

Discovering Properties about Arrays in Simple Programs^{*}

Nicolas Halbwachs Mathias Péron

Verimag - CNRS[†], Grenoble – France

Nicolas.Halbwachs@imag.fr/Mathias.Peron@imag.fr

Abstract

Array bound checking and array dependency analysis (for parallelization) have been widely studied. However, there are much less results about analyzing properties of array *contents*. In this paper, we propose a way of using abstract interpretation for *discovering* properties about array contents in some restricted cases: one-dimensional arrays, traversed by simple “for” loops. The basic idea, borrowed from [GRS05], consists in partitioning arrays into symbolic intervals (e.g., $[1, i - 1]$, $[i, i]$, $[i + 1, n]$), and in associating with each such interval I and each array A an abstract variable A_I ; the new idea is to consider *relational* abstract properties $\psi(A_I, B_I, \dots)$ about these abstract variables, and to interpret such a property pointwise on the interval I : $\forall \ell \in I, \psi(A[\ell], B[\ell], \dots)$. The abstract semantics of our simple programs according to these abstract properties has been defined and implemented in a prototype tool. The method is able, for instance, to discover that the result of an insertion sort is a sorted array, or that, in an array traversal guarded by a “sentinel”, the index stays within the bounds.

Categories and Subject Descriptors D.2.4 [Software/Program Verification]; F.3.1 [Specifying and Verifying and Reasoning about Programs]; F.3.2 [Semantics of Programming Languages]; Program analysis

General Terms Verification, Theory

Keywords Program verification, abstract interpretation, invariant synthesis, arrays, sorting algorithms, “sentinel”

1. Introduction

Although array bound checking was a motivation of the very first work on abstract interpretation [CC76], analyzing properties of array *contents* was considered only recently. The reason is, of course, that the general problem is difficult: array indexing induces complex semantics, and in particular the possibility of aliasing; moreover, since the size of an array can be large or unknown, it represents a large or unbounded number of variables. In this paper, we

^{*} This work has been partially supported by the APRON project of the “ACI Sécurité Informatique” of the French Ministry of Research.

[†] Verimag is a joint laboratory of Université Joseph Fourier, CNRS and INPG.

This is the author’s version of the work. It is posted here by permission of ACM for your personal use. Not for redistribution. The definitive version was published in SIGPLAN Notices, vol. 43, iss. 6, (june 2008)

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’08, June 7–13, 2008, Tucson, Arizona, USA.

Copyright © 2008 ACM 978-1-59593-860-2/08/06...\$5.00

propose a way of using abstract interpretation for discovering properties about array contents in some restricted cases. First, we restrict ourselves to one-dimensional arrays and “simple programs”, which manipulate arrays only by sequential traversal: typically, *for* loops incrementing (or decrementing) their index at each iteration, and accessing arrays by simple expressions (constant translations) of the loop index. Fig. 1 shows several examples of such “simple programs”, which will be used throughout the paper. Second, we consider a restricted class of properties: while a quite general kind of property about arrays A_1, A_2, \dots, A_m could be written

$$\forall \ell \in D, \psi(A_1[f_1(\ell)], \dots, A_m[f_m(\ell)], x_1, \dots, x_p) \quad (1)$$

where D is some set of values for indices, ψ is some scalar property about content values, f_1, \dots, f_m are general index functions, and x_1, \dots, x_p are scalar variables, we will only consider properties of the form

$$\forall \ell \in I, \psi(A_1[\ell + k_1], \dots, A_m[\ell + k_m], x_1, \dots, x_p) \quad (2)$$

where I is an *interval*, and k_1, \dots, k_m are integer constants.

For instance, our method will discover automatically the following properties:

- at the end of “Array maximum” (Fig. 1.a):

$$\forall \ell \in [1, n], A[\ell] \leq \max \quad (3)$$

- at the end of “Array copy” (Fig. 1.b):

$$\forall \ell \in [1, n], A[\ell] = B[\ell] \quad (4)$$

- at the end of “Insertion sort” (Fig. 1.c):

$$\forall \ell \in [2, n], A[\ell] \geq A[\ell - 1] \quad (5)$$

So, in spite of severe restrictions both on programs and properties, our method allows interesting properties to be found about non trivial programs.

Related work: The automatic analysis of properties of array contents was considered only recently. [BMS06, IHV08] study decidable logics for expressing such properties. If we restrict ourselves to automatic analysis, an important track initiated by [FQ02] concerns verification of programs with arrays using predicate abstraction [LBC03, LB04], possibly improved with counter-example guided refinement [BHMR07] and Craig interpolants [JM07]. All these approaches make use of the property to be proved, while we aim at discovering properties.

Concerning automatic program analysis methods based on a abstract interpretation, a general common approach [BCC⁺03, GMD⁺04] is by *summarizing* a collection S of variables with one auxiliary variable, say s , managed to satisfy the disjunction of properties of variables in the collection: if s satisfies a property ψ so do all the variables $v \in S$.

In [BCC⁺03], this approach is called “*array smashing*”: all the cells of an array A are subsumed by one variable a , of the same

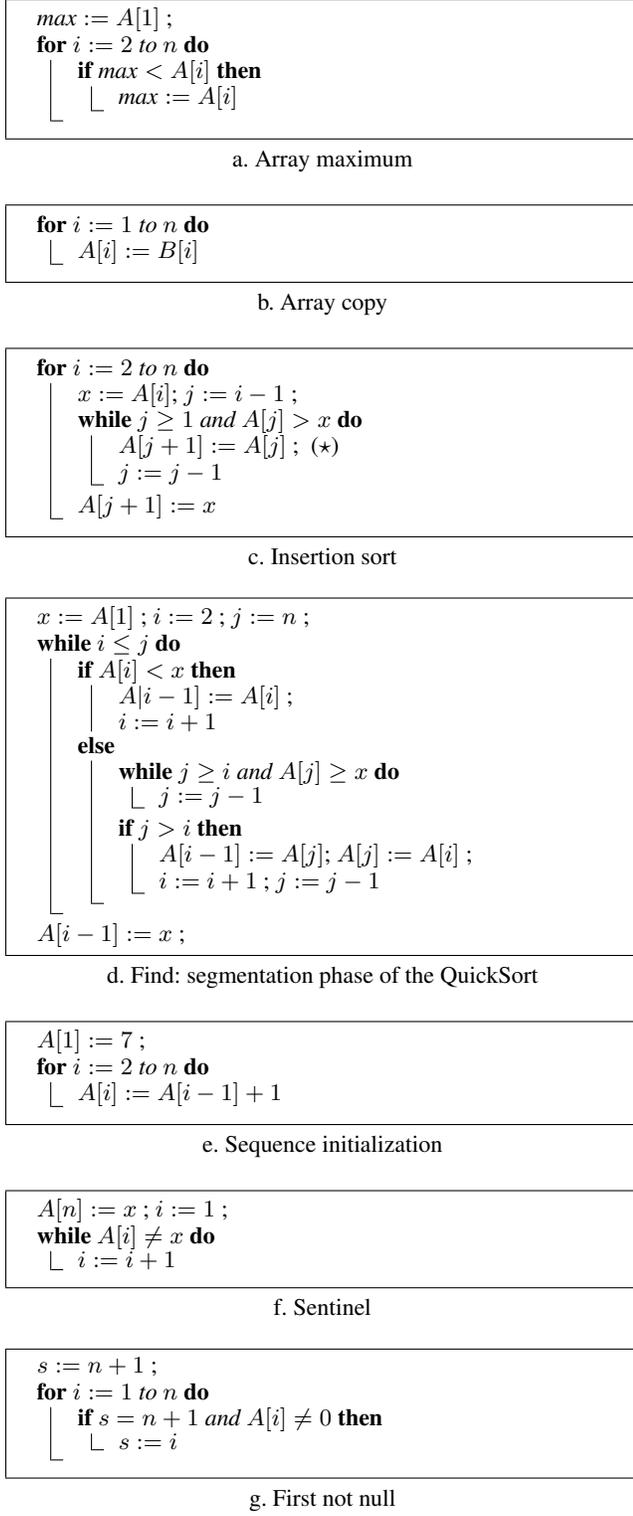


Figure 1. Some simple programs

type as the cells. Initially, a is given the strongest known property satisfied by all the initial values of the cells of A . Each assignment “ $A[i] := e$ ” to an array cell is replaced by a *weak update* of the auxiliary variable: the weak update of an expression e to a variable a (it will be noted $a \sqsubseteq e$) can be interpreted as a non deterministic choice between keeping a unchanged and actually performing the assignment (or *strong update*) $a := e$. For instance, if all the cells of an array A are known to satisfy some property ψ , and if the expression e is known to satisfy ψ' , then after an assignment $A[i] := e$ all the cells of A are known to satisfy $\psi \sqcup \psi'$ (where \sqcup is the least upper bound operator on properties), which is exactly the effect of a weak update $a \sqsubseteq e$, if a is known to satisfy ψ before. The problem with this approach is that the weak assignment can only lose information; moreover tests on individual cells don't bring any information. One needs to know (i.e., the user has to provide) an initial property satisfied by all the array cells, and the analysis can only weaken this initial knowledge. As a consequence, the results are generally very unprecise.

