

A Scalable Segmented Decision Tree Abstract Domain

Patrick Cousot^{2,3}, Radhia Cousot^{3,1}, and Laurent Mauborgne^{3,4}

¹ Centre National de la Recherche Scientifique

² Courant Institute of Mathematical Sciences, New York University

³ École Normale Supérieure, Paris

⁴ Instituto Madrileño de Estudios Avanzados, Madrid

Dedicated to the memory of Amir Pnueli

1 Introduction

The key to precision and scalability in all formal methods for static program analysis and verification is the handling of disjunctions arising in relational analyses, the flow-sensitive traversal of conditionals and loops, the context-sensitive inter-procedural calls, the interleaving of concurrent threads, etc. Explicit case enumeration immediately yields to combinatorial explosion. The art of scalable static analysis is therefore to abstract disjunctions to minimize cost while preserving weak forms of disjunctions for expressivity.

Building upon packed binary decision trees to handle disjunction in tests, loops and procedure/function calls and array segmentation to handle disjunctions in array content analysis, we introduce *segmented decision trees* to allow for more expressivity while mastering costs via widenings.

2 Semantic Disjunctions in Abstract Interpretation

The main problem in applying abstract interpretation [2, 5, 6] to static analysis is to abstract a non-computable fixpoint collecting semantics $\mathbf{lfp}_{\perp}^{\sqsubseteq} F$ for a concrete transformer $F \in \mathcal{C} \mapsto \mathcal{C}$, partial order \sqsubseteq , and infimum \perp into an abstract semantics $\mathbf{lfp}_{\perp^{\#}}^{\sqsubseteq^{\#}} F^{\#}$ for an abstract transformer $F^{\#} \in \mathcal{A} \mapsto \mathcal{A}$, abstract order $\sqsubseteq^{\#}$, and abstract infimum $\perp^{\#}$ where the existence of fixpoints is guaranteed by Tarski's theorem [21] on complete lattices or its extension to complete partial orders (cpo). The collecting semantics is the specification of the undecidable properties we want to collect about programs. The abstract semantics is an effective approximation of the collecting semantics. For soundness, $\mathbf{lfp}_{\perp}^{\sqsubseteq} F \sqsubseteq \gamma(\mathbf{lfp}_{\perp^{\#}}^{\sqsubseteq^{\#}} F^{\#})$ where $\gamma \in \mathcal{A} \mapsto \mathcal{C}$ is the concretization function. Particular cases involve a Galois connection $\langle \mathcal{C}, \sqsubseteq \rangle \xleftrightarrow[\alpha]{\gamma} \langle \mathcal{A}, \sqsubseteq^{\#} \rangle$ such that $\forall x \in \mathcal{C} : \forall y \in \mathcal{A} : \alpha(x) \sqsubseteq^{\#} y \iff x \sqsubseteq \gamma(y)$ and the case of completeness requiring $\alpha(\mathbf{lfp}_{\perp}^{\sqsubseteq} F) = \mathbf{lfp}_{\perp^{\#}}^{\sqsubseteq^{\#}} F^{\#}$.

In general the concrete domain $\langle \mathcal{C}, \sqsubseteq, \perp, \sqcup \rangle$ is a complete lattice or cpo and the concrete transformer F is in disjunctive form $F \triangleq \bigsqcup_{i \in \Delta} F_i$ and often

completely distributive ($\forall i \in \Delta : F_i(\bigsqcup_{j \in \Delta'} X_j) = \bigsqcup_{j \in \Delta'} F_i(X_j)$) or continuous (completely distributive on increasing chains).

In that most usual case, the iterative fixpoint computation $X^0 = \perp, \dots, X^{n+1} = F(X^n), \dots, X^\omega = \bigsqcup_{n \geq 0} X^n = \mathbf{lfp}_{\perp}^{\square} F$ is $X^1 = F(\perp) = \bigsqcup_{i \in \Delta} F_i(\perp), X^2 = F(X^1) = \bigsqcup_{i \in \Delta} F_i(\bigsqcup_{j \in \Delta} F_j(\perp)) = \bigsqcup_{i, j \in \Delta^2} F_i \circ F_j(\perp), \dots, X^{n+1} = F(X^n) = \bigsqcup_{i \in \Delta} F_i(\bigsqcup_{i_1, \dots, i_n \in \Delta^n} F_{i_1} \circ \dots \circ F_{i_n}(\perp)) = \bigsqcup_{i_1, \dots, i_n, i_{n+1} \in \Delta^{n+1}} F_{i_1} \circ \dots \circ F_{i_n} \circ F_{i_{n+1}}(\perp), \dots$, so that passing to the limit $\mathbf{lfp}_{\perp}^{\square} F = X^\omega = \bigsqcup_{n \geq 0} X^n = \bigsqcup_{n \geq 0} \bigsqcup_{i_1, \dots, i_n \in \Delta^n} F_{i_1} \circ \dots \circ F_{i_n}(\perp)$. This shows that the disjunctive explosion problem appears in the concrete iterative fixpoint definition.

If the abstraction is a Galois connection, the abstraction preserves existing joins. It follows that $\alpha(\mathbf{lfp}_{\perp}^{\square} F) = \alpha(\bigsqcup_{n \geq 0} \bigsqcup_{i_1, \dots, i_n \in \Delta^n} F_{i_1} \circ \dots \circ F_{i_n}(\perp)) = \bigsqcup_{n \geq 0} \bigsqcup_{i_1, \dots, i_n \in \Delta^n} \alpha(F_{i_1} \circ \dots \circ F_{i_n}(\perp))$ which is most often over approximated as $\bigsqcup_{n \geq 0} \bigsqcup_{i_1, \dots, i_n \in \Delta^n} \alpha \circ F_{i_1} \circ \gamma \circ \alpha \circ F_{i_2} \circ \gamma \circ \dots \circ \alpha \circ F_{i_n} \circ \gamma(\alpha(\perp)) \sqsubseteq^{\sharp} \bigsqcup_{n \geq 0} \bigsqcup_{i_1, \dots, i_n \in \Delta^n} F_{i_1}^{\sharp} \circ F_{i_2}^{\sharp} \circ \dots \circ F_{i_n}^{\sharp}(\perp^{\sharp}) = \mathbf{lfp}_{\perp^{\sharp}}^{\square} F^{\sharp}$ where $\forall i \in \Delta : \alpha \circ F_i \circ \gamma \sqsubseteq^{\sharp} F_i^{\sharp}, F^{\sharp} \triangleq \bigsqcup_{i \in \Delta} F_i^{\sharp}$ and $\perp^{\sharp} \triangleq \alpha(\perp)$. This shows that the disjunctive explosion problem does also exist in the abstract.

The situation is even worst in absence of best abstraction, that is of a Galois connection, since the concrete transformers F_i may have many, possibly non-comparable, abstractions F_i^{\sharp} . In absence of minimal abstractions (as shown by the abstraction of a disk by polyhedra [12]), infinitely many potential abstractions may exist. Choosing which abstraction should better be used during the analysis is another source of potential combinatorial explosion.

3 Handling Disjunctions in Abstract Interpretation

Contrary to purely enumerative or symbolic encodings of program properties, abstract interpretation offers solutions to the combinatorial explosion of disjunctions so as to minimize computational costs. The key idea is to abstract away irrelevant properties of the collecting semantics.

