortunately, the dependency relation is highly redundant, making it a good application of BDDs. For example,  $\langle c_1, c_3, l \rangle \in (\rightsquigarrow)$  and  $\langle c_2, c_3, l \rangle$  are different but share the common suffix, and  $\langle c_1, c_2, l_1 \rangle$  and  $\langle c_1, c_2, l_2 \rangle$  are different but share the common prefix. BDDs can effectively share such common suffixes and prefixes. We treat each relation  $\langle c_1, c_2, l \rangle$ , by bit-encoding each control point and abstract location, as a boolean function that is naturally represented by BDDs. This way of using BDDs greatly reduced memory costs. For example, for vim60 (227 KLOC), set-based representation of data dependencies required more than 24 GB of memory but BDD-implementation just required 1 GB. No particular dynamic variable ordering was necessary in our case.

## 6. Experiments

In this section, we evaluate sparse non-relational and relational static analyses designed in Section 3 and Section 4, respectively. The evaluation was performed on top of SPARROW [21, 23, 24, 26, 32–35], an industrial-strength static analyzer for C programs.

For the non-relational analysis, we use the interval domain [8], a representative non-relational domain that is widely used in practice [1, 2, 4, 12, 24]. For the relational analysis, we use the octagon domain [31], a representative relational domain whose effectiveness is well-known in practice [4, 12, 26, 40].

We have analyzed 18 software packages. Table 1 shows characteristics of our benchmark programs. The benchmarks are various open-source applications, and most of them are from GNU open-source projects. Standard library calls are summarized using hand-crafted function stubs. For other unknown procedure calls to external code, we assume that the procedure returns arbitrary values and has no side-effect. Procedures that are unreachable from the main procedure, such as callbacks, are made to be explicitly called from the main procedure. All experiments were done on a Linux 2.6 system running on a single core of Intel 3.07 GHz box with 24 GB of main memory.

### 6.1 Interval Domain-based Sparse Analysis

**Setting** The baseline analyzer,  $\text{Interval}_{\text{base}}$ , is the global abstract interpretation engine of SPARROW. The abstract domain of the analysis is an extension of the one defined in Section 3 to support additional C features such as arrays and structures. The analysis abstracts an array by a set of tuples of base address, offset, and size. Abstraction of dynamically allocated array is similarly handled except that base addresses are abstracted by their allocation-sites. A structure is abstracted by a tuple of base address and set of field locations (the analysis is field-sensitive). The fixpoint is computed

by a worklist algorithm using the conventional widening operator [8] for interval domain. Details of the analysis can be found in [21, 24, 26, 32–35].

The baseline analyzer is not a straw-man but much engineering effort has been put to its implementation. In particular, the analysis exploits the technique of localization [35, 38, 42], which localizes the analysis so that each code block is analyzed with only the to-be-accessed parts of the input state. We use the access-based technique [35], which was shown to be faster by up to 50x than conventional, reachability localization-based technique [35].

From the baseline, we made  $\text{Interval}_{\text{vanilla}}$  and  $\text{Interval}_{\text{sparse}}$ .  $\text{Interval}_{\text{vanilla}}$  is identical to  $\text{Interval}_{\text{base}}$  except that  $\text{Interval}_{\text{vanilla}}$  does not perform the access-based localization. We compare the performance between  $\text{Interval}_{\text{vanilla}}$  and  $\text{Interval}_{\text{base}}$  just to check that our baseline analyzer is not a straw-man.  $\text{Interval}_{\text{sparse}}$  is the sparse version derived from the baseline. The sparse analysis consists of three steps: pre-analysis (to approximate def-use sets), data dependency generation, and actual fixpoint computation. As described in Section 3, we use a flow-insensitive pre-analysis. The fixpoint of sparse abstract transfer function is computed by a worklist-based fixpoint algorithm. The analyzers are written in OCaml. We use the BuDDy library [28] for BDD implementation.

**Result** Table 2 gives the analysis time and peak memory consumption of the three analyzers. Because three analyzers share a common frontend, we report only the analysis time. For  $\text{Interval}_{\text{base}}$ , the time represents entire analysis time, including the pre-analysis [35]. For  $\text{Interval}_{\text{sparse}}$ , **Dep** includes times for pre-analysis and generation of data dependencies and **Fix** represent the time for fixpoint computation.