[GRS05, Gop07] proposes a significant improvement, by partitioning the index domain (say, $[1..n]$) into several symbolic intervals (e.g., $I_1 = [1..i-1]$, $I_2 = [i, i]$, $I_3 = [i+1..n]$), and associating with each subarray $A[I_k]$ a summary auxiliary variable a_k , managed so that

$$\psi(a_k) \Rightarrow \psi(A[\ell]), \forall \ell \in I_k \quad (6)$$

In order to reduce the loss of information due to weak updates, each cell $A[i]$ appearing in the left-hand side of an assignment or in a test constitutes a singleton in the partition, so the assignment can be interpreted as a strong update of the corresponding summary variable. This technique is able to discover our property (3) at the end of the “Array maximum” program. Concerning the “Array copy”, if we know that all the cells of B are positive, it can deduce that so are all the cells of A after the copy. So, the method is much more effective in discovering global properties of array contents. However, it cannot discover relations between the contents of different cells. For instance, it is not able to discover that $A[\ell] = B[\ell]$ for all ℓ (property (4)) at the end of “Array copy”. The paper proposes a technique to check such properties, which succeeds for our properties (4) and (5), but these candidate invariants must be provided by the user.

While in our property scheme (2), the indices will be quantified over intervals chosen in a fixed partition like in [GRS05, Gop07], [GMT08] considers more general pointwise properties of the form $\forall \ell, \varphi(\ell) \Rightarrow \psi(A[\ell] \dots)$. This generality involves difficulties, since in such properties, assumptions φ on indices must be under-approximated. Moreover, the user must provide some template for the properties that should be discovered.

[BBH⁺06] deals with the same kind of problems in a quite different context, which is the analysis of data-sensitive programs manipulating single linked lists. The considered data structure is rather different, and the method is specialized to ordering relations. Moreover, as for Cousot's *parametric abstract domains* [Cou03], the considered properties express relations between *all* the elements of two data collections, while the essence of our approach is to express *pointwise* relations.

Contribution: Our main contribution, in the present paper, is to propose a fully automatic method to *discover relations* between array cells. For that, we use the same partitioning approach as [GRS05, Gop07], but the auxiliary variables a_k are not summary variables (let us call them *slice variables*), they are interpreted in a much more restrictive sense: we generalize the interpretation (6) by giving the following sense to *relations* between slice variables:

$$\psi(a_k, b_k) \Rightarrow \psi(A[\ell], B[\ell]), \forall \ell \in I_k \quad (7)$$

Moreover, we will introduce *shift variables*, representing fixed translations of array slices, in order to be able to express relations like

$$\psi(A[\ell], A[\ell + k_1], B[\ell + k_2]), \forall \ell \in I$$

as announced at (2).

The paper is organized as follows: in Section 2, we give a better intuition of the method, by dealing informally with the ‘‘Array copy’’ example. Section 3 makes precise the kind of programs we consider, and Section 4 defines the abstract properties we shall deal with. All necessary operations on these properties are defined in Section 5. Section 6 describes our prototype implementation and the experiments performed. We conclude the paper with some perspectives.

2. An Intuitive Example

We first give a very informal intuition of the method. Let us consider the program ‘‘Array copy’’ of Fig. 1.b. As in [GRS05], since there is an assignment to $A[i]$, the set $[1, n]$ of index values is split into three intervals:

$$I_1 = [1, (i - 1)], I_2 = [i, i], I_3 = [(i + 1), n]$$

and with each array, we associate three slice variables, say:

$$A \rightarrow (a_1, a_2, a_3) \quad B \rightarrow (b_1, b_2, b_3)$$

which take their values in the same set as array contents. As said before, a property $\psi(a_k)$ should be understood as

$$\forall \ell \in I_k, \psi(A[\ell]).$$

Remarks:

1. If I_k is empty (e.g., I_1 when $i = 1$), $\psi(a_k)$ is true for all ψ (in particular *false*).
2. Intervals I_k are symbolic; in particular, the emptiness of I_k depends on the value of the index i .

With these auxiliary variables, interpreted in that way:

- the assignment ‘‘ $A[i] := B[i]$ ’’ can be abstracted into ‘‘ $a_2 := b_2$ ’’
- when i is incremented, its previous value moves from I_2 to I_1 , and its current value is extracted from I_3 to become the only element of I_2 . So, the index incrementation involves
 - a ‘‘weak update’’ $(a_1, b_1) \sqcup = (a_2, b_2)$; the assignment is weak because I_1 is not a singleton.
 - an assignment $(a_2, b_2) := (a_3, b_3)$, which is strong because I_2 is a singleton (This just an intuitive justification. In fact, such a strong assignment would result in the property $a_2 = a_3$, which makes no sense according to formula (7) since it relates variables corresponding to distinct intervals with different size. The actual interpretation is more complex, see §5.5).

In summary, instead of considering the initial program with arrays, we can analyze the following program without arrays:

```

1  $i := 1$ ;
2 while  $i \leq n$  do
3    $a_2 := b_2$ ;
4    $i := i + 1$ ;
5    $(a_1, b_1) \sqcup = (a_2, b_2)$ ;
6    $(a_2, b_2) := (a_3, b_3)$ 

```

and interpret, in the results, the properties of (a_k, b_k) as defined by formula (7).

Let us assume that classical analyses are available, which take into account simple inequalities of indices (e.g., difference bound

$\langle \text{program} \rangle$::=	$\langle \text{statement} \rangle$
$\langle \text{statement} \rangle$::=	$\langle \text{left-part} \rangle := \langle C\text{-exp} \rangle$
		$i := \langle Iexp \rangle$
		$\text{for}(i := \langle Iexp \rangle; \langle \text{cond} \rangle; \langle \text{progress} \rangle)$
		$\langle \text{statement} \rangle$
		$\text{if } \langle \text{cond} \rangle \langle \text{statement} \rangle \langle \text{statement} \rangle$
		$\langle \text{statement} \rangle ; \langle \text{statement} \rangle$
$\langle \text{left-part} \rangle$::=	$x \mid A[\langle Iexp \rangle]$
$\langle Iexp \rangle$::=	$k \mid n \mid i \mid \langle Iexp \rangle + k$
$\langle \text{cond} \rangle$::=	$\langle Icond \rangle \mid \langle Ccond \rangle$
		$\langle \text{cond} \rangle \text{ and } \langle \text{cond} \rangle \mid \langle \text{cond} \rangle \text{ or } \langle \text{cond} \rangle$
$\langle \text{progress} \rangle$::=	$++ \mid --$

Figure 2. The syntax of simple programs

matrices [Dil89] or octagons [Min01]), and equalities of array contents: so, after the assignment ‘‘ $a_2 := b_2$ ’’, we know that $a_2 = b_2$.

Now, the analysis of our program without arrays provides:

- At the first iteration:
 - After line 2: $i = 1$, which implies $I_1 = \emptyset$, so $\text{false}(a_1, b_1)$ (cf. remark 2 above).
 - after line 3: $i = 1, \text{false}(a_1, b_1), (a_2 = b_2)$
 - after line 6: $i = 2, (a_1 = b_1), (a_2 = a_3), (b_2 = b_3)$, since $\text{false}(a_1, b_1) \sqcup (a_1 = b_1) = (a_1 = b_1)$;
- At the second iteration:
 - at line 2, the property on i is widened to $1 \leq i \leq n + 1$, and a least upper bound is taken on other properties, giving $1 \leq i \leq n + 1, (a_1 = b_1)$.
 - after line 3: $2 \leq i \leq n, (a_1 = b_1), (a_2 = b_2)$
 - after line 6: $2 \leq i \leq n, (a_1 = b_1), (a_2 = a_3), (b_2 = b_3)$, and the iteration converges;
- So, the final result at the end of the program is $(i = n + 1), (a_1 = b_1)$, which can be interpreted as the expected result

$$(i = n + 1), \forall \ell \in [1, i - 1], A[\ell] = B[\ell].$$

3. Simple Programs

For simplicity, we assume that programs manipulate only data types ‘‘integer’’ (for indices), ‘‘content’’ (an arbitrary type), and ‘‘array of content’’. The following sets will be considered (with associated meta-variables):

- Integer constants: $\mathbb{Z} (\ni k)$
- Integer parameters (or non assignable integer variables, used for parametric array size): $Params (\ni n)$
- Integer variables: $Indices (\ni i, j)$
- Content constants: \mathbb{C}
- Content variables: $CVars (\ni x, y)$
- Array variables: $Arrays (\ni A, B)$

3.1 Syntax

Fig. 2 gives the abstract syntax of our simple programs. Some of the imposed restrictions are just for simplicity, others are necessary for the analysis to give good results. As usual in abstract interpretation, a more general language can be considered, either by defining specific extensions to the analysis, or by abstracting away the extended features.