The abstract domain $\langle \mathcal{A}, \sqsubseteq^{\sharp}, \perp^{\sharp}, \perp^{\sharp} \rangle$ can be chosen as finite (e.g. predicate abstraction [3, 13]) or better of finite height (e.g. constant propagation [15]) to bound n in $\mathbf{lfp}_{\perp^{\sharp}}^{\square} F^{\sharp} = \bigsqcup_{n \geq 0} F^{\sharp n}(\perp^{\sharp})$.

However this solution has been shown to have intrinsic limitations [8] that can be eliminated thanks to infinite abstract domains not satisfying the ascending chain condition together with widenings ∇ and narrowings Δ [2, 4, 5] (including the common practice of including the widening in the transformers $F_i^{\sharp}, i \in \Delta$ mentioned in [9], by choosing $\lambda X \cdot X \nabla (\alpha \circ F_i \circ \gamma(X)) \sqsubseteq^{\sharp} F_i^{\sharp}, i \in \Delta$, which is not such a good idea since it precludes the later use of a narrowing).

Moreover, in absence of a best abstraction, that is of a Galois connection, a choice is usually made among the possible non-comparable abstractions F_i^{\sharp} of the concrete transformers F_i to minimize costs [7].

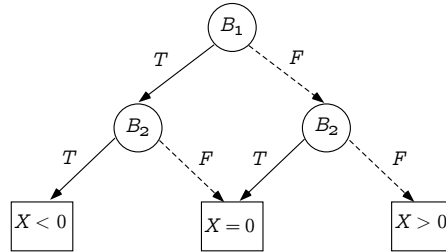
The objective of minimizing computational costs is antagonist to that of precision necessary for making proofs, so that a lot of work in abstract interpretation is on the design of abstract domains offering cost effective yet precise enough ways of abstracting infinite disjunctions.

In this line of work, we consider the examples of binary decision trees in Sect. 4 and segmented arrays in Sect. 6 which common generalization is segmented decision trees as introduced in Sect. 7, a generalization which is more precise while remaining scalable.

4 Binary Decision Trees

Inspired by the use of Binary Decision Diagrams in abstract interpretation [17, 18], binary decision trees can express disjunctive properties depending on the values of binary variables with opportunistic sharing at the leaves [16, 19]. They are an instance of the reduced cardinal power of abstract domains [6, Sect. 10.2] mapping the values of boolean variables (represented in decision nodes) to an abstraction of the other variables (represented in the leaf nodes) for the values of the boolean variables along the path leading to the leaf.

Example 1. The following binary decision tree



written $\llbracket B_1 : \llbracket B_2 : (X < 0), (X = 0) \rrbracket, \llbracket B_2 : (X = 0), (X > 0) \rrbracket \rrbracket$
 encodes
 $(B_1 \wedge B_2 \wedge X < 0) \vee ((B_1 \wedge \neg B_2) \vee (\neg B_1 \wedge B_2)) \wedge X = 0 \vee (\neg B_1 \wedge \neg B_2 \wedge X > 0)$.
 □

The parenthesized representation of trees uses (\dots) for leaves and $\llbracket x : \dots \rrbracket$ for left to right decision nodes on variable x .

Example 2. In the following example, ASTRÉE [10] discovers a relation between the boolean variable B and the unsigned integer variable X.

```
% cat -n decisiontree.c
1 typedef enum {FALSE = 0, TRUE = 1} BOOLEAN;
2 BOOLEAN B;
3 void main () {
4   unsigned int X, Y;
5   while (TRUE) {
6     __ASTREE_analysis_log();
```

```

7     B = (X == 0);
8     if (!B) {
9         Y = 1 / X;
10    };
11 };
12 }
astree --exec-fn main --print-packs decisiontree.c
...
<boolean relations:decisiontree.c@9@5=
  if B then <integers (intv+cong+bitfield+set): X in {0} >
  else <integers (intv+cong+bitfield+set): X in [1, 4294967295] >
  >
...
%
```

At line 9, the relation between B and X is expressed by a binary decision tree with B as only decision node and X in the leaf nodes. This binary decision tree states that if B is TRUE then X is zero else X is positive. Since B is checked to be false there is no division by zero at line 9. \square

5 Variable Packing

Relational abstractions such as octagons [20], polyhedra [12] or binary decision trees in Sect. 4 can express relations between values of variables hence complex disjunctions. However their cost may grow polynomially or even exponentially in size and time. Even for abstract domains such as the octagon domain which has a very light cubic cost compared to many relational domains (such as exponential costs in the number of variables for polyhedra [12]), this would still be too high for very large embedded programs with tens of thousands variables as found in aeronautics.

Variable packing is a way to control memory and computation costs by limiting the number of variables that can be related at a given program point. The decision of which variables can be considered in such abstract relations at each program control point can be taken during the analysis e.g. by widenings (deciding which variables to eliminate by projection from the relation). Variable packing is an even more radical solution which consists in controlling costs a priori, by statically restricting the variables which can be considered in such abstract relations at each program control point. Such a more radical solution is necessary when the cost of even one iteration with all variables in the relation is prohibitive.

The idea is to make small packs of variables that are related while no attempt is made to relate variables appearing in different packs. A simple and cheap pre-analysis groups variables that are interdependent (used together in expressions, loop indices, etc.) in a way that can be represented by the relation. Relations are established among variables associated to each pack, but no relation is kept between variables of distinct packs. The cost thus becomes linear, as it is linear in the number of packs (which is linear in the code size, and so, in the number

of variables) and, e.g. for octagons, cubic in the size of packs (which depends on the size of the considered packs, and is a constant).

Example 3. In *Ex. 2*, ASTRÉE has packed the boolean variable B and the unsigned integer variable X together (but the pack does not contain variable Y). The abstract domain used in the leaves is the reduced product [6] of several abstract domains which can be specified as an option of the analysis and by default are the following

```
/* Domains: ... Packed (Boolean relations (based on Absolute value
equality relations, and Symbolic constant propagation (max_depth=
20), and Linearization, and Integer intervals, and Integer
congruences, and Integer finite sets, and Integer bitfields, and
Float intervals)), and ... */
...
Boolean relations :
List of packs
  decisiontree.c@9@5 { B X }
...
<boolean relations:decisiontree.c@9@5=
if B then <integers (intv+cong+bitfield+set): X in {0} >
else <integers (intv+cong+bitfield+set): X in [1, 4294967295] >
>
...

```

The output in *Ex. 2* only prints informative results, specifically intervals [4, 5], simple congruences [14], bit fields (recording an invariant bit pattern in the binary representation of the variable, if any) and sets of small values (recording the set of possible values of the variable in a small range of values near zero) for that example. \square

Candidates for packing in a binary decision tree are the boolean variables to which a boolean expression is assigned or which are involved in a test as well as the variables which depend directly or indirectly on such a boolean variable, with a maximum number of boolean variables which can be set by an option `--max-bool-var` (3 by default). The option `--print-packs` allows printing packs of variables determined by the pre-analysis for the binary decision trees. In case of insufficient precision, ASTRÉE can be asked to create binary decision tree packs containing given variables by the `_ASTREE_boolean_pack` directive inserted in the program code. Of course putting all boolean program variables in a single pack would certainly achieve great precision but is also the best way to get a combinatorial explosion.

6 Array Segmentation

Array segmentation was introduced in [11] to abstract array content (as opposed to array bounds [4]). The array is divided into consecutive segments. The content of each segment is abstracted uniformly but different segments can have different

abstractions. Starting from a single initial segment with uninitialized content, segments are split and filled by assignments to arrays elements and joined when merging control flows. A widening may be necessary to merge segments so as to avoid the degenerescence to the case of one element per segment.