The results show that  $\text{Interval}_{\text{base}}$  already has a competitive performance: it is faster than  $\text{Interval}_{\text{vanilla}}$  by 8–55x, saving peak memory consumption by 54–85%.  $\text{Interval}_{\text{vanilla}}$  scales to 35 KLOC before running out of time limit (24 hours). In contrast,  $\text{Interval}_{\text{base}}$  scales to 111 KLOC.

$\text{Interval}_{\text{sparse}}$  is faster than  $\text{Interval}_{\text{base}}$  by 5–110x and saves memory by 3–92%. In particular, the analysis’ scalability has been remarkably improved:  $\text{Interval}_{\text{sparse}}$  scales to 1.4M LOC, which is an order of magnitude larger than that of  $\text{Interval}_{\text{base}}$ .

There are some counterintuitive results. First, the analysis time for  $\text{Interval}_{\text{sparse}}$  does not strictly depend on program sizes. For example, analyzing emacs-22.1 (399 KLOC) requires 10 hours, taking six times more than analyzing ghostscript-9.00 (1,363 KLOC). This is mainly due to the fact that some real C programs have unexpectedly large recursive call cycles [25, 34, 43]. Column **maxSCC** in Table 1 reports the sizes of the largest recursive cycle (precisely speaking, strongly connected component) in programs. Note that some programs (such as nethack-3.3.0, vim60, and emacs-22.1) have a large cycle that contains hundreds or even thousands of procedures. Such non-trivial SCCs markedly increase analysis cost because the large cyclic dependencies among procedures make data dependencies much more complex.

Second, data dependency generation takes much longer time than actual fixpoint computation. For example, data dependency generation for ghostscript-9.00 takes 14,116 s but the fixpoint is computed in 698 s. The seemingly unbalanced timing results are partly because of the uses of BDDs in dependency construction. While BDD dramatically saves memory costs, set operations for BDDs such as addition and removal are noticeably slower than usual set operations.

### 6.2 Octagon Domain-based Sparse Analysis

**Setting** We implemented octagon domain-based static analyzers  $\text{Octagon}_{\text{vanilla}}$ ,  $\text{Octagon}_{\text{base}}$ , and  $\text{Octagon}_{\text{sparse}}$  by replacing interval domains of SPARROW with octagon domains. Non-numerical values (such as pointers, array, and structures) are handled in the

Program	LOC	Functions	Statements	Blocks	maxSCC	AbsLocs
gzip-1.2.4a	7K	132	6,446	4,152	2	1,784
bc-1.06	13K	132	10,368	4,731	1	1,619
tar-1.13	20K	221	12,199	8,586	13	3,245
less-382	23K	382	23,367	9,207	46	3,658
make-3.76.1	27K	190	14,010	9,094	57	4,527
wget-1.9	35K	433	28,958	14,537	13	6,675
screen-4.0.2	45K	588	39,693	29,498	65	12,566
a2ps-4.14	64K	980	86,867	27,565	6	17,684
bash-2.05a	105K	955	107,774	27,669	4	17,443
lsh-2.0.4	111K	1,524	137,511	27,896	13	31,164
sendmail-8.13.6	130K	756	76,630	52,505	60	19,135
nethack-3.3.0	211K	2,207	237,427	157,645	997	54,989
vim60	227K	2,770	150,950	107,629	1,668	40,979
emacs-22.1	399K	3,388	204,865	161,118	1,554	66,413
python-2.5.1	435K	2,996	241,511	99,014	723	51,859
linux-3.0	710K	13,856	345,407	300,203	493	139,667
gimp-2.6	959K	11,728	1,482,230	286,588	2	190,806
ghostscript-9.00	1,363K	12,993	2,891,500	342,293	39	201,161

**Table 1.** Benchmarks: lines of code (LOC) is obtained by running `wc` on the source before preprocessing and macro expansion. **Functions** reports the number of functions in source code. **Statements** and **Blocks** report the number of statements and basic blocks in our intermediate representation of programs (after preprocessing). **maxSCC** reports the size of the largest strongly connected component in the callgraph. **AbsLocs** reports the number of abstract locations that are generated during the interval domain-based analysis.