$$\begin{aligned} \llbracket \cdot \rrbracket : & \langle \text{statement} \rangle \mapsto \text{States} \mapsto \text{States} \\ & \langle \text{Cexp} \rangle \mapsto \text{States} \mapsto \mathbb{C} \\ & \langle \text{Iexp} \rangle \mapsto \text{States} \mapsto \mathbb{Z} \\ & \langle \text{cond} \rangle \mapsto \text{States} \mapsto \mathbb{B} \end{aligned}$$

$$\begin{aligned} \llbracket x := \langle \text{C-exp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= (\mathcal{I}, \mathcal{C}[\llbracket \langle \text{C-exp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A})/x], \mathcal{A}) \\ \llbracket i := \langle \text{Iexp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= (\mathcal{I}[\llbracket \langle \text{Iexp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A})/i], \mathcal{C}, \mathcal{A}) \\ \llbracket A[\langle \text{Iexp} \rangle] := \langle \text{C-exp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= (\mathcal{I}, \mathcal{C}, \mathcal{A}[\mathcal{F}/A]) \\ & \text{where } \mathcal{F} = \lambda z. \begin{cases} \mathcal{A}(A)(z) & \text{if } z \neq \llbracket \langle \text{Iexp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) \\ \llbracket \langle \text{C-exp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) & \text{otherwise} \end{cases} \\ \llbracket \text{for}(i := \langle \text{Iexp} \rangle; \langle \text{cond} \rangle; \langle \text{prog} \rangle) \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= \\ & \llbracket \text{while} \langle \text{cond} \rangle \langle \text{stat} \rangle i \langle \text{prog} \rangle \rrbracket (\mathcal{I}[\llbracket \langle \text{Iexp} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A})/i], \mathcal{C}, \mathcal{A}) \\ \llbracket \text{while} \langle \text{cond} \rangle \langle \text{stat} \rangle i \langle \text{prog} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= \\ & \begin{cases} (\mathcal{I}, \mathcal{C}, \mathcal{A}) & \text{if } \llbracket \langle \text{cond} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) = \text{false} \\ \llbracket \langle \text{stat} \rangle; i \langle \text{prog} \rangle; \text{while} \langle \text{cond} \rangle \langle \text{stat} \rangle i \langle \text{prog} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) & \text{otherwise} \end{cases} \\ \llbracket i ++ \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= (\mathcal{I}[(\mathcal{I}(i) + 1)/i], \mathcal{C}, \mathcal{A}) \\ \llbracket i -- \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= (\mathcal{I}[(\mathcal{I}(i) - 1)/i], \mathcal{C}, \mathcal{A}) \\ \llbracket \text{if} \langle \text{cond} \rangle \langle \text{stat1} \rangle \langle \text{stat2} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= \\ & \begin{cases} \llbracket \langle \text{stat1} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) & \text{if } \llbracket \langle \text{cond} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) = \text{true} \\ \llbracket \langle \text{stat2} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) & \text{otherwise} \end{cases} \\ \llbracket \langle \text{stat1} \rangle; \langle \text{stat2} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A}) &= \llbracket \langle \text{stat2} \rangle \rrbracket (\llbracket \langle \text{stat1} \rangle \rrbracket (\mathcal{I}, \mathcal{C}, \mathcal{A})) \end{aligned}$$

Figure 3. Semantics of simple programs

We ignore the declarations, so a program is simply a (sequence of) statement(s). Statements can be assignments to content variables or array elements (the syntax of content expressions $\langle \text{C-exp} \rangle$ is left unspecified, since it depends on the content type), assignment to integer variables (which must *not* be loop indices), “for” loops (which are restrictions of C-like “for” constructs: in the initialization of the loop index “ $i := \langle \text{Iexp} \rangle$ ”, the expression $\langle \text{Iexp} \rangle$ may not depend on i ; the loop progress statement can only be an index incrementation ($++$) or decrementation ($--$); the loop index may not be assigned inside the loop), conditionals, and sequences. The syntax of “for” loops is convenient to express the wanted restrictions, but for detailing the analysis of examples, such loops will often be decomposed into “ $i := \langle \text{Iexp} \rangle; \text{while} \langle \text{cond} \rangle \{ \langle \text{statement} \rangle; i \langle \text{progress} \rangle \}$ ”. Index expressions are restricted to constants or parameters and sums of an index or parameter and a constant. Conditions are conjunctions and/or disjunctions of atomic conditions, whose syntax will depend on the lattices used in the analysis (see §4.1): the conditions on contents are supposed to express at least equalities (e.g., $x = A[i+1]$), those on indices are supposed to express at least potentials (e.g., $i \leq j - 3$).

3.2 Semantics

Arrays will be indexed from 1 to m , where m is the size of the array. We are concerned with the analysis of array contents, not with array bound checking, which we assume to be solved by other means. As a consequence, we don’t want to bother about access out of bounds. This is reflected by considering an array value to be a function $\mathbb{Z} \mapsto \mathbb{C}_\perp$ from all relative integers to the domain of contents completed with a \perp element: of course, an array value A is restricted to return a non \perp value exactly on an interval $[1, m]$.

Let States denote the set of states of a program. A state is a triple $(\mathcal{I}, \mathcal{C}, \mathcal{A})$, where $\mathcal{I} : (\text{Indices} \cup \text{Params}) \mapsto \mathbb{Z}$ is a valuation for indices and parameters, $\mathcal{C} : \text{CVars} \mapsto \mathbb{C}$ is a valuation for content variables, and $\mathcal{A} : \text{Arrays} \mapsto (\mathbb{Z} \mapsto \mathbb{C}_\perp)$ is a valuation for arrays. The semantics of statements is described in Fig. 3 as functions from States to States .

4. Array Content Properties

4.1 Lattices

Throughout the paper, we assume the existence of two analyses, the former concerning the behavior of indices, and the later concerning contents. In some sense, our method is parameterized by these analyses:

- The analysis of indices is based on a lattice $(L_{\mathbb{Z}}, \sqsubseteq_{\mathbb{Z}}, \sqcap_{\mathbb{Z}}, \sqcup_{\mathbb{Z}}, \top_{\mathbb{Z}}, \perp_{\mathbb{Z}})$ of properties over the set Indices of index variables. Elements of $L_{\mathbb{Z}}$ will be noted ϕ . $L_{\mathbb{Z}}$ must be a relational lattice, at least as powerful as potential constraints (i.e., systems of inequalities of the form $i - j \leq k$, $k_1 \leq i \leq k_2$, often implemented as “Difference Bound Matrices” [Dil89, ACD93]), and defining convex properties. Candidates for $L_{\mathbb{Z}}$ are potential constraints, octagons [Min01], octahedra [CC04], or polyhedra [CH78]. In this paper we will consider $L_{\mathbb{Z}}$ to be the lattice of potential constraints. We will also use an extension $L'_{\mathbb{Z}}$ of $L_{\mathbb{Z}}$, expressing properties over $\text{Indices} \cup \{\ell\}$, where ℓ is a new special variable (used for quantification). Elements of $L'_{\mathbb{Z}}$ will be noted φ . The same notations will be used for operations in $L_{\mathbb{Z}}$ and $L'_{\mathbb{Z}}$, and $L_{\mathbb{Z}}$ will often be implicitly plugged into $L'_{\mathbb{Z}}$.
- The analysis of contents is based on a lattice $(L_{\mathbb{C}}, \sqsubseteq_{\mathbb{C}}, \sqcap_{\mathbb{C}}, \sqcup_{\mathbb{C}}, \top_{\mathbb{C}}, \perp_{\mathbb{C}})$, of which we only assume that it is able to express equality relations. Elements of $L_{\mathbb{C}}$ will be noted ψ .

In many classical examples, it happens that array contents are numbers. In that case, we can choose $L_{\mathbb{C}} = L_{\mathbb{Z}}$, but this choice is not compulsory.

4.2 Partitions

Following [GRS05], our method relies on the choice of a symbolic partition of the index domain. We formalize such a partition as a finite set $\mathcal{P} = (\varphi_p)_{p \in P}$ of properties in $L'_{\mathbb{Z}}$, such that

$$\bigsqcup_{p \in P} \varphi_p = \top_{\mathbb{Z}}, \forall p, p' \in P, (p \neq p') \Rightarrow \varphi_p \sqcap_{\mathbb{Z}} \varphi_{p'} = \perp_{\mathbb{Z}}$$

Elements φ_p of a partition will be called *slices*.

An example of partition is

$$\begin{aligned} \varphi_0 &= (i < 1) & \varphi_1 &= (\ell < 1 \leq i \leq n) \\ \varphi_2 &= (1 \leq \ell < i \leq n) & \varphi_3 &= (1 \leq \ell = i \leq n) \\ \varphi_4 &= (1 \leq i < \ell \leq n) & \varphi_5 &= (1 \leq i \leq n < \ell) \\ \varphi_6 &= (n < i) \end{aligned}$$

A different partition will be considered for each loop in the program. The choice of the partition, for a given loop, is performed automatically from the text of the program according to the following rules:

- (i) the partition of a nested loop should refine the partition of the outer loop;
- (ii) for a loop “for($i := \text{Iexp}$; cond; $++$)”, the partition should distinguish the cases ($\ell < \text{Iexp}$) and ($\ell \geq \text{Iexp}$); conversely, for a loop “for($i := \text{Iexp}$; cond; $--$)”, the partition should distinguish the cases ($\ell > \text{Iexp}$) and ($\ell \leq \text{Iexp}$).
- (iii) for each “ $A[\text{Iexp}]$ ” appearing either in the left part of an assignment, or in a condition, the partition should distinguish between ($\ell < \text{Iexp}$), ($\ell = \text{Iexp}$), and ($\ell > \text{Iexp}$).