The bounds of the segments are specified by a set of side-effect free expressions which all have the same concrete value (maybe unknown in the abstract). Segments can be empty thus allowing an array segmentation to encode a disjunction of cases in a way that avoids the explosion of cases.

Example 4. The segmentation $\{0\} \top \{1\} 0\{i\} ? \top \{n\}$ of an array A states that $0 < 1 \leq i < n$, that the values of the array elements $A[0]$, $A[i]$, $A[i+1]$, \dots , $A[n-1]$ are all unknown, while the values of $A[1]$, $A[2]$, \dots , $A[i-1]$, if any, are all initialized to zero. It is possible that $i = 1$ in which case the segment $A[1]$, \dots , $A[i-1]$ is empty, and $i < n$ so that the segment $A[i]$, \dots , $A[n-1]$ cannot be empty (it contains at least one element).

The segmentation $\{0\} \top \{1\} 0\{i\} \top \{n\}$ is similar but for the fact that $i > 1$ so that the segment $A[1]$, \dots , $A[i-1]$ cannot be empty. So $\{0\} \top \{1\} 0\{i\} ? \top \{n\}$ is a compact encoding of the disjunctive information $(\{0\} \top \{1, i\} \top \{n\}) \vee (\{0\} \top \{1\} 0\{i\} \top \{n\})$ distinguishing the first case $i = 1$ from the second $i > 1$.

Similarly $\{0\} \top \{1\} 0\{i\} ? \top \{n\} ?$ is the disjunction $(\{0\} \top \{1, i, n\}) \vee (\{0\} \top \{1, i\} \top \{n\}) \vee (\{0\} \top \{1\} 0\{i, n\}) \vee (\{0\} \top \{1\} 0\{i\} \top \{n\})$. Of course, expressivity is limited since it is not possible to express that either $i = 1$ or $i = n$ but not both (although this might be separately expressible by the abstraction of the simple variables i and n). \square

Note that there are no holes in the segmentation since any such hole is just a segment which content is unknown.

An enriched semantics of arrays is used viewing an array value as the pair of an index and the value of the corresponding array element. In this way the uniform abstraction used in a segment can relate array values within each segment to their index included between the segment bounds.

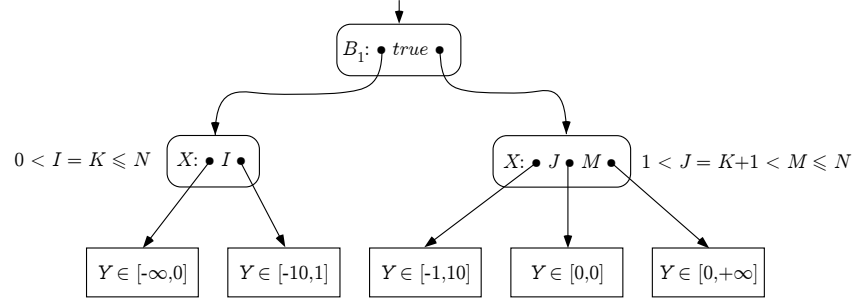
The array segmentation is implemented as a functor which means that the abstract domains representing sets of expressions and array index-elements abstractions should be passed as parameters to the functor to create a particular instance. The advantage of this approach is that the abstract domain parameters can be changed without having to rewrite the static analyzer.

Moreover the static analyzer can optionally perform a reduced product [6] between an instance of the array segmentation functor and the abstract domains that are used for variables appearing in the expressions of the segment bounds. It follows that the array segmentation takes into account both the operations on arrays (array element access and assignment) and operations on variables related to array indexes.

7 Segmented Decision Trees

Segmented decision trees are decision trees where the choices are made on the values of variables according to ranges specified by a symbolic segmentation.

Example 5. The segmented decision tree (where *false* < *true* for booleans)



can be written in the parenthesized form

$$\begin{aligned} & \llbracket B_1 : \llbracket X \{0 < I = K \leq N\} : (Y \in [-\infty, 0]) I (Y \in [-10, 1]) \rrbracket \rrbracket \text{true} \\ & \llbracket X \{1 < J = K + 1 < M \leq N\} : \\ & \quad (Y \in [-1, 10]) J (Y \in [0, 0]) M (Y \in [0, +\infty]) \rrbracket \rrbracket \end{aligned}$$

This segmented decision tree encodes the fact that if B_1 is false (i.e. $B_1 < \text{true}$) then if $X < I$ then Y is non-positive while if $X \geq I$ then $-10 \leq Y \leq 1$. Similarly, if B_1 is true (i.e. $B_1 \geq \text{true}$) then either $X < J$ and $-1 \leq Y \leq 10$, or $J \leq X < M$ and Y is null, or $X > M$ and Y is non-negative. So the leaf nodes specify abstract properties of Y while the decision nodes on B_1 and X specifying conditions for these properties to hold. Attached to each decision node, is a side relation on expressions that holds in the concrete under the condition that this node is reached. For example $(B_1 \wedge 1 < J = K + 1 < M \leq N) \vee (\neg B_1 \wedge 0 < I = K \leq N)$. These expressions are usually in a restricted normal form. In this example the normal form of expressions is an integer constant, a variable, or a variable plus an integer constant and the side relations are those expressible with the octagon abstract domain [20]. The segment bounds are any representative of the equivalent class of expressions which have equal concrete values (so that we could have chosen K for I and $K + 1$ for J). The abstract domain of side relations is assumed to be expressive enough to maintain such equality information between expressions in normal form (i.e. $I = K$ and $J = K + 1$). \square

As for boolean decision trees, an ordering is imposed on all decision variables. That allows binary operations on decisions trees to operate on the same variable⁵. But unlike binary decision trees, the number of choices for a given variable is not bounded *a priori* and the choices may be on different criteria (the bounds in the symbolic segmentations) at each node, even if they have the same decision variables.

As for simple array segmentation, the ordering of the bounds of each segment describes an order on expressions. That means that segments could describe

⁵ In addition it may allow to eliminate the nodes with only one child, an optimization we will not apply in this paper.

two kinds of informations: a series of tests deciding what choice to make and a pre-order on some expressions. Unlike for array segmentation, that information would now be relative to the choices above a node in a tree. So we decided to separate the two notions, such that at each node, we have a decision variable, a pre-order and a segmentation that respects the pre-order (instead of prescribing it) and leads to subtrees. A direct consequence is that segments will be much more simple, as that removes the necessity of a ? tag for emptiness or of the lowest and highest bounds of the segmentation. Also, it allows the use of single expressions as bounds instead of sets of expressions.

Separating the pre-order information from the decision process allows us to use much more precise domains for the pre-order and leads to more precise unifications of segmentations. But storing a pre-order on all possible expressions that can be used at a given node in the tree might lead to very expensive representations. So we chose instead to rely on reduction with other abstract domains for pre-order information valid at every node in the tree and store at each node the ordering information that we can add to what holds for its father. In that way, if the abstraction of an instruction implies an ordering that is not relevant to the tree, it will not increase the complexity of the tree, but further operations might still use that information from reduction with the other abstract domains. In the drawings and examples, we will represent that information as a set of inequations between expressions, but in practice, depending on the canonical expressions we use, we can implement it in a more efficient way, using for example small octagons if the canonical expressions are of the form $\pm X + c$ with X a variable and c a constant. Also, in order to simplify the presentation, in the schemata, we put all global pre-order information at the root ⁶.