Programs	Interval <sub>vanilla</sub>		Interval <sub>base</sub>		Spd $\uparrow_1$	Mem $\downarrow_1$	Interval <sub>sparse</sub>						Spd $\uparrow_2$	Mem $\downarrow_2$
	Time	Mem	Time	Mem			Dep	Fix	Total	Mem	$\hat{D}(c)$	$\hat{U}(c)$		
gzip-1.2.4a	772	240	14	65	55 x	73 %	2	1	3	63	2.4	2.5	5 x	3 %
bc-1.06	1,270	276	96	126	13 x	54 %	4	3	7	75	4.6	4.9	14 x	40 %
tar-1.13	12,947	881	338	177	38 x	80 %	6	2	8	93	2.9	2.9	42 x	47 %
less-382	9,561	1,113	1,211	378	8 x	66 %	27	6	33	127	11.9	11.9	37 x	66 %
make-3.76.1	24,240	1,391	1,893	443	13 x	68 %	16	5	21	114	5.8	5.8	90 x	74 %
wget-1.9	44,092	2,546	1,214	378	36 x	85 %	8	3	11	85	2.4	2.4	110 x	78 %
screen-4.0.2	$\infty$	N/A	31,324	3,996	N/A	N/A	724	43	767	303	53.0	54.0	41 x	92 %
a2ps-4.14	$\infty$	N/A	3,200	1,392	N/A	N/A	31	9	40	353	2.6	2.8	80 x	75 %
bash-2.05a	$\infty$	N/A	1,683	1,386	N/A	N/A	45	22	67	220	3.0	3.0	25 x	84 %
lsh-2.0.4	$\infty$	N/A	45,522	5,266	N/A	N/A	391	80	471	577	21.1	21.2	97 x	89 %
sendmail-8.13.6	$\infty$	N/A	$\infty$	N/A	N/A	N/A	517	227	744	678	20.7	20.7	N/A	N/A
nethack-3.3.0	$\infty$	N/A	$\infty$	N/A	N/A	N/A	14,126	2,247	16,373	5,298	72.4	72.4	N/A	N/A
vim60	$\infty$	N/A	$\infty$	N/A	N/A	N/A	17,518	6,280	23,798	5,190	180.2	180.3	N/A	N/A
emacs-22.1	$\infty$	N/A	$\infty$	N/A	N/A	N/A	29,552	8,278	37,830	7,795	285.3	285.5	N/A	N/A
python-2.5.1	$\infty$	N/A	$\infty$	N/A	N/A	N/A	9,677	1,362	11,039	5,535	108.1	108.1	N/A	N/A
linux-3.0	$\infty$	N/A	$\infty$	N/A	N/A	N/A	26,669	6,949	33,618	20,529	76.2	74.8	N/A	N/A
gimp-2.6	$\infty$	N/A	$\infty$	N/A	N/A	N/A	3,751	123	3,874	3,602	4.1	3.9	N/A	N/A
ghostscript-9.00	$\infty$	N/A	$\infty$	N/A	N/A	N/A	14,116	698	14,814	6,384	9.7	9.7	N/A	N/A

**Table 2.** Performance of interval analysis: time (in seconds) and peak memory consumption (in megabytes) of the various versions of analyses.  $\infty$  means the analysis ran out of time (exceeded 24 hour time limit). **Dep** and **Fix** reports the time spent during data dependency analysis and actual analysis steps, respectively, of the sparse analysis. **Spd $\uparrow_1$**  is the speed-up of Interval<sub>base</sub> over Interval<sub>vanilla</sub>. **Mem $\downarrow_1$**  shows the memory savings of Interval<sub>base</sub> over Interval<sub>vanilla</sub>. **Spd $\uparrow_2$**  is the speed-up of Interval<sub>sparse</sub> over Interval<sub>base</sub>. **Mem $\downarrow_2$**  shows the memory savings of Interval<sub>sparse</sub> over Interval<sub>base</sub>.  $\hat{D}(c)$  and  $\hat{U}(c)$  show the average size of  $\hat{D}(c)$  and  $\hat{U}(c)$ , respectively.

same way as the interval analysis. Semantic functions are appropriately changed. Octagon<sub>base</sub> performs the access-based localization [35] in terms of variable packs. Octagon<sub>vanilla</sub> is the same except for the localization. Octagon<sub>sparse</sub> is the sparse version of Octagon<sub>base</sub>. To represent octagon domain, we used Apron library [20].