In practice, a partition doesn’t have to cover the whole domain of indices: it is enough that it covers all the valuations that are reachable during the program execution; in particular, the cases where ℓ is outside the array bounds are not considered. A preliminary analysis of the indices generally provides a restricted domain for each

$\phi = 0 \leq j \leq i - 1, 2 \leq i \leq n$	
$\varphi_1 = (1 = \ell < j < i)$	$\psi_1 = \top_{\mathcal{C}}$
$\varphi_2 = (1 = j = \ell < i)$	$\psi_2 = \top_{\mathcal{C}}$
$\varphi_3 = (1 = j + 1 = \ell < i)$	$\psi_3 = (a^0 > x)$
$\varphi_4 = (2 \leq \ell < j < i)$	$\psi_4 = (a^0 \geq a^{-1})$
$\varphi_5 = (2 \leq j = \ell < i)$	$\psi_5 = (a^0 \geq a^{-1})$
$\varphi_6 = (2 \leq j + 1 = \ell < i)$	$\psi_6 = (a^0 > x, a^0 \geq a^{-1})$
$\varphi_7 = (1 \leq j + 1 < \ell < i)$	$\psi_7 = (a^0 \geq a^{-1} > x)$
$\varphi_8 = (2 \leq \ell = j + 1 = i)$	$\psi_8 = (a^0 = x)$
$\varphi_9 = (1 \leq j + 1 < \ell = i)$	$\psi_9 = (a^0 \geq a^{-1} > x)$
$\varphi_{10} = (1 \leq j + 1 \leq i, 2 \leq i < \ell)$	$\psi_{10} = \top_{\mathcal{C}}$

Figure 4. Example of abstract value for insertion sort

loop. The left column of Fig. 4 shows a realistic example of partition, where, since it is known that $0 \leq j < i \leq n + 1$, situations like $i \leq 0$ or $j \geq i$ are not considered.

Remarks:

- Although a slice φ_p in a partition is intended to specify an interval for the specific variable ℓ , this interval depends on the valuation \mathcal{I} of indices.
- Moreover, φ_p may involve constraints on \mathcal{I} . We will note $\overline{\varphi_p}$ the property $\exists \ell. \varphi_p$ which summarizes the constraints on indices in the slice φ_p .
- As soon as an equation ($\ell = \text{Iexp}$) appears in the definition of a slice, the slice is a *singleton*: for any \mathcal{I} there is at most one ℓ in the slice. Singleton slices play an important role, because their corresponding slice variables may be dealt with as scalar variables. We note $\text{Single}(\mathcal{P})$ the set of singletons in \mathcal{P} .

Slice compatibility: Let's recall that $\overline{\varphi_p}$ denotes the constraints on indices involved by a slice φ_p . Now, if $\overline{\varphi_p} \cap \overline{\varphi_{p'}} = \perp_{\mathcal{Z}}$, the slices φ_p and $\varphi_{p'}$ induce contradictory constraints, so these slices cannot be both non-empty in a same state. There are said to be *incompatible*. We shall use a notion of compatibility relative to some formula ϕ on indices: two slices φ_p and $\varphi_{p'}$ will be said to be ϕ -compatible, if and only if $\phi \cap \overline{\varphi_p} \cap \overline{\varphi_{p'}} \neq \perp_{\mathcal{Z}}$.

4.3 Abstract Values

As announced in the introduction, our analysis will make use of *slice variables*, possibly *shifted*. The choice of these variables will be explicated later, but we already assume the existence of a finite set $\{a^z\}$ of slice variables: this notation associates uppercase letters (e.g., A, B) for arrays, with lowercase letters (a, b) for slice variables; the exponent z is a relative integer, called the *shift*. For a given valuation \mathcal{I} of index variables and a given slice φ_p , the slice variable a^z represents the subarray $\{A[\ell + z] \mid \varphi_p(\mathcal{I}, \ell)\}$.

Remark: the disappearance of the subscript p , referring to a symbolic slice, on slice variables is *not* a shortcut. As we saw during the analysis of our intuitive example (§2), we are interested in the join of properties linked to different symbolic slices (e.g., weak update). If these properties use the same slice variables, the join operation can be performed in a classical lattice (here $L_{\mathcal{C}}$). In the example, in place of joining $a_1 = b_1$ and $a_2 = b_2$ the join will be between two identical properties $a^0 = b^0$.

Given a partition $(\varphi_p)_{p \in P}$ and a set of slice variables $\{a^z\}$, an *abstract value* Ψ consists of

- a property $\phi \in L_{\mathcal{Z}}$ of index variables
- a tuple $(\psi_p)_{p \in P}$ of properties ($\in L_{\mathcal{C}}$) of slice variables and content variables.

The concretization $\gamma(\Psi)$ of an abstract value Ψ is a set of states ($\gamma(\Psi) \subseteq \text{States}$) defined as follows:

$$(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \gamma(\Psi) \Leftrightarrow \begin{cases} \bullet \phi(\mathcal{I}) \\ \bullet \forall p \in P, \forall \ell \text{ such that } \varphi_p(\mathcal{I}, \ell) \\ \quad \psi_p[A[\ell + z]/a^z](\mathcal{C}, \mathcal{A}) \end{cases}$$

(in the formula above, $\psi_p[A[\ell + z]/a^z](\mathcal{C}, \mathcal{A})$ means that $(\mathcal{C}, \mathcal{A})$ satisfy the formula ψ_p where each variable a^z has been replaced by $A[\ell + z]$).

Example: Fig. 4 shows the abstract value associated with the entry of the inner loop of the insertion sort (Fig. 1.c) at the end of the analysis. For instance, (φ_7, ψ_7) expresses that

$$\forall \ell, (1 \leq j + 1 < \ell < i) \Rightarrow A[\ell] \geq A[\ell - 1] > x.$$

We note $\text{Vars}(\psi_p)$ the set of variables appearing in (i.e., constrained by) ψ_p .

For each partition \mathcal{P} , the set of abstract values forms a lattice $(L_{\mathbb{A}}(\mathcal{P}), \sqsubseteq_{\mathbb{A}}, \sqcap_{\mathbb{A}}, \sqcup_{\mathbb{A}}, \top_{\mathbb{A}}, \perp_{\mathbb{A}})$. The lattice operations and all other necessary operations on abstract values are described in the following section.

5. Operations on Abstract Values

5.1 Normalization

An abstract value $\Psi = (\phi, (\varphi_p)_{p \in P})$ is a complex object which must be kept consistent. For instance, if its concretization is empty, it should be normalized to $\perp_{\mathbb{A}}$. There are other needs for normalization, which are taken into account by the following conditional rewriting rules:

1. (empty slice) As said before, when a slice φ_p is empty, the corresponding ψ_p can be any formula, in particular the strongest one $\perp_{\mathcal{C}}$. Hence the rule:

$$(\phi \cap_{\mathcal{Z}} \overline{\varphi_p} = \perp_{\mathcal{Z}}) \Rightarrow \psi_p \rightarrow \perp_{\mathcal{C}}.$$

2. (unsatisfiable index property) If ϕ is not satisfiable, $\gamma(\Psi)$ is empty:

$$\phi = \perp_{\mathcal{Z}} \Rightarrow \Psi \rightarrow \perp_{\mathbb{A}}.$$

3. (consistency of constraints on scalars) The knowledge about scalar contents variables in each ψ_p should be propagated to all slices ϕ -compatible slices with φ_p . For instance, after the assignment “ $x = A[i]$ ”, in the slice $[i, i]$ we know that x has the same properties than $A[\ell]$. The part of this knowledge that does not depend on ℓ (e.g., $x \geq 0$) must be propagated to slices $[1, i - 1]$ and $[i + 1, n]$. Let $\text{Sc}(\psi_p)$ be the projection of ψ obtained by existential quantification of all slice variables a^z . Then

$$(\overline{\varphi_p} \cap_{\mathcal{Z}} \phi) \sqsubseteq_{\mathcal{Z}} \overline{\varphi_{p'}} \Rightarrow \psi_p \rightarrow \psi_p \cap_{\mathcal{C}} \text{Sc}(\psi_{p'}).$$

Fig. 5 shows an abstract value on partition $\{\varphi_1 = (1 \leq \ell < i \leq n), \varphi_2 = (1 \leq \ell = i \leq n), \varphi_3 = (1 \leq i < \ell \leq n)\}$, followed by the result of the application of the rule. Such an abstract value will appear for instance, in the program

$$i := 1; x := 5; \text{for } i := 1 \text{ to } n \text{ do } x := x + 1$$

at the head of the loop.

4. (consistency of shifts) If some z -shift of a slice φ_p is included inside another slice $\varphi_{p'}$, the constraints on a^{-z} in $\varphi_{p'}$ should apply to a^0 in φ_p . For instance, if we know that $\forall \ell \in [2, n], A[\ell] \geq A[\ell - 1]$, we may conclude that $\forall \ell \in [1, 1], A[\ell] \leq A[\ell + 1]$; that is, for $\varphi_1 = (\ell = 1)$ and $\varphi_2 = (2 \leq \ell \leq n)$, if ψ_2 implies $(a^{-1} \leq a^0)$ then φ_1 should imply $(a^0 \leq a^1)$. Let us note $\varphi_p \oplus z = \varphi_p[\ell - z/\ell]$ and

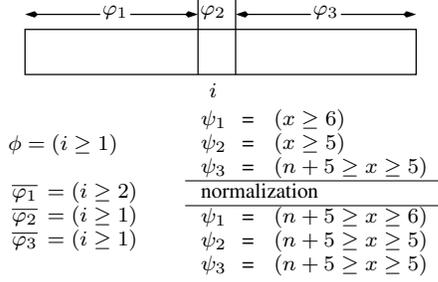


Figure 5. Example for consistency of scalars properties

$\psi_{p'} \boxplus z = \psi_{p'}[a^{y-z}/a^y]$. Then the first rule for consistency of shift is

$$\varphi_p \boxplus z \sqsubseteq_{\mathbb{Z}} \varphi_{p'} \Rightarrow \psi_p \rightarrow (\psi_p \sqcap_{\mathbb{C}} (\psi_{p'} \boxplus z)).$$

In the example above, $\varphi_1 \boxplus 1 = (\ell = 2) \sqsubseteq_{\mathbb{Z}} \varphi_2$, so ψ_2 must be strengthened with $(a^{-1} \leq a^0) \boxplus -1 = (a^0 \leq a^1)$, as desired.