Definition 1. A segmented decision tree $t \in \mathbb{T}((\mathbb{D}, <_{\mathbb{D}}), \mathbb{E}, D_c, D_\ell)$ over decision variables in the totally ordered set $(\mathbb{D}, <_{\mathbb{D}})$, canonical expressions in \mathbb{E} , ordering abstract domain D_c (with concretization γ_c) and leaf abstract domain D_ℓ (with concretization γ_ℓ) is either $\langle p \rangle$ with p an element of D_ℓ or $\llbracket \mathbf{x} \{C\} : t_0 b_1 t_1 \dots b_n t_n \rrbracket$ such that \mathbf{x} is the smallest variable in \mathbb{D} , each b_i ($1 \leq i \leq n$) is an element of \mathbb{E} , C is an element of D_c and each $t_i \in \mathbb{T}((\mathbb{D} \setminus \{x\}, <_{\mathbb{D}}), \mathbb{E}, D_c, D_\ell)$ ($0 \leq i \leq n$). \square

To define the concretization of a segmented decision tree, we will write ρ for concrete environments assigning concrete values $\rho(\mathbf{x})$ to variables \mathbf{x} and $\llbracket e \rrbracket \rho$ for the concrete value of the expression e in the concrete environment ρ . The concretization of a segmented decision tree reduced to a leaf is

$$\gamma_t(\langle p \rangle) \triangleq \gamma_\ell(p)$$

and the concretisation of a segmented decision tree rooted at a decision node is

⁶ For clarity, some redundancy is sometimes preserved in segmented decision trees while, for brevity, some repetitive information is omitted in fixpoint computations.

$$\begin{aligned}
\gamma_t.(\llbracket \mathbf{x} \{C\} : t_0 b_1 t_1 \dots b_n t_n \rrbracket) &\triangleq \\
\{\rho \in \gamma_c(C) \mid \forall i \in [1, n) : \llbracket b_i \rrbracket \rho \leq \llbracket b_{i+1} \rrbracket \rho \wedge \\
&(n = 0 \vee \rho(\mathbf{x}) < \llbracket b_1 \rrbracket \rho) \implies \rho \in \gamma_t(t_0) \wedge \\
&\forall i \in [1, n) : (\llbracket b_i \rrbracket \rho \leq \rho(\mathbf{x}) < \llbracket b_{i+1} \rrbracket \rho) \implies \rho \in \gamma_t(t_i) \wedge \\
&(n > 0 \wedge \rho(\mathbf{x}) \geq \llbracket b_n \rrbracket \rho) \implies \rho \in \gamma_t(t_n)\}
\end{aligned}$$

We introduce also the notation $\perp_{\mathbb{D}}$ for the decision tree in $\mathbb{T}((\mathbb{D}, <_{\mathbb{D}}), \mathbb{E}, D_c, D_\ell)$ such that each node is of the form $\llbracket \mathbf{x} \{\top_C\} : t \rrbracket$ and the only leaf is $\llbracket \perp_\ell \rrbracket$, where \top_C is the top element of D_c and \perp_ℓ is the bottom element of D_ℓ . When \mathbb{D} is clear from the context, we simply write \perp .

7.1 Segmented Decision Tree Abstract Functor

The segmented decision tree abstract functor $\mathbb{T}((\mathbb{D}, <_{\mathbb{D}}), \mathbb{E}, D_c, D_\ell)$ is a parameterized abstract domain taking as a parameter a totally ordered set $(\mathbb{D}, <_{\mathbb{D}})$ of decision variables, a set \mathbb{E} of canonical expressions, an ordering abstract domain D_c and a leaf abstract domain D_ℓ for the leaves.

The abstract domain D_ℓ for the leaves is usually the reduced product [6] of several abstract domains, as was the case for binary decision trees in Sect. 4. The list of abstract domains appearing in this reduced product at the leaves is assumed to be an option of the static analyzer constructor. Therefore this option specifies a particular instance of the segmented decision tree abstract functor used to build a particular instance of the static analyzer for that option. The advantage of this modular approach is that the static analyzer can be changed by changing the options, without any re-programming.

The maximal height of the segmented decision trees is a parameter of the static analysis which can therefore be changed before each run of the static analyzer. A variable packing pre-analysis is used to determine which variables \mathbb{D} are chosen to appear in the decision and leaf nodes. The number of variables in the decision nodes is bounded by this maximal height. Following [11], the choice of which expressions $b_1, \dots, b_n \in \mathbb{E}$, $n \geq 0$ do appear in decision nodes is made during the static analysis.

7.2 Reduction of an Abstract Property by a Segmented Decision Tree

Given a segmented decision tree t and an abstract property $p \in D$ of the variables in abstract domain $\langle D, \sqsubseteq, \perp, \sqcup, \sqcap \rangle$ with concretization γ , $t \sqcap_D p$ is the abstraction of the conjunction $\gamma_t(t) \cap \gamma(p)$ in the abstract domain D . It is the intersection of p with the join of the abstract properties obtained along paths of t feasible for p .

Example 6. In *Ex. 5*, the hypotheses that B_1 is *true* and $X < M$ imply that $Y \in [-1, 10] \sqcup [0, 0] = [-1, 10]$. The implied condition collects information along

the path and at the leaves and is therefore $B_1 \wedge (X < M) \wedge (1 < J = K + 1 < M \leq N) \wedge Y \in [-1, 10]$. \square

The operation $t \sqcap_D p$ is used in the definition of the reduced product of the abstract domain $\langle D, \sqsubseteq \rangle$ by the abstract domain of segmented decision trees.

A path of t is feasible for p if it corresponds, at each level in the tree, to a segment of the node which, according to p and the conjunction of side conditions collected from the root to this node, is not empty. For each decision variable, the conjunction of the hypothesis p and the collected side conditions is used to determine to which segments the value of the variable does belong. The information available on this variable is thus the join of the information available at the leaves for each segment plus the fact that the variable value is between the extreme bounds of these segments along this path. More formally (\sqsubseteq is the abstract implication, $(\text{true} ? a : b) \triangleq a$, $(\text{false} ? a : b) \triangleq b$ is the conditional, $D(p)$ is the best approximation of $p \in D_c \cup D_\ell$ in D (in absence of best abstraction, an over-approximation must be used)).

$$\begin{aligned}
\perp_{\mathbb{D}} \sqcap_D p &\triangleq \perp \\
t \sqcap_D \perp &\triangleq \perp \\
(p') \sqcap_D p &\triangleq D(p') \sqcap p \\
\llbracket \mathbf{x} \{C\} : t_0 b_1 t_1 \dots b_n t_n \rrbracket \sqcap_D p &\triangleq t_0 \sqcap_D (p \sqcap D(C) \sqcap (n = 0 ? \top : D(\mathbf{x} < b_1))) \\
&\sqcup \bigsqcup_{i=1}^{n-1} t_i \sqcap_D (p \sqcap D(C) \sqcap D(b_i \leq \mathbf{x} < b_{i+1})) \\
&\sqcup (n > 0 ? t_n \sqcap_D (p \sqcap D(C) \sqcap D(\mathbf{x} \geq b_n)) : \perp)
\end{aligned}$$

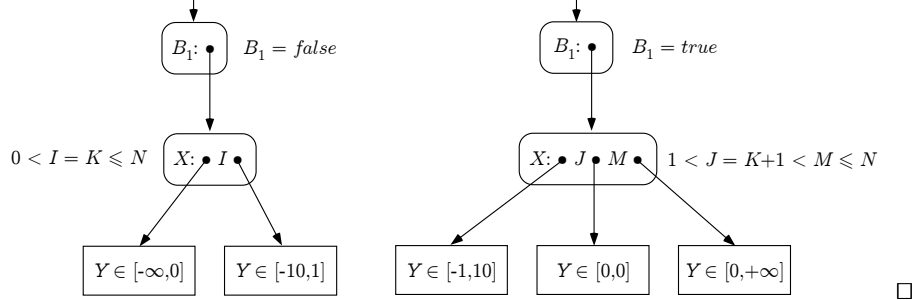
Observe that $p \sqsubseteq q$ implies $\gamma(p) \subseteq \gamma(q)$ whereas $p \not\sqsubseteq q$ does not implies that $\gamma(p) \not\subseteq \gamma(q)$ whereas the sufficient condition $p \sqcap q = \perp$ implies $\gamma(p) \cap \gamma(q) = \emptyset$ and so $\gamma(p) \not\subseteq \gamma(q)$.