In all experiments, we used a syntax-directed packing strategy. Our packing heuristic is similar to Miné’s approach [12, 31], which groups abstract locations that have syntactic locality. For examples, abstract locations involved in the linear expressions or loops are grouped together. Scope of the locality is limited within each of syntactic C blocks. We also group abstract locations involved in actual and formal parameters, which is necessary to capture relations across procedure boundaries. Large packs whose sizes exceed a threshold (10) were split down into smaller ones.

**Result** While Octagon<sub>vanilla</sub> requires extremely large amount of time and memory space but Octagon<sub>base</sub> makes the analysis realistic by leveraging the access-based localization. Octagon<sub>base</sub> is able to analyze 20 KLOC within 6 hours and 588MB of memory. With

the localization, analysis speed of Octagon<sub>base</sub> increases by 10x–20x and memory consumption decreases by 50%–76%. Though Octagon<sub>base</sub> saves a lot of memory, the analysis is still not scalable at all. For example, bc-1.06 requires 5 times more memory than gzip-1.2.4a. This memory consumption is not reasonable considering program size and interval analysis result.

Thanks to sparse analysis technique, Octagon<sub>sparse</sub> becomes more practical and scales to 130 KLOC within 25 mins and 9.8 GB of memory consumption. Octagon<sub>sparse</sub> is 30–377x faster than Octagon<sub>base</sub> and saves memory consumption by 84%–95%.

### 6.3 Discussion

**Sparsity** We discuss the relation between performance and sparsity. Column  $\hat{D}(c)$  and  $\hat{U}(c)$  in Table 2 and Table 3 show how many abstract locations are defined and used for each basic block on average. It clearly shows the key observation to sparse analysis in real programs; only a few abstract locations are defined and used in each program point. For example, the interval domain-based analysis of

Programs	Octagon <sub>vanilla</sub>		Octagon <sub>base</sub>		Spd $\uparrow$ <sub>1</sub>	Mem $\downarrow$ <sub>1</sub>	Octagon <sub>sparse</sub>						Spd $\uparrow$ <sub>2</sub>	Mem $\downarrow$ <sub>2</sub>
	Time	Mem	Time	Mem			Dep	Fix	Total	Mem	D( <i>c</i> )	U( <i>c</i> )		
gzip-1.2.4a	9,649	5,744	483	1,355	20 x	76 %	2	13	15	211	2.3	2.5	30 x	84 %
bc-1.06	15,027	10,090	1,454	5,065	10 x	50 %	4	17	21	381	3.1	3.4	55 x	92 %
tar-1.13	$\infty$	N/A	21,125	13,810	N/A	N/A	6	46	52	588	2.5	2.5	377 x	95 %
less-382	$\infty$	N/A	$\infty$	N/A	N/A	N/A	5	126	131	405	7.0	7.1	N/A	N/A
make-3.76.1	$\infty$	N/A	$\infty$	N/A	N/A	N/A	16	23	39	554	4.0	4.1	N/A	N/A
wget-1.9	$\infty$	N/A	$\infty$	N/A	N/A	N/A	8	181	189	1,098	2.0	2.1	N/A	N/A
screen-4.0.2	$\infty$	N/A	$\infty$	N/A	N/A	N/A	724	6,445	7,169	19,143	38.8	39.7	N/A	N/A
a2ps-4.14	$\infty$	N/A	$\infty$	N/A	N/A	N/A	31	763	794	1,229	2.4	2.6	N/A	N/A
bash-2.05a	$\infty$	N/A	$\infty$	N/A	N/A	N/A	45	362	407	1,875	2.5	2.6	N/A	N/A
lsh-2.0.4	$\infty$	N/A	$\infty$	N/A	N/A	N/A	391	1,162	1,553	4,449	5.4	5.5	N/A	N/A
sendmail-8.13.6	$\infty$	N/A	$\infty$	N/A	N/A	N/A	517	946	1,463	9,881	15.9	16.0	N/A	N/A

**Table 3.** Performance of octagon analysis: all columns are same as those in Table 2

a2ps-4.14 defines and uses only 0.1% of all abstract locations in one program point.