There is a second rule when the reverse inclusion also holds:

$$\varphi_{p'} \sqsubseteq_{\mathbb{Z}} \varphi_p \boxplus z \Rightarrow \psi_{p'} \rightarrow (\psi_{p'} \sqcap_{\mathbb{C}} (\psi_p \boxplus -z)).$$

For instance, this rule will be applied in the normalization phase of Fig. 8, where

$$\begin{aligned} \varphi_1 &= (\ell = i) & \varphi_2 &= (\ell = i + 3) \\ \psi_1 &= (a^0 = x) & \psi_2 &= (a^0 = a^{-3}) \end{aligned}$$

we have $\varphi_2 \boxplus -3 = \varphi_1$, so both ψ_1 and ψ_2 are rewritten:

$$\psi_1 \rightarrow (a^0 = x, a^3 = a^0), \quad \psi_2 \rightarrow (a^0 = a^{-3}, a^{-3} = x).$$

To make the rules readable, we discarded the use of ϕ . However, better results are obtained if the slices appearing in premises of the rules are intersected with ϕ . For instance, consider the partition of Fig. 5 with $\phi = (i = 2)$. Then we have $(\varphi_1 \sqcap_{\mathbb{Z}} \phi) \sqsubseteq_{\mathbb{Z}} ((\varphi_2 \boxplus -1) \sqcap_{\mathbb{Z}} \phi)$ and ψ_1 could be strengthened.

Of course, these rules influence each other. By chance they can be applied step by step, in reverse order w.r.t. the list above (i.e., apply thoroughly the rule (4), then apply thoroughly the rule (3), etc.). Henceforth, we will assume that all abstract values are normalized.

Proposition. Normalization does not change the concretization of formulas: $\Psi \rightarrow \Psi' \Rightarrow \gamma(\Psi) = \gamma(\Psi')$

5.2 Lattice Operations

After normalization, the lattice operations are straightforward. They are defined on abstract values based on the same partition:

$$(\phi, (\varphi_p)_{p \in P}) \sqcap_{\mathbb{A}} (\phi', (\varphi'_p)_{p \in P}) = (\phi \sqcap_{\mathbb{Z}} \phi', (\varphi_p \sqcap_{\mathbb{C}} \varphi'_p)_{p \in P})$$

$$(\phi, (\varphi_p)_{p \in P}) \sqcup_{\mathbb{A}} (\phi', (\varphi'_p)_{p \in P}) = (\phi \sqcup_{\mathbb{Z}} \phi', (\varphi_p \sqcup_{\mathbb{C}} \varphi'_p)_{p \in P})$$

$$\Psi \sqsubseteq_{\mathbb{A}} \Psi' \iff \phi \sqsubseteq_{\mathbb{Z}} \phi' \wedge \forall p \in P, \psi_p \sqsubseteq_{\mathbb{C}} \psi'_p$$

5.3 Change of Partition

When entering a loop, a new index i is introduced, and the partition is refined according to this index. Conversely, when a loop is exited, its index is forgotten, and the partition is simplified accordingly. We need operations for transforming an abstract value Ψ on a partition \mathcal{P} to an abstract value $\Psi' = [\mathcal{P} \rightarrow \mathcal{P}'](\Psi)$ on a partition \mathcal{P}' that refines \mathcal{P} , and the converse operation, noted $[\mathcal{P}' \xrightarrow{i} \mathcal{P}]$ which also forgets about the index i . If $\mathcal{P}' = (\varphi'_p)_{p \in P'}$ is a refinement of

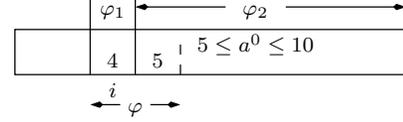


Figure 6. Example for slice property

$\mathcal{P} = (\varphi_p)_{p \in P}$, then, for every $p \in P'$, there exists $f(p) \in P$ such that $\varphi_p \sqsubseteq_{\mathbb{Z}} \varphi_{f(p)}$. Then,

- for $\Psi = (\phi, (\psi_p)_{p \in P}) \in L_{\mathbb{A}}$, we define

$$[\mathcal{P} \rightarrow \mathcal{P}'](\Psi) = (\phi, (\psi_{f(p)})_{p \in P'})$$

- conversely, for $\Psi' = (\phi', (\psi'_p)_{p \in P'})$ we define

$$[\mathcal{P} \xrightarrow{i} \mathcal{P}'](\Psi') = (\exists i. \phi', (\bigsqcup_{p=f(p')} \psi'_p)_{p \in P})$$

5.4 Slice Property

Given a symbolic interval $\varphi \in L'_{\mathbb{Z}}$ and an abstract formula $\Psi \in L_{\mathbb{A}}(\mathcal{P})$, we want to extract the strongest property $\psi^{\Psi}(\varphi)$ implied by Ψ on φ . Of course, if there exists $p \in P$ such that $\varphi = \varphi_p$, $\psi^{\Psi}(\varphi)$ is simply ψ_p . However, the operation is more tricky and interesting when φ intersects several slices in \mathcal{P} . The first idea is to take the least upper bound of the properties of these slices. However we shall see that additional information can be gained from shifted variables.

We define:

$$Inter^{\Psi}(\varphi) = \{p \in P \mid \varphi_p \sqcap_{\mathbb{Z}} \varphi \sqcap_{\mathbb{Z}} \phi \neq \perp_{\mathbb{Z}}\}$$

$$Trans^{\Psi}(\varphi) = \{(p, z) \mid a^z \in Vars(\psi_p) \text{ and } \varphi_p \boxplus z \sqcap_{\mathbb{Z}} \varphi \sqcap_{\mathbb{Z}} \phi \neq \perp_{\mathbb{Z}}\}$$

Then we define

$$\psi^{\Psi}(\varphi) = \bigsqcup_{p \in Inter^{\Psi}(\varphi)} \psi_p \sqcap \bigsqcap_{(q, z) \in Trans^{\Psi}(\varphi)} \psi_q \boxplus -z$$

We explain this complex formula by means of a small example (see Fig. 6): the partition is $\{\varphi_1 = (i = \ell), \varphi_2 = (i < \ell)\}$, the abstract formula $\Psi = (1 \leq i, \{\psi_1 = (a^0 = 4, a^1 = 5), \psi_2 = (5 \leq a^0 \leq 10)\})$. For $\varphi = (i \leq \ell \leq i + 1)$, we have $Inter^{\Psi}(\varphi) = \{1, 2\}$ and $Trans^{\Psi}(\varphi) = \{(1, 1)\}$. We get $\bigsqcup_{p \in Inter^{\Psi}(\varphi)} \psi_p = (4 \leq a^0 \leq 10)$. If we take also into account the translated slices intersected by φ (we have $(\varphi \sqcap_{\mathbb{Z}} \varphi_2) \sqsubseteq_{\mathbb{Z}} \varphi_1 \boxplus 1$) we get the more precise result $\psi = (4 \leq a^0 \leq 5)$.

5.5 Index Progression

Let us consider index incrementing (the decrementing being symmetrical). When an index i is incremented, the subarray represented by a slice φ_p involving i changes. Let us note $\varphi'_p = \varphi_p[i + 1/i]$. Thanks to the operation defined in the previous subsection, we can compute the property $\psi'_p = \psi^{\Psi}(\varphi'_p)$, which is the property to be satisfied within φ_p after the assignment. So, we can define

$$[i + +]_{\mathbb{A}}(\Psi) = \left([i + +]_{\mathbb{Z}}(\phi), \left(\psi^{\Psi}(\varphi_p[i + 1/i]) \right)_{p \in P} \right).$$

For instance, Fig. 7 represents a property with partition $\{\varphi_1 = (1 \leq \ell < i \leq n), \varphi_2 = (1 \leq \ell = i \leq n), \varphi_3 = (1 \leq i < \ell \leq n)\}$, and $\phi = (1 \leq i \leq n)$, $\psi_1 = (0 \leq a^0 \leq 4)$, $\psi_2 = (a^0 = 6)$, $\psi_3 = (a^0 > a^{-1})$. Then, $\varphi'_1 = (1 \leq \ell < i + 1 \leq n)$, $\varphi'_2 = (1 \leq \ell = i + 1 \leq n)$, $\varphi'_3 = (1 \leq i + 1 < \ell \leq n)$. We get $[i + +]_{\mathbb{A}}(\Psi) = (\phi = (2 \leq i \leq n + 1), \psi_1 = \psi^{\Psi}(\varphi'_1) = (0 \leq a^0 \leq 6), \psi_2 = \psi^{\Psi}(\varphi'_2) = (a^0 > a^{-1}), \psi_3 = \psi^{\Psi}(\varphi'_3) = (a^0 > a^{-1}))$.