7.3 Reduction of a Segmented Decision Tree by an Abstract Property, Tests

In a test, all paths that are feasible in the segmented decision tree are preserved while all the paths that, for sure, can never be followed according to the tested condition are disregarded.

Example 7. Assume that in *Ex. 5*, the test is on B_1 . On the *true* branch of the test, the *false* subtree is disregarded while on the *false* branch of the test the

true subtree is disregarded. We get:



Formally, we define the restriction $t \sqcap_t p$ of the segmented decision tree t by condition $p \in D$ as

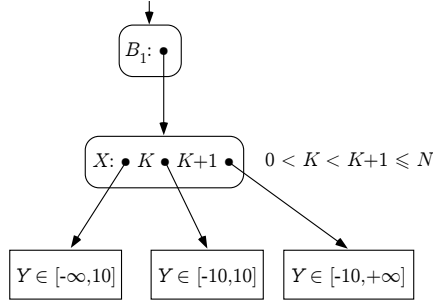
$$\begin{aligned}
\perp_{\mathbb{D}} \sqcap_t p &\triangleq \perp_{\mathbb{D}} \\
t \sqcap_t \perp &\triangleq \perp_{\mathbb{D}} \\
\langle p' \rangle \sqcap_t p &\triangleq \langle p' \sqcap_{D_\ell} D_\ell(p) \rangle \\
\llbracket \mathbf{x} \{C\} : t_0 b_1 t_1 \dots b_n t_n \rrbracket \sqcap_t p &\triangleq \\
&\text{let } t'_0 \triangleq t_0 \sqcap_t (p \sqcap D(C) \sqcap (n = 0 ? \top : D(\mathbf{x} < b_1))) : \\
&\text{and for } i = 1, \dots, n-1 : t'_i \triangleq t_i \sqcap_t (p \sqcap D(C) \sqcap D(b_i \leq \mathbf{x} < b_{i+1})), \\
&\text{and if } n > 0, t'_n \triangleq t_n \sqcap_t (p \sqcap D(C) \sqcap D(\mathbf{x} \geq b_n)) \\
&\text{in } ((\forall i \in [0, n] : t'_i = \perp_{\mathbb{D}}) ? \perp_{\mathbb{D}} : \\
&\quad \llbracket \mathbf{x} \{\text{relax}_C(C, p)\} : L'_i b_{i+1} t'_{i+1} \dots b_m t'_m \rrbracket)
\end{aligned}$$

Note that the information of p that changes the ordering is kept global. On the contrary, it is possible to relax the incremental information on pre-order at each node: each constraint implied by both C and p can safely be removed from C , leading to a more compact representation. The choice of computing that relaxation depends on the domain for C . This choice is noted $\text{relax}_C(C, p)$ in the algorithm. In addition, we can perform another optimization when a node is left with only one child. In that case, we can join the pre-order of the node and its child and put top as the incremental pre-order for the child.

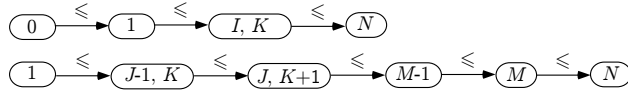
7.4 Segments Unification, Tree Merges and Binary Operations

When performing a binary operation on two decision trees (joins, widening, ordering...), we have to go through the two trees at the same time. Because the variables are ordered, for each pair of subtree during the traversal, the root is on the same decision variable, but the segments may have no bound in common. So, in order to push the binary operation to subtrees, we need to find a common refinement for the two segments (as given by the refinement algorithm in [11]), which we call segments unification.

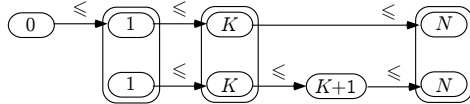
Example 8. Consider the random assignment $B_1 = ?$ to the boolean variable B_1 in the context of *Ex. 5*. The subtrees of the two segments of B_1 must be merged, as follows



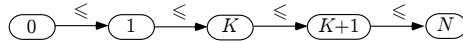
The pre-orders on the decision variable X in *Ex. 5* involve the bounds of the decision variable and the equivalence classes of the expressions appearing in the segmentation.



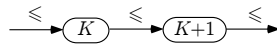
The union of these pre-orders eliminates the variables I , J , and M since they are not comparable in both pre-orders. For example the first pre-order may correspond to a program context where $I = K$ but this might not hold in the program context corresponding to the second pre-order. However although K and N might have different values in these two program contexts, the relation $K \leq N$ is valid in both program contexts and so is preserved in the union of the pre-orders.



This union contains only one maximal chain



which yields a relation between (classes of equal) expressions which is valid in both pre-orders and can therefore be attached to the merged node for X . The segmentation for X in the merged tree is the subchain obtained by considering classes of expressions with representatives appearing in either of the original segmentations (that is $K = I$ and $K + 1 = J$ while 0 , 1 and N did not appear).



- The subtree in the refined first segment $X < K$ is the merge of the subtrees of the corresponding segment $Y \in [-\infty, 0]$ on the left ($X < I = K$) and $Y \in [-1, 10]$ on the right ($X < J = K + 1$), that is $Y \in [-\infty, 0] \sqcup [-1, 10] = [-\infty, 10]$.

- The subtree in the refined second segment $K \leq X < K + 1$ is the merge of the subtrees of the corresponding segments on the left ($Y \in [-10, 1]$ when $X \geq I = K$) and on the right ($Y \in [-1, 10]$ when $X < J = K + 1$), that is $Y \in [-10, 1] \sqcup [-1, 10] = [-10, 10]$.
- The subtree in the refined third and last segment $X \geq K + 1$ is the merge of the subtrees of the corresponding segments on the left ($Y \in [-10, 1]$ for $X \geq I = K$) and on the right ($Y \in [0, 0] \sqcup [0, +\infty]$ for $J = K + 1 \leq X < M$ or $M \leq X$), that is $Y \in [-10, 1] \sqcup [0, 0] \sqcup [0, +\infty] = [-10, +\infty]$. \square