One interesting observation from the experiment results is that the analysis performance is more dependent on the sparsity than the program size. For instance, even though ghostscript-9.00 is 3.5 times bigger than emacs-22.1 in terms of LOC, ghostscript-9.00 takes 2.6 times less time to analyze. Behind this phenomenon, there is a large difference on sparsity; average  $\hat{D}(c)$  size (and  $\hat{U}(c)$  size) of emacs-22.1 is 30 times bigger than the one of ghostscript-9.00.

**Variable Packing** For maximal precision, packing strategy should be more carefully devised for each target program. However, note that our purpose of experiments is to show relative performance of Octagon<sub>sparse</sub> over Octagon<sub>base</sub>, and we applied the same packing strategy for all analyzers. Though our general-purpose packing strategy is not specialized to each program, the packing strategy reasonably groups logically related variables. The average size of packs is 5–7 for our benchmarks.

## 7. Related Work

There exists a clear separation in existing sparse analyses:

- Fined-grained sparse analyses in particular settings
- Coarse-grained sparse analyses in general settings

We combine both to obtain fine-grained sparse analyses in a general setting.

**Scalable Sparse Pointer Analysis** Sparse pointer analysis techniques [17, 18, 27] are not general enough to be used for arbitrarily complicated semantic analysis. Recently, scalability of flow-sensitive pointer analysis has been greatly improved using sparse analysis; in 2009, Hardekopf et al. [17] presented a pointer analysis algorithm that scales to large code bases (up to 474 KLOC) for the first time, and after that, flow-sensitive pointer analysis becomes scalable even to millions of lines of code via a sparse analysis technique [18, 27]. We already showed that our framework subsumes two scalable sparse pointer analysis presented in [17, 18]. In addition, the techniques are tightly coupled with pointer analysis and it is not obvious how to generalize them and prove their correctness. We provide a general framework that enables a family of abstract interpretation to be automatically turned into sparse analysis versions.

**Sparse Analysis Techniques for Dataflow Analysis** Traditional sparse analysis techniques are in a simpler setting than the one postulated in our framework. Sparse analysis techniques were first pioneered for optimizing dataflow analysis [14, 37, 41]. Reif and Lewis [37] developed a sparse analysis algorithm for constant propagation and Wegman et al. [41] extended it to conditional constant propagation. Dhamdhere et al. [14] showed how to sparse partial redundancy elimination. These algorithms are fully sparse in that

precise def-use chains are syntactically identifiable and values are always propagated along to def-use chains (in an SSA form). However, these techniques only consider the programs without pointers. Without pointers, precise sparse analysis is trivial.

**Sparse Evaluation Graphs** Sparse evaluation techniques [7, 13, 19, 36] are generally applicable but have limitations in sparseness. The goal of sparse evaluation [7, 13, 19, 36] is to remove statements whose abstract semantic functions are identity function. For example, in typical pointer analyses, statements for numerical computation are considered as identity and we can remove those statements before analysis begins. Sparse evaluation techniques are not effective when the underlying analysis does not have many identity semantics, which is the case for static analyses that consider “full” semantics, including numbers and pointers.

**Localization** Localization [30, 35, 38, 42] is used in general settings but not powerful enough. When analyzing code blocks such as procedure bodies, localization attempts to remove irrelevant parts of abstract states that will not be used during the analysis. It is widely used as a key cost-reduction technique in many semantics-based static analysis, such as shape analysis [38, 42], higher-order flow analysis [30], and numeric abstract interpretation [35]. However, localization cannot avoid unnecessary propagation of abstract values along control flows.