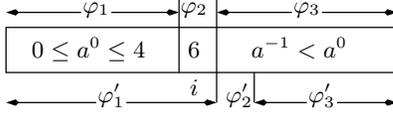


Figure 7. Example for index progression

5.6 Loop Index Initialization

Index initialization “ $i := Iexp$ ” only occurs when entering a loop. So, it is applied just after the refinement of the partition (cf. §5.3), and nothing is known about i before the assignment. So the transformation is just

$$[i := Iexp]_{\mathbb{A}}(\Psi) = (\phi \sqcap (i = Iexp), (\psi_p)_{p \in P}).$$

5.7 Content Assignment

A content assignment of an expression exp , affects several slice formulas, in which the modified slice variable will be different, the update will be applied strongly or weakly, and exp will be translated in terms of slice variables in accordance. First, we address left part issues:

Left part:

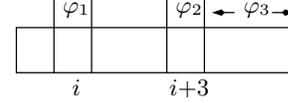
- In case of a scalar assignment $x := exp$, the variable x must be strongly updated in all ψ_p .
- In case of an array cell assignment $A[i + k] := exp$, we know that there are singleton slices φ_s where $(\ell = i + k)$. In the corresponding ψ_s , a^0 must be strongly assigned. We call $Strong_1$ the set of such s . Moreover, for each $s \in Strong_1$, there can be also other slices φ_p , ϕ -compatible with φ_s , such that there is some $a^z \in Vars(\psi_p)$ and $\varphi_p \oplus z$ intersects (in fact, contains) φ_s . We call $Affect(\varphi_s)$ the set of such p . If $p \in Affect(\varphi_s)$, ψ_p refers to a shift variable which may be aliased with the assigned variable: if φ_p is a singleton, a^z must be strongly assigned, otherwise it must be weakly assigned in ψ_p . So, we define

$$\begin{aligned} Strong_1 &= \{s \in P \mid \varphi_s \sqsubseteq_{\mathbb{Z}} (\ell = i + k)\} \\ Affect(\varphi_s) &= \{(p, z) \mid a^z \in Vars(\psi_p), \varphi_s \sqsubseteq_{\mathbb{Z}} \varphi_p \oplus z\} \\ Strong &= \{s' \in Single(\mathcal{P}) \mid \exists s \in Strong_1, \exists z \in \mathbb{Z}, \\ &\quad (s', z) \in Affect(\varphi_s)\} \\ Weak &= \{p \in P \setminus Single(\mathcal{P}) \mid \exists s \in Strong_1, \\ &\quad \exists z \in \mathbb{Z}, (p, z) \in Affect(\varphi_s)\} \end{aligned}$$

Scalar right part: When the right part exp does not contain any array access, the postcondition of the assignment is easy to compute. Let us note ψ'_p the slice formulas of $[Left = exp]_{\mathbb{A}}(\Psi)$. If the left-hand side is scalar, then $\forall p, \psi'_p = [x := exp]_{\mathbb{C}}(\psi_p)$. Otherwise, with the notations defined above:

$$\begin{aligned} \forall s \in Strong_1, \psi'_s &= [a^0 := exp]_{\mathbb{C}}(\psi_s) \\ \forall s \in Strong, (s, z) \in Affect(\varphi_{s'}) &, \psi'_s = [a^z := exp]_{\mathbb{C}}(\psi_s) \\ \forall p \in Weak, (p, z) \in Affect(\varphi_{s'}) &, \psi'_p = \psi_p \sqcup_{\mathbb{C}} [a^z := exp]_{\mathbb{C}}(\psi_p) \end{aligned}$$

Right part with array accesses: Without loss of generality assume the right part of the assignment is an expression “ $B[j + k]$ ” (more complex expressions may be reduced to this case using auxiliary scalar variables and successive assignments). If the left-hand side is a scalar variable x , its new value must be computed for each slice. If it is an array cell, we have to perform weak or strong assignments in some slice property ψ_p to some slice variable a^z , p belonging either to $Strong_1$, or to $Strong$, or to $Weak$. So, we want to get the strongest information about $B[j + k]$ available in some slice φ_p . Two cases occur:



$$\begin{aligned} \psi_1 &= (a^0 = x) \\ \psi_2 &= (a^0 > x) \\ \psi_3 &= (a^0 > x, a^{-1} > x, a^0 \geq a^{-1}) \\ \hline \text{content assignment } A[i + 3] &:= A[i] \\ \hline \psi_1 &= (a^0 = x) \\ \psi_2 &= (a^0 = a^{-3}) \\ \psi_3 &= (a^0 > x, a^{-1} \geq x, a^0 \geq a^{-1}) \\ \hline \text{normalization} \\ \psi_1 &= (a^0 = x, a^0 = a^3, a^3 = x) \\ \psi_2 &= (a^{-3} = x, a^0 = x, a^0 = a^{-3}) \\ \psi_3 &= (a^0 > x, a^{-1} \geq x, a^0 \geq a^{-1}) \end{aligned}$$

Figure 8. Example for content assignment

- either $\varphi_p(\ell)$ implies $j + k = \ell + w$, for some w , in which case, $B[j + k]$ is represented by a shift variable b^w : if $\varphi_p \sqcap_{\mathbb{Z}} (\ell' = j + k)$ implies $\ell' = \ell + w$ for some w , then

$$[A[i] := B[j + k]]_{\mathbb{A}}(\psi_p) = [a^z := b^w]_{\mathbb{C}}(\psi_p)$$

- otherwise, we extract the scalar information about $B[j + k]$ using our “slice property” operation (§5.4), as follows: the formula $\psi = \psi^{\Psi}((\ell = j + k) \sqcap_{\mathbb{Z}} \phi \sqcap_{\mathbb{Z}} \overline{\varphi_p})$ gives for b^0 the best information about $B[j + k]$ for the slice φ_p ($(\ell = j + k)$ is ϕ -compatible with φ_p , so $\psi' = Sc(\psi[y/b^0])$, where y is a fresh scalar variable, gives the wanted information for y . Then the result is obtained in a similar way as for a scalar right part reduced to y :

$$[x := B[j + k]]_{\mathbb{A}}(\psi_p) = \exists y. [x := y]_{\mathbb{C}}(\psi_p \sqcap_{\mathbb{C}} \psi')$$

Let $\varphi' = (\ell = j + k)$; For $s \in Strong_1$,

$$[A[i] := B[j + k]]_{\mathbb{A}}(\psi_s) = \exists y. [a^0 := y]_{\mathbb{C}}(\psi_s \sqcap_{\mathbb{C}} \psi')$$

where $\psi' = Sc(\psi^{\Psi}(\varphi' \sqcap_{\mathbb{Z}} \phi \sqcap_{\mathbb{Z}} \overline{\varphi_s})[y/b^0])$.

For $s \in Strong$, $(s, z) \in Affect(\varphi_{s'})$,

$$[A[i] := B[j + k]]_{\mathbb{A}}(\psi_s) = \exists y. [a^z := y]_{\mathbb{C}}(\psi_s \sqcap_{\mathbb{C}} \psi')$$

where $\psi' = Sc(\psi^{\Psi}(\varphi' \sqcap_{\mathbb{Z}} \phi \sqcap_{\mathbb{Z}} \overline{\varphi_s})[y/b^0])$.

For $p \in Weak$, $(p, z) \in Affect(\varphi_{s'})$,

$$[A[i] := B[j + k]]_{\mathbb{A}}(\psi_p) = \psi_p \sqcup_{\mathbb{C}} [a^z := y]_{\mathbb{C}}(\psi_p \sqcap_{\mathbb{C}} \psi')$$

where $\psi' = Sc(\psi^{\Psi}(\varphi' \sqcap_{\mathbb{Z}} \phi \sqcap_{\mathbb{Z}} \overline{\varphi_p})[y/b^0])$.

Example: On the example shown in Fig 8, we perform the assignment $A[i + 3] := A[i]$. We have $s = 2$, $Affect(\varphi_s) = \{(2, 0), (3, -1)\}$, $Strong_1 = \{2\}$, $Strong = \{2\}$, $Weak = \{3\}$. So, $\psi'_2 = \psi_2[a^0 := a^{-3}]_{\mathbb{C}}(\psi_2) = (a^0 = a^{-3})$. To compute ψ_3 , we have first to compute $\psi' = Sc(\psi^{\Psi}((\ell = i) \sqcap_{\mathbb{Z}} \phi \sqcap_{\mathbb{Z}} \overline{\varphi_3})[y/a^0]) = (a^0 = x)[y/a^0] = (y = x)$, and then $\psi'_3 = \psi_3 \sqcup_{\mathbb{C}} [a^{-1} := y]_{\mathbb{C}}(\psi_3 \sqcap_{\mathbb{C}} \psi') = \psi_3 \sqcup_{\mathbb{C}} [a^{-1} := y]_{\mathbb{C}}(a^0 > x, a^{-1} > x, a^0 \geq a^{-1}, y = x, a^0 > y, a^{-1} > y) = \psi_3 \sqcup_{\mathbb{C}} (a^0 > x, y = x, a^0 > y, a^{-1} = y, a^0 > a^{-1}, a^{-1} = x) = (a^0 > x, a^{-1} \geq x, a^0 \geq a^{-1})$. The rest of the work is done by the normalization. Notice that, in spite of the weak assignment, the property that the slice φ_3 is sorted is preserved.