In contrast with simple array segmentation, we may have richer ordering informations, and it will be useful to provide precise segments unification. Given two nodes $\llbracket \mathbf{x} \{C^0\} : t_0^0 b_1^0 \dots b_{n_0}^0 t_{n_0}^0 \rrbracket$ and $\llbracket \mathbf{x} \{C^1\} : t_0^1 b_1^1 \dots b_{n_1}^1 t_{n_1}^1 \rrbracket$, assuming that we collected all the preordering informations in C^0 and C^1 , we must compute two nodes $\llbracket \mathbf{x} \{C^0\} : t_0^0 b_1 \dots b_n t_n^0 \rrbracket$ and $\llbracket \mathbf{x} \{C^1\} : t_0^1 b_1 \dots b_n t_n^1 \rrbracket$ which do share the same bounds and are sound (and precise) over-approximations of the input nodes. For these nodes to be valid, the bounds $b_1 \dots b_n$ must respect the ordering in both C^0 and C^1 . So, the first step is to compute $C^0 \cup_C C^1$. Then the resulting nodes will be more precise if the new bounds can all be compared to the previous bounds. That is, for all $k \in \{0, 1\}$, for all $i \leq n^k$ and all $j \leq n$, either $\forall \rho \in \gamma_c(C^k) \llbracket b_i^k \rrbracket \rho \leq \llbracket b_j \rrbracket \rho$ or $\forall \rho \in \gamma_c(C^k) \llbracket b_i^k \rrbracket \rho \geq \llbracket b_j \rrbracket \rho$. Then solving the segment unification problem consists in finding a chain in the pre-order defined by $C^0 \cup_C C^1$ such that all elements of the chain respect that property. In general there is no best chain, but the longer the chain the better, and it is certainly best to maximize the number of bounds b_i^k such that there is a b_j such that $\forall \rho \in \gamma_c(C^k) \llbracket b_i^k \rrbracket \rho = \llbracket b_j \rrbracket \rho$, meaning that a maximum number of bounds are preserved.

The algorithm to find the bounds $b_1 \dots b_n$ will thus starts by building a graph whose vertices are elements of each segments and the expressions in E that are equal to a bound in each segment according to the C^k 's (including extrema for the decision variable, if any). Then on that graph, we merge the strongly connected components. Each vertex of this graph can be colored by a couple of bounds (b_i^0, b_j^1) or by one bound of a segment. Then we find the path in the graph with maximal number of colors. A couple of optimizations can be used to compute that path, and any path is a correct, although maybe imprecise, answer, so we can also stop that algorithm at any point based on a time limit.

Once we have computed the bounds $b_1 \dots b_n$, we compute the new subtrees

$$\begin{aligned} & \llbracket \mathbf{x} \{C^0\} : t_0^0 b_1^0 \dots b_{n_0}^0 t_{n_0}^0 \bowtie_{C^0} b_1 \dots b_n \rrbracket \\ \text{and } & \llbracket \mathbf{x} \{C^1\} : t_0^1 b_1^1 \dots b_{n_1}^1 t_{n_1}^1 \bowtie_{C^1} b_1 \dots b_n \rrbracket \end{aligned}$$

where:

$$t_0 a_1 t_1 \dots a_m t_m \bowtie_C b_1 \dots b_n =$$

- if $m = 0$ then
 - if one of the b_i is such that $\mathbf{x} < b_i$ in C , let l be the smallest such i . The result is $t_0 b_1 \dots b_{l-1} t_0 b_l \perp b_{l+1} \dots b_m \perp$
 - else $t_0 b_1 \dots b_n t_0$

- else if $n = 0$ then $t_0 \cup \dots \cup t_m$
- else if $b_1 = a_1$ in C , we must look at b_2 , if any, in case of segment creation:
 - if $n > 1$ and $b_2 = b_1$ in C , then we create a segment. The value of that segment depends on the binary operation (as in [11]). If the operation is a join, a widening or the segment is the first argument of inclusion testing then the value is \perp . If the operation is a meet, or a narrowing, or the second argument of an inclusion testing then the value is \top . Let us call \top that value. the result is then $t_0 b_1 (\top a_1 t_1 \dots a_m t_m \bowtie_C b_2 \dots b_n)$
 - else the result is $t_0 b_1 (t_1 a_2 t_2 \dots a_m t_m \bowtie_C b_2 \dots b_n)$
- else if $a_1 \leq b_1$ in C then
 - if $m > 1$ and $b_1 \leq a_2$ in C then $(t_0 \cup t_1) b_1 t_1 a_2 \dots a_m t_m \bowtie_C b_1 \dots b_n$
 - else $(t_0 \cup t_1) a_2 t_2 \dots a_m t_m \bowtie_C b_1 \dots b_n$
- else ($b_1 \leq a_1$),
 - if $b_1 \leq \mathbf{x}$ in C then $\perp b_1 (t_0 a_1 t_1 \dots a_m t_m \bowtie_C b_2 \dots b_n)$
 - else $t_0 b_1 (t_0 a_1 t_1 \dots a_m t_m \bowtie_C b_2 \dots b_n)$

Once we have computed two new trees which agree on the bounds, we can perform the binary operation. The binary operation must be carried on the incremental pre-orders and on each pair of subtrees recursively, until reaching the leaves where we can compute the binary operation on D_ℓ .

Join There is a special case when joining two trees: in case we create new segments, there is an opportunity to add incremental pre-order informations. If we merge $\llbracket \mathbf{x} \{C\} : \dots \rrbracket$ in a context of pre-order information C_1 , and \perp in the context of C_2 , then the result is $\llbracket \mathbf{x} \{C \sqcap_c (\text{relax}_C(C_1, C_2))\} : \dots \rrbracket$ because we know that on that segment only the pre-order in C_1 is true. For that mechanism to be precise, we need to keep track of the pair of possible pre-orders during the computation of the joins.

Widening Because segments are split, their number in a decision node can grow indefinitely so that the segmented decision tree can explode in breadth. This must be prevented by a segment widening. We are in a situation where one abstract domain (the segments of a node) is a basis for another one. We can use a mechanism similar to [22] and widen on the segmentations, then on that common segmentation apply the widening child by child.

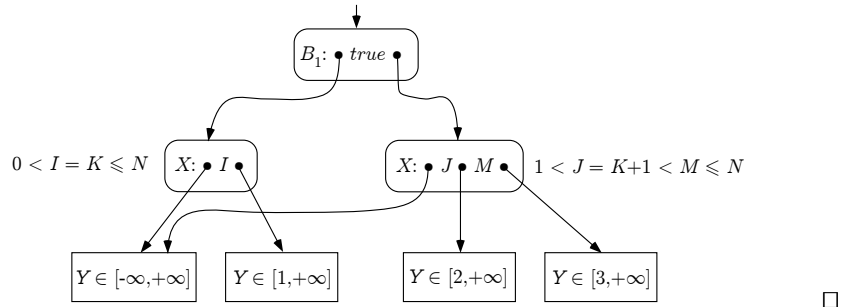
A first possibility for the segmentation widening is the segments unification of [11] which is based on the use of the common expressions in the two segmentations. In that way, the number of expressions that appear in segments can only decrease. We can achieve the same property by keeping all expressions that can occur as bounds of the first argument of the widening and only those expressions. This is easily implemented by keeping only those expressions in the graph where we look for a best chain of expressions.

7.5 Assignments

As variables may occur in leaves, expression bounds or as decision variables, we have to consider all cases, keeping in mind that those cases are not exclusive.

Assignments to leaf variables An assignment to a variable appearing in the leaf nodes only will determine the feasible paths to the leaves where it appears and perform the assignments in each of these leaves (in the abstract domain of the leaves).