**Scalable Global Analyzers** Our interval and octagon domain-based analyzers achieve higher scalability (up to 1 MLOC and 130 KLOC, respectively) than the previous general-purpose global analyzers. Zitser et al. [44] report that PolySpace C Verifier [29], a commercial tool for detection of runtime errors, cannot analyze sendmail because of scalability problem. Both our interval and octagon domain-based analyzers can analyze sendmail. Airac [24, 32], a general-purpose interval domain-based global static analyzer, scales only to 30 KLOC in global analysis. Recently, a significant progress has been reported by Oh et al. [35], but it still does not scale over 120 KLOC. Other similar (interval domain-based) analyzers are also not scalable to large code bases [1, 2]. Nevertheless, there have been scalable domain-specific static analyzers, like Astrée [4, 12] and CGS [40], which scale to hundreds of thousands lines of code. However, Astrée targets on programs that do not have recursion and backward gotos, which enables a very efficient interpretation-based analysis [12], and CGS is not fully flow-sensitive [40]. There are other summary-based approaches [15, 16] for scalable global analysis, which are independent of our abstract interpretation-based approach.

**BDDs** We propose a new usage of Binary Decision Diagram (BDD) [5] in program analysis. We represent data dependency relation in BDDs. Most of the previous uses are limited to compact representations of points-to sets in pointer analysis [3, 17, 18].

## Acknowledgments

We thank Lucas Brutschy and Yoonseok Ko for their contribution to the implementation of our analyzer. We thank Deokhwan Kim and all members of Programming Research Laboratory for their useful comments and suggestions. We would also like to thank the anonymous PLDI reviewers for their constructive feedback on this paper. This work was supported by the Engineering Research Center of Excellence Program of Korea Ministry of Education, Science and Technology(MEST) / National Research Foundation of Korea(NRF) (Grant 2011-0000971) and the Brain Korea 21 Project, School of Electrical Engineering and Computer Science, Seoul National University in 2011.