5.8 Conditional statements

Let us note $Cond(c, \Psi)$ the strengthening of a formula $\Psi = (\phi, (\psi_p)_{p \in P})$ by the knowledge that the condition “ c ” is true. According to the syntax of our simple programs, a condition is either a formula on indices, or a formula on contents (including array cells) or a conjunction or a disjunction of conditions. We assume that an elementary condition only on indices is abstracted by a formula $\phi_c \in L_{\mathbb{Z}}$. Then, $Cond(\phi_c, \Psi) = (\phi \sqcap_{\mathbb{Z}} \phi_c, (\psi_p)_{p \in P})$. For elementary conditions only on contents, let us note $Cond(c, \Psi) = (\phi, (Cond(c, \psi_p)_{p \in P}))$. If c does not involve any array expression, it is supposed to be a formula $\psi_c \in L_C$ on content variables, and $Cond(c, \psi_p) = \psi_p \sqcap_C \psi_c$ for all p . Otherwise, if $A[i+k]$ appears in c , there are singleton slices $\varphi_s \sqsubseteq_{\mathbb{Z}} (\ell = i+k)$, for which a represents $A[i+k]$ in ψ_s . As for the assignment, either other array expressions in c can be expressed as translations of $A[i+k]$, and replaced by some a^z in c , or their properties can be extracted using “slice property” operation and reported in c . Anyway, c is transformed into a formula $\psi_c^s \in L_C$. Then, for slices s such that $\varphi_s \sqsubseteq_{\mathbb{Z}} (\ell = i+k)$, $Cond(c, \psi_s) = \psi_s \sqcap_C \psi_c^s$.

Finally, for conditions mixing conditions on indices and on contents, we use the classical lattice operators:

$$\begin{aligned} Cond(c \text{ or } c', \Psi) &= Cond(c, \Psi) \sqcup_{\mathbb{A}} Cond(c', \Psi) \\ Cond(c \text{ and } c', \Psi) &= Cond(c, \Psi) \sqcap_{\mathbb{A}} Cond(c', \Psi) \end{aligned}$$

5.9 Soundness of Operations

Proposition. All the operations defined so far are sound, i.e., $\forall \Psi, \forall (\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \gamma(\Psi)$,

- $\varphi(\mathcal{I}, \ell) \Rightarrow \psi^{\Psi}(\varphi)[A[\ell+z]/a^z](\mathcal{C}, \mathcal{A})$
- $\llbracket i++ \rrbracket(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \llbracket i++ \rrbracket_{\mathbb{A}}(\Psi)$
- $\llbracket i := Iexp \rrbracket(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \llbracket i := Iexp \rrbracket_{\mathbb{A}}(\Psi)$
- $\llbracket x := exp \rrbracket(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \llbracket x := exp \rrbracket_{\mathbb{A}}(\Psi)$
- $\llbracket A[i+k] := exp \rrbracket(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \llbracket A[i+k] := exp \rrbracket_{\mathbb{A}}(\Psi)$
- if $\llbracket c \rrbracket(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in \gamma(\Psi) = true$, then $(\mathcal{I}, \mathcal{C}, \mathcal{A}) \in Cond(c, \Psi)$.

5.10 Widening

We assume that a widening $\nabla_{\mathbb{Z}}$ is available in $L_{\mathbb{Z}}$. For L_C we note ∇_C either the widening in L_C if any, or the least upper bound \sqcup_C if L_C is of finite depth. Then, a natural definition for a widening operator in $L_{\mathbb{A}}$ is

$$(\phi, (\psi_p)_{p \in P}) \nabla (\phi', (\psi'_p)_{p \in P}) = (\phi \nabla_{\mathbb{Z}} \phi', (\psi_p \nabla_C \psi'_p)_{p \in P}).$$

Now, this operator is not completely satisfactory for the following reason: when it happens that some slice φ_p is empty in Ψ (i.e., $\phi \sqcap_{\mathbb{Z}} \varphi_p = \perp_{\mathbb{Z}}$) and not empty in Ψ' , the corresponding slice property is widened from \perp_C in the result (since $\psi_p = \perp_C$). In that case, it is better to wait for two meaningful iterates before performing the global widening. So, we define $(\phi, (\psi_p)_{p \in P}) \nabla_{\mathbb{A}} (\phi', (\psi'_p)_{p \in P})$ as

$$\begin{cases} (\phi', (\psi'_p)_{p \in P}) & \text{if } \exists p, \phi \sqcap_{\mathbb{Z}} \varphi_p = \perp_{\mathbb{Z}} \\ & \text{and } \phi' \sqcap_{\mathbb{Z}} \varphi_p \neq \perp_{\mathbb{Z}} \\ (\phi \nabla_{\mathbb{Z}} \phi', (\psi_p \nabla_C \psi'_p)_{p \in P}) & \text{otherwise} \end{cases}$$

Proposition. $\nabla_{\mathbb{A}}$ is a widening (i.e., $\Psi \sqcup_{\mathbb{A}} \Psi' \sqsubseteq_{\mathbb{A}} \Psi \nabla_{\mathbb{A}} \Psi'$, and $\nabla_{\mathbb{A}}$ satisfies the chain condition).

5.11 About Shift Variables

The analysis introduces slice and shift variables at different steps: slice variables are created at the creation of new partitions (first entry in a loop); shift variables are introduced during normalization (consistency of shifts) and assignment (array expression in right part). It is important to notice that, even if these variables are created during the analysis, their number remains bounded: they can

only appear during the first traversal of a loop, the next iterations can only remove variables.

6. Implementation and Experiments

The method has been implemented in OCaml from a generic analyzer due to B. Jeannot (see <http://bjeannot.gforge.inria.fr/fixpoint/>). When invoking the prototype, one can choose the lattice for indices (intervals or potentials) or a domain for arrays which is a functor on the abstract domains ($L_{\mathbb{Z}}$ and L_C) assuming a fixed set of functions on these domains. Of course, classical operations on lattices must be available, but also normalization functions, postcondition of assignment and condition, support (set of involved variables), and support extension/reduction. The program to be analyzed is given as a control-flow graph, in a language inspired from the input language of FAST [BLP06].

The array content analysis takes place after some preliminary steps:

1. First, an index analysis is performed, which checks the array accesses w.r.t. array bounds, and computes invariants about index values at each control point of the program.
2. Second, an analysis of live variables is performed to determine the live indexes in each point
3. Finally, a traversal of strongly connected subcomponents (SCSCs) builds the partition attached with each SCSC, according to the strategy described in §4.2.

The results of these three analyses are merged to get, at each control point, the needed actions for changing the partition (cf. §5.3). The results of analysis (1) are ϕ formulas that are used for simplifying the partition. Only live variables are considered for partitioning.

Then the analysis described in this paper is performed. Notice that for examples 1.f and 1.g, which use disequations over array cells in their conditions, we used the abstract domain of d DBM [PH07] as the lattice of properties over contents (L_C). This domain is an extension of the DBM one, handling disequations between pairs of variables.

We give now the results of this analysis on some examples: Insertion sort (Fig. 1.c), a version of the famous “Find” program (Fig. 1.d) used for segmenting arrays in QuickSort, a simple initialization sequence (Fig. 1.e) showing properties mixing indices and contents when the same lattice is used (numerical arrays), a version of the “Sentinel” program (Fig. 1.f), which is a well-known challenge for array bound checking, and a program involving an index assignment: “First not null”, extracted from [FQ02]. Of course, simpler programs like “Array maximum” and “Array copy” (Fig 1) give the announced results, which we don’t detail here.

Results for “Find”: Fig. 1.d shows a version of the famous “Find” program [Hoa71], used for segmenting an array according to its first element, and used at each step of the “QuickSort”. The final results attached by our analyzer to the very end of the program are exactly those expected (i.e., the array is segmented):

- $\phi = (2 \leq i \leq n+1, j = i-1)$
- either $n = 1$, then $i = 2$ and $A[1] = x$
- or $n \geq 2$ and $A[i-1] = x$ and $\forall \ell, 1 \leq \ell < i-1 \Rightarrow A[\ell] < x$ and $\forall \ell, i \leq \ell \leq n \Rightarrow x \leq A[\ell]$

Results for “Insertion sort”: The result of the analysis at the end of the program is as simple as the partition at that point: $\forall \ell, 2 \leq \ell \leq n \Rightarrow A[\ell-1] \leq A[\ell]$.

A more interesting result is after the array assignment in the nested loop ($(*)$ point in Fig.1.c). The situation is: sorting is not terminated ($i \leq n$) and current value in cell j is greater than the key x , so we assign its value to cell $j+1$. We have:

	# vert. × # edg.	# φ_p	# shifts avg (max)	# iter.	time (s)	# norm.	time (%)	norm. avg. time (s) /(# φ_p + #shifts)
array copy	3 × 3	3	0 (0)	5	0.02	30	-	-
seq. init.	3 × 3	4	0.8 (2)	5	0.05	32	54	0.00016
max. search	5 × 6	4	0.8 (2)	5	0.10	58	50	0.00017
sentinel	3 × 3	9	0 (1)	5	0.21	28	22	0.00019
first n. null	6 × 8	13	0 (1)	6	2.25	89	50	0.00102
insert. sort	9 × 11	4-10	4.6 (11)	7	5.38	169	85	0.00153
find	9 × 13	14	6.7 (14)	6	22.87	171	74	0.00453

Figure 9. Performance results

- $\phi = (1 \leq j \leq i - 1 \leq n - 1)$
- either $j = 1$ then $A[1] = A[2] > x$
 - moreover if $i \geq 3$ then $\forall \ell, 3 \leq \ell \leq i, A[\ell] \geq A[\ell - 1] > x$
- or $j \geq 2$, then $\forall \ell, 2 \leq \ell < j, A[\ell] \geq A[\ell - 1]$ and $A[j] > x$ and $A[j + 1] = A[j] \geq A[j - 1]$
 - moreover if $i > j + 1$ then $\forall \ell, j + 2 \leq \ell \leq i, A[\ell] \geq A[\ell - 1] > x$

Results for “Sequence init.”: At the end of the program, we get the expected bound on array contents ($\leq n + 6$):

- $\phi = (1 \leq n = i - 1)$
- $n \geq 1 \Rightarrow A[1] = 7 \leq n + 6$
- moreover if $n \geq 2$ then $\forall \ell, 2 \leq \ell \leq n \Rightarrow 8 \leq A[\ell] \leq n + 6, 7 \leq A[\ell - 1] \leq n + 5, A[\ell] = A[\ell - 1] + 1$

Results for “Sentinel”: A more surprising success is the discovery, in the program of Fig. 1.f, that the index cannot exceed the bound n : the “sentinel” x is either found before or at n . At the end of the program we get:

- $\phi = 1 \leq i \leq n$
- $A[i] = x$ and $\forall \ell, 1 \leq \ell < i \Rightarrow A[\ell] \neq x$

It is an interesting case where a property on contents involves a property on indices.