Example 9. Assuming in *Ex. 5* that nothing is known on the upper bound of I , J , K , M , and N in the variable environment, the assignment $Y = X$ will determine that either $\neg B_1$ in which case if $X < I$ then else $X \geq I > 0$ so $Y \in [1, +\infty]$ or B_1 holds and so either $X < J$ in which case $Y \in [-\infty, +\infty]$, or $1 < J \leq X$ so $Y \in [2, +\infty]$, or else $1 < J < M \leq X$ and so $Y \in [2, +\infty]$. We get



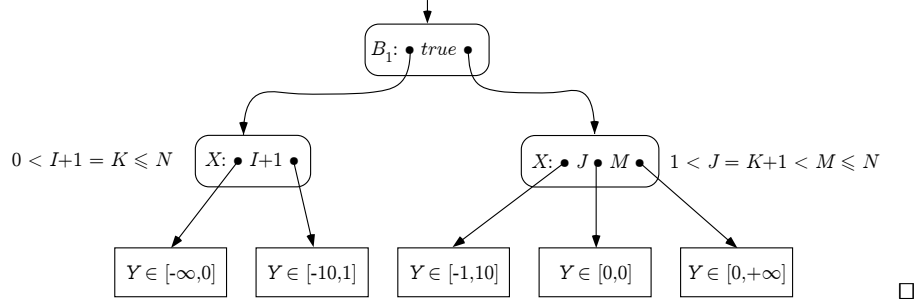
In general the assignment of an expression to a variable involves some conditions (such as absence of overflow, positiveness, non-nullity, etc) that have to be taken into account by pruning the tree as in Sect. 7.3. In case where we have to do such pruning, we can follow the same algorithm, but performing the assignment at the leaves in addition to imposing the test.

Assignments to segment bound variables An assignment to a variable appearing in segment bounds may be invertible, in which case segments which were based on the old value of the variable can be expressed based on the new value, or not invertible, in which case it is not possible to keep the segments bounds when they are only expressible in terms of the old values of the assigned variable.

More precisely, if the assignment can be expressed as $b = f(b)$ and f invertible, we can replace the variable b by $f^{-1}(b)$ in each expression appearing in bounds of the decision tree, and that encodes the same property. To complete the assignment, we must also carry it to incremental pre-orders at each node.

Example 10. Consider the assignment $I = I - 1$ in the context of *Ex. 5*. After this assignment the old value I_o of I (to which *Ex. 5* is referring to) can be expressed in terms of the new value I_n as $I_n = I_o - 1$ so $I_o = I_n + 1$ by inversion. So, we get the post-condition of the assignment $I = I - 1$ by replacing I by $I + 1$

in the segmented decision tree of *Ex. 5*.

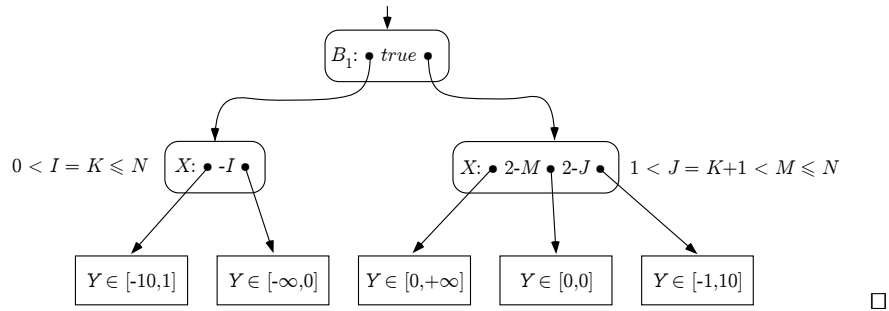


In case of non-invertible assignment to a variable b appearing in a bound, we look at all bounds where b appears. If at that bound the pre-order information can provide another expression that is known to be equal to the bound but that does not contain b , we can replace the bound by that expression. Otherwise, we drop the bound from the segmentation and merge the two consecutive subtrees that were separated by that bound. As for the tests, if that results in only one child for a node, we can push up the incremental pre-order information of that child.

In addition, the assignment must be carried also in the incremental pre-orders.

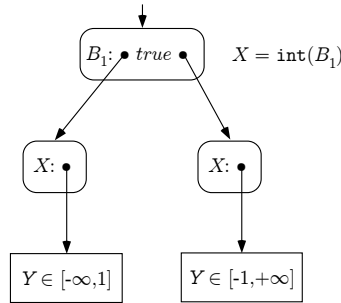
Assignments to decision variables For an invertible assignment $X = f(X)$ to a decision variable X where X_o (resp. X_n) is the value of X before (resp. after) the assignment, we have $X_n = f(X_o)$ such that $X_o = f^{-1}(X_n)$. The segment conditions $b \leq X_o$ are transformed into $b \leq f^{-1}(X_n)$ that is $f(b) \leq X_n$ when f is increasing and $f(b) \geq X_n$ when f is decreasing. Similarly $X_o < b$ are transformed into $f^{-1}(X_n) < b$ that is $X_n < f(b)$ when f is strictly increasing and $x_n > f(b)$ when f is strictly decreasing. If f also depends on other decision variables, their abstract value must be evaluated along the path to the nodes for X .

Example 11. Consider the invertible assignment is $X = \text{int}(B_1) - X$ in *Ex. 5* where $\text{int}(\text{false}) = 0$ and $\text{int}(\text{true}) = 1$ that is $X = -X$ when B_1 is *false* and $X = 1 - X$ when B_1 is *true*, which are both invertible assignments.



In case of non-invertible assignment, we cannot keep the segments related to the assigned variable. So we merge the children of that variable. In case where the assignment $\mathbf{x}=\mathbf{e}$ can be represented exactly in the pre-order domain, then we are as precise as possible. Otherwise, it is still possible to add some information in the tree. For example if the expression is a monotone function over another variable y smaller than x (for $<_D$) then we can store the inequalities implied by the segmentation over y into the incremental pre-orders associated with y . If the variable is greater than x , we can do the same if all the nodes on y share a bound.

Example 12. Consider the non-invertible assignment is $X = \text{int}(B_1)$. The post-condition is preserved while selectively merging the children. Assuming $\text{int}(b)$ to be a canonical integer expression for canonical Boolean expressions b , we get:



8 Abstracting Functions (and Array Contents)

Binary Decision Diagrams were originally developed to represent boolean functions [1]. In the same way, segmented decision trees can be used to approximate functions over totally ordered domains: we make decisions for each parameter of the function, and the leaves of the tree represent the possible values of the function for that constraint on the parameters.

Example 13. The function $\sin x$, $x \in [0, 2\pi]$ could be approximated by the segmented decision tree $\llbracket x \{0 \leq x \leq 2\pi\} : (\sin x : [0, 1]) \ \pi \ (\sin x : [-1, 0]) \rrbracket$. \square

Formally, a function $f(x_0, \dots, x_n)$ can be seen as a set of vectors of size $n + 1$ of the form $\langle v_0, \dots, v_n, f(v_0, \dots, v_n) \rangle$. Then a property over functions is a set of sets of vectors. The first abstraction we perform is to go back to sets of vectors, by taking the union of the sets, then we are in a setting where we can use segmented decision trees directly, with decision variables the first n variables ordered by the order on parameters of the function. Such abstraction could be very powerful to summarize functions and perform modular analyzes.

Multi-dimensional arrays can be seen as functions from index values to array content. So we can use the same combination of abstractions and obtain precise representations. Because arrays don't have formal parameters, we just need a convention to name to variables which will correspond to the array dimensions

and to the array content. One possibility is to subscript the array name with the number of the dimension for the indexes and with v for the content.

One more important difference between functions and arrays is the assignment. It is easily implemented as an assignment to the content variable under the appropriate condition for the dimension variables. May-assign can arise when such conditions cannot be represented exactly in the expressions allowed as bounds, but in general this mechanism will be very precise.