## References

- [1] X. Allamigeon, W. Godard, and C. Hymans. Static analysis of string manipulations in critical embedded C programs. In *SAS*, 2006.
- [2] G. Balakrishnan and T. Reps. Analyzing memory accesses in x86 binary executables. In *CC*, 2004.
- [3] M. Berndt, O. Lhoták, F. Qian, L. Hendren, and N. Umanee. Points-to analysis using bdds. In *PLDI*, 2003.
- [4] B. Blanchet, P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, D. Monniaux, and X. Rival. A static analyzer for large safety-critical software. In *PLDI*, 2003.
- [5] R. E. Bryant. Graph-based algorithms for boolean function manipulation. *IEEE TC*, 1986.
- [6] D. R. Chase, M. Wegman, and F. K. Zadeck. Analysis of pointers and structures. In *PLDI*, 1990.
- [7] J.-D. Choi, R. Cytron, and J. Ferrante. Automatic construction of sparse data flow evaluation graphs. In *POPL*, 1991.
- [8] P. Cousot and R. Cousot. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL*, 1977.
- [9] P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *POPL*, 1979.
- [10] P. Cousot and R. Cousot. Abstract interpretation frameworks. *J. Log. Comput.*, 1992.
- [11] P. Cousot and N. Halbwachs. Automatic discovery of linear restraints among variables of a program. In *POPL*, 1978.
- [12] P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, and X. Rival. Why does astrée scale up? *Formal Methods in System Design*, 2009.
- [13] R. K. Cytron and J. Ferrante. Efficiently computing  $\phi$ -nodes on-the-fly. *TOPLAS*, 1995.
- [14] D. M. Dhamdhere, B. K. Rosen, and F. K. Zadeck. How to analyze large programs efficiently and informatively. In *PLDI*, 1992.
- [15] I. Dillig, T. Dillig, and A. Aiken. Sound, complete and scalable path-sensitive analysis. In *PLDI*, 2008.
- [16] I. Dillig, T. Dillig, and A. Aiken. Precise reasoning for programs using containers. In *POPL*, 2011.
- [17] B. Hardekopf and C. Lin. Semi-sparse flow-sensitive pointer analysis. In *POPL*, 2009.
- [18] B. Hardekopf and C. Lin. Flow-sensitive pointer analysis for millions of lines of code. In *CGO*, 2011.
- [19] M. Hind and A. Pioli. Assessing the effects of flow-sensitivity on pointer alias analyses. In *SAS*, 1998.
- [20] B. Jeannot and A. Miné. Apron: A library of numerical abstract domains for static analysis. In *CAV*, 2009.
- [21] Y. Jhee, M. Jin, Y. Jung, D. Kim, S. Kong, H. Lee, H. Oh, D. Park, and K. Yi. Abstract interpretation + impure catalysts: Our Sparrow experience. Presentation at the Workshop of the 30 Years of Abstract Interpretation, San Francisco, ropas.snu.ac.kr/~char\~kwang/paper/30yai-08.pdf, January 2008.
- [22] R. Johnson and K. Pingali. Dependence-based program analysis. In *PLDI*, 1993.
- [23] Y. Jung and K. Yi. Practical memory leak detector based on parameterized procedural summaries. In *ISMM*, 2008.
- [24] Y. Jung, J. Kim, J. Shin, and K. Yi. Taming false alarms from a domain-unaware C analyzer by a bayesian statistical post analysis. In *SAS*, 2005.
- [25] C. Lattner, A. Lenharth, and V. Adve. Making Context-Sensitive Points-to Analysis with Heap Cloning Practical For The Real World. In *PLDI*, 2007.
- [26] W. Lee, W. Lee, and K. Yi. Sound non-statistical clustering of static analysis alarms. In *VMCAI*, 2012.
- [27] L. Li, C. Cifuentes, and N. Keynes. Boosting the performance of flow-sensitive points-to analysis using value flow. In *FSE*, 2011.
- [28] J. Lind-Nielsen. BuDDy, a binary decision diagram package.
- [29] MathWorks. Polyspace embedded software verification. <http://www.mathworks.com/products/polyspace/index.html>.
- [30] M. Might and O. Shivers. Improving flow analyses via GCFA: Abstract garbage collection and counting. In *ICFP*, 2006.
- [31] A. Miné. The Octagon Abstract Domain. *HOSC*, 2006.
- [32] H. Oh. Large spurious cycle in global static analyses and its algorithmic mitigation. In *APLAS*, 2009.
- [33] H. Oh and K. Yi. An algorithmic mitigation of large spurious interprocedural cycles in static analysis. *SPE*, 2010.
- [34] H. Oh and K. Yi. Access-based localization with bypassing. In *APLAS*, 2011.
- [35] H. Oh, L. Brutschy, and K. Yi. Access analysis-based tight localization of abstract memories. In *VMCAI*, 2011.
- [36] G. Ramalingam. On sparse evaluation representations. *Theoretical Computer Science*, 2002.
- [37] J. H. Reif and H. R. Lewis. Symbolic evaluation and the global value graph. In *POPL*, 1977.
- [38] N. Rinetzkly, J. Bauer, T. Reps, M. Sagiv, and R. Wilhelm. A semantics for procedure local heaps and its abstractions. In *POPL*, 2005.
- [39] T. B. Tok, S. Z. Guyer, and C. Lin. Efficient flow-sensitive interprocedural data-flow analysis in the presence of pointers. In *CC*, 2006.
- [40] A. Venet and G. Brat. Precise and efficient static array bound checking for large embedded c programs. In *PLDI*, 2004.
- [41] M. N. Wegman and F. K. Zadeck. Constant propagation with conditional branches. *TOPLAS*, 1991.
- [42] H. Yang, O. Lee, J. Berdine, C. Calcagno, B. Cook, D. Distefano, and P. O'Hearn. Scalable shape analysis for systems code. In *CAV*, 2008.
- [43] H. Yu, J. Xue, W. Huo, X. Feng, and Z. Zhang. Level by level: making flow- and context-sensitive pointer analysis scalable for millions of lines of code. In *CGO*, 2010.
- [44] M. Zitser, D. E. S. Group, and T. Leek. Testing static analysis tools using exploitable buffer overflows from open source code. In *FSE*, 2004.