Results for “First not null”: The program of Fig. 1.g, was given in [FQ02]. The standard method gives poor results. However, if we enrich the partition by distinguishing the singleton ($\ell = s$) we get the expected property at then end of the program:

- $\phi = (1 \leq s \leq i = n + 1)$
- $s \leq n + 1 \Rightarrow (\forall \ell, 1 \leq \ell \leq s - 1 \Rightarrow A[\ell] = 0)$
- moreover if $s \leq n$ then $A[s] \neq 0$

Performances: Experiments were driven on a Core2 Duo 1.6 GHz, with 2Mo of RAM. In Fig. 9, the first part of the table shows, for each example, the size of the control-flow graph (number of vertices and edges), the size of the partition, the (average and maximum) number of shift variables presents in abstract values during the analysis, the number of iterations before convergence, and the total computation time.

If small examples take less than a half of second to be analyzed, analysis time sharply increases when more slices and shift variables are required. In order to confirm this point, the second part of the table shows the number of normalizations performed during the analysis. This number is directly linked to the number of abstract operations done, which is a good measure of the complexity of the program under analysis. But it does not explain the high time of the analysis of “Find”. The table indicates also the amount of time spent doing normalizations, which most often need more than half the total analysis time. Finally, it gives the average of normalization

time divided by the sum of the size of the partition and the number of shift variables. As expected, the cost to maintain the coherence between slices is dominating, and seems to behave exponentially with respect to the number of slices (compare “Sentinel” and “First not null”), and the number of shift variables (compare “Find” and “First not null”).

We know there is room for improvements, particularly focusing on the design of a clever data structure for the partition.

7. Conclusion and Future Work

We have presented a method for automatically discovering properties involving array contents in simple programs. The key idea is to synthesize pointwise relations between array segments. The method is able to deal with problems which are well-known to be difficult. A prototype tool has been implemented and applied to some classical examples, most of which, to our knowledge, were not managed by previously existing methods.

This work deserves to be extended in several directions. Our first task will be to improve our prototype implementation, both concerning the performances and the flexibility; in particular, it should be made fully compatible with the APRON interface (see <http://apron.cri.enscm.fr/library/>), which would allow all compatible lattices to be used. A more clever choice of the partition would reduce the number of slices and variables. An improved widening is also likely to reduce the number of iterations.

Of course, more general programs must be considered: on one hand, non convex slices should be dealt with to take into account more than simple index incrementation/decrementation; on the other hand, multi-dimensional arrays should be considered as well. We could make use of non-convex partitioning of multi-dimensional arrays as proposed in [Mas92]. Also, we are not far from analyzing programs like QuickSort, but our current prototype does not deal with recursive programs.

Concerning our examples, notice that generally we did not fully verify programs: for “array maximum”, we prove that the result is greater than all the array elements, it remains to prove that it belongs to the set of elements; for “insertion sort”, we prove that the result is sorted, it remains to prove that it is a permutation of the initial array, and similarly for “find”. So, it would be nice to design an analysis dealing with the (multi-)sets of array cell contents.

A longer term perspective would be to generalize our abstract values for other uses: an array is a special case of function, so an interesting question is whether, using similar principles, we could design an abstract domain for expressing function properties.

Acknowledgments

We are indebted to Tom Reps and Bertrand Jeannot for helpful discussions, and to the anonymous referees for their constructive comments.

References

- [ACD93] R. Alur, C. Courcoubetis, and D. L. Dill. Model-checking in dense real-time. *Information and Computation*, 104(1):2–34, 1993. Preliminary version appears in the Proc. of 5th LICS, 1990.
- [BBH⁺06] A. Bouajjani, M. Bozga, P. Habermehl, R. Iosif, P. Moro, and T. Vojnar. Programs with lists are counter automata. In *Computer Aided Verification (CAV 2006)*, pages 517–531. LNCS 4144, Springer Verlag, July 2006.
- [BCC⁺03] 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, ACM SIGPLAN SIGSOFT Conference on Programming Language Design and Implementation*, pages 196–207, San Diego (Ca.), June 2003.
- [BHMR07] D. Beyer, T. A. Henzinger, R. Majumdar, and A. Rybalchenko. Path invariants. In J. Ferrante and K. S. McKinley, editors, *PLDI 2007*, pages 300–309. ACM, 2007.
- [BLP06] S. Bardin, J. Leroux, and G. Point. Tool presentation: Fast extended release. In *18th Conf. Computer Aided Verification (CAV'2006)*, pages 63–66, Seattle (Washington), 2006. LNCS 4144, Springer-Verlag.
- [BMS06] A. R. Bradley, Z. Manna, and H. B. Sipma. What's decidable about arrays? In E. A. Emerson and K. S. Namjoshi, editors, *VMCAI 06*, pages 427–442. LNCS 3855, Springer Verlag, 2006.
- [CC76] P. Cousot and R. Cousot. Static determination of dynamic properties of programs. In *2nd Int. Symp. on Programming*. Dunod, Paris, 1976.
- [CC04] R. C. Claris and J. Cortadella. Verification of parametric timed circuits using octahedra. In *Designing correct circuits, DCC'04*, Barcelona, March 2004.
- [CH78] P. Cousot and N. Halbwachs. Automatic discovery of linear restraints among variables of a program. In *5th ACM Symposium on Principles of Programming Languages, POPL'78*, Tucson (Arizona), January 1978.
- [Cou03] P. Cousot. Verification by abstract interpretation. In N. Dershowitz, editor, *Proc. Int. Symp. on Verification – Theory & Practice – Honoring Zohar Manna's 64th Birthday*, pages 243–268, Taormina, Italy, June 29 – July 4 2003. © Springer-Verlag, Berlin, Germany.
- [Dil89] D. L. Dill. Timing assumptions and verification of finite state concurrent systems. In *Workshop on Automatic Verification Methods for Finite State Systems, Grenoble*. LNCS 407, Springer Verlag, June 1989.
- [FQ02] C. Flanagan and S. Qadeer. Predicate abstraction for software verification. In *POPL 2002*, pages 191–202. ACM, 2002.
- [GMD⁺04] D. Gopan, F. Di Maio, N. Dor, T. Reps, and M. Sagiv. Numeric domains with summarized dimensions. In *TACAS'04*, pages 512–529, Barcelona, 2004.
- [GMT08] S. Gulwani, B. McCloskey, and A. Tiwari. Lifting abstract interpreters to quantified logical domains. In G. C. Necula and P. Wadler, editors, *POPL 2008*, pages 235–246. ACM, 2008.
- [Gop07] D. Gopan. *Numeric program analysis techniques with applications to array analysis and library summarization*. PhD thesis, Computer Science Department, University of Wisconsin, Madison, WI, August 2007.
- [GRS05] D. Gopan, T. Reps, and M. Sagiv. A framework for numeric analysis of array operations. In *Proc. of POPL'2005*, pages 338 – 350, Long Beach, CA, 2005.
- [Hoa71] C. A. R. Hoare. Proof of a program: Find. *CACM*, 14(1):39–45, 1971.
- [IHV08] R. Iosif, P. Habermehl, and T. Vojnar. What else is decidable about arrays? In R. Amadio, editor, *FOSSACS 2008*. LNCS, Springer Verlag, 2008.
- [JM07] R. Jhala and K. L. McMillan. Array abstractions from proofs. In W. Damm and H. Hermanns, editors, *CAV 2007*, pages 193–206. LNCS 4590, Springer Verlag, 2007.
- [LB04] S. K. Lahiri and R. E. Bryant. Indexed predicate discovery for unbounded system verification. In R. Alur and D. Peled, editors, *CAV 2004*, pages 135–147. LNCS 3114, Springer Verlag, 2004.
- [LBC03] S. K. Lahiri, R. E. Bryant, and B. Cook. A symbolic approach to predicate abstraction. In W. A. Hunt Jr. and F. Somenzi, editors, *CAV 2003*, pages 141–153. LNCS 2725, Springer Verlag, 2003.
- [Mas92] F. Masdupuy. Array abstractions using semantic analysis of trapezoid congruences. In *ICS '92: Proceedings of the 6th international conference on Supercomputing*, pages 226–235, New York, NY, USA, 1992. ACM.
- [Min01] A. Miné. The octagon abstract domain. In *AST 2001 in WCRE 2001*, IEEE, pages 310–319. IEEE CS Press, October 2001.
- [PH07] M. Péron and N. Halbwachs. An abstract domain extending Difference-Bound Matrices with disequality constraints. In B. Cook and A. Podelski, editors, *8th International Conference on Verification, Model-checking, and Abstract Interpretation, VMCAI'07*, Nice, France, January 2007.