If we specialize that abstraction to arrays of dimension 1, we have an abstraction that is equivalent to [11] where the leaves consisted of abstractions of couples index-content.

9 Examples

9.1 Conditional Computation

```

        int x1, x2, y, z;
        struct s;
/* 0: */ x1 = 0;
        y = z = 1;
        s = INIT;
/* 1: */ while /* 2: */( z < 100 ) {
/* 3: */   if (x1 < y) s = INIT; /* 4: */
           else { /* 5: */
               if (x2 > y) s = SPECIAL; /* 6: */
               else s = computation(s, x2); /* 7: */
           }
/* 8: */   }
/* 9: */   z++;
/* 10: */   if (?) y++; /*11: */
/* 12: */   if (?) x1++; /*13: */
/* 14: */ }
/* 15: */ if (x1>y && x2>y) /* 16: */ assert(s==SPECIAL);

```

In that abstract program, the structure s is the output and the variable $x2$ is the input. On some variation, $x2$ may change at each loop iteration but that does not change the invariants here, so we just consider that its value is fixed.

A pre-analysis can fix the decision variables to be $x1$ and $x2$ as they are in guards to compute values and then in a guard to test those values. The structure s will be at the leaves (the result of `computation` is abstracted by `COMP`).

The fixpoint iteration with widening will go as follows (we write $\llbracket v : \dots \rrbracket$ for $\llbracket v \{true\} : \dots \rrbracket$):

```

1-9:  $\llbracket x1 \{x1 = 0, y = z = 1\} : \llbracket x2 : (\text{INIT}) \rrbracket \rrbracket$       {Initialization with input  $x2$ }
14:  $\llbracket x1 \{x1 \in [0, 1], y \in [1, 2], z = 2\} : \llbracket x2 : (\text{INIT}) \rrbracket \rrbracket$ 
2:  $\llbracket x1 \{0 \leq x1 \leq y \leq z, x1 < z, 0 < y\} : \llbracket x2 : (\text{INIT}) \rrbracket \rrbracket$ 
3:  $\llbracket x1 \{0 \leq x1 \leq y \leq z, x1 < z, 0 < y\} : \llbracket x2 : (\text{INIT}) \rrbracket \rrbracket$ 
4:  $\llbracket x1 \{0 \leq x1 < y \leq z < 100\} : \llbracket x2 : (\text{INIT}) \rrbracket \rrbracket$ 

```

- 5: $\llbracket x1 \{0 < y = x1 < z < 100\} : \llbracket x2 : (\text{INIT}) \rrbracket \rrbracket$
- 6: $\llbracket x1 \{0 < y = x1 < z < 100, y < x2\} : \llbracket x2 : (\text{SPECIAL}) \rrbracket \rrbracket$
- 7: $\llbracket x1 \{0 < y = x1 < z < 100, x2 \leq y\} : \llbracket x2 : (\text{COMP}) \rrbracket \rrbracket$
- 8: $\llbracket x1 \{0 < y = x1 < z < 100\} : \llbracket x2 : (\text{COMP}) \rrbracket y + 1 (\text{SPECIAL}) \rrbracket$
- 9: $\llbracket x1 \{0 \leq x1 \leq y \leq z < 100, x1 < z, 0 < y\} : \llbracket x2 : (\text{INIT}) \rrbracket y \llbracket x2 : (\text{COMP}) \rrbracket y + 1 (\text{SPECIAL}) \rrbracket$
- 10: $\llbracket x1 \{0 \leq x1 \leq y < z < 101, x1 < z - 1, 0 < y\} : \llbracket x2 : (\text{INIT}) \rrbracket y \llbracket x2 : (\text{COMP}) \rrbracket y + 1 (\text{SPECIAL}) \rrbracket$
- 11: $\llbracket x1 \{0 \leq x1 < y \leq z < 101, x1 < z - 1, 1 < y\} : \llbracket x2 : (\text{INIT}) \rrbracket y - 1 \llbracket x2 : (\text{COMP}) \rrbracket y (\text{SPECIAL}) \rrbracket$
- 12: $\llbracket x1 \{0 \leq x1 \leq y \leq z < 101, x1 < z - 1, 0 < y\} : \llbracket x2 : (\text{INIT}) \rrbracket y - 1 \llbracket x2 : (\text{COMP} \cup \text{INIT}) y (\text{COMP} \cup \text{SPECIAL}) \rrbracket y \llbracket x2 : (\text{COMP}) y (\text{COMP} \cup \text{SPECIAL}) y + 1 (\text{SPECIAL}) \rrbracket \rrbracket$
- 13: $\llbracket x1 \{1 < x1 < z < 101, 0 < y \leq z\} : \llbracket x2 : (\text{INIT}) \rrbracket y \llbracket x2 : (\text{COMP} \cup \text{INIT}) y (\text{COMP} \cup \text{SPECIAL}) \rrbracket y + 1 \llbracket x2 : (\text{COMP}) y (\text{COMP} \cup \text{SPECIAL}) y + 1 (\text{SPECIAL}) \rrbracket \rrbracket$
- 14: $\llbracket x1 \{0 \leq x1 < z < 101, 0 < y \leq z\} : \llbracket x2 : (\text{INIT}) \rrbracket y - 1 \llbracket x2 : (\text{COMP} \cup \text{INIT}) y (\text{COMP} \cup \text{SPECIAL} \cup \text{INIT}) \rrbracket y \llbracket x2 : (\text{COMP} \cup \text{INIT}) y (\text{COMP} \cup \text{SPECIAL}) \rrbracket y + 1 \llbracket x2 : (\text{COMP}) y (\text{COMP} \cup \text{SPECIAL}) y + 1 (\text{SPECIAL}) \rrbracket \rrbracket$
- 2: = 14: without $z < 101$
 {as the union of 14: and 2: is 14: here and this is the abstract loop invariant}
- 15: = 2:
- 16: $\llbracket x1 \{0 < y < x1 < z, y < x2, 99 < z\} : \llbracket x2 : (\text{SPECIAL}) \rrbracket \rrbracket$
 {The assertion in 16: is proved correct.}

9.2 Partial Array Initialization

The program below partially initializes an array **A**.

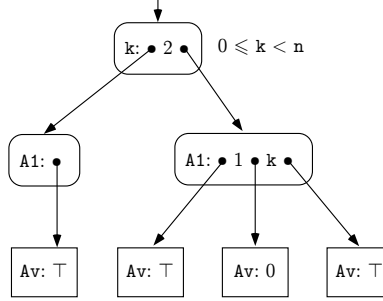
```

    int n; /* n > 0 */
    int k, A[n];
/* 0: */ k = 0;
/* 1: */ while /* 2: */ (k < n) {
/* 3: */     if (k > 0) {
/* 4: */         A[k] = 0;
/* 5: */     };
/* 6: */     k = k+1;
/* 7: */ };
/* 8: */

```

The ordering abstract domain D_c is assumed to be the octagon abstract domain [20]). Following Sect. 8, an array **A** is abstracted by two fresh variables $A1 \in \mathbb{D}$ to segment

indices $A1$ of array A , $A1 \in [A.low, A.high]$ and a variable $Av \in \mathbb{D}$ standing for any value of the array in a given segment such that $Av <_{\mathbb{D}} A1$ and Av is a leaf. For leaves we use constant propagation [15]. The loop invariant found at point 3 is



The fixpoint iteration with widening is the following:

- 0: $\llbracket k \{0 < n, 0 \leq A1 < n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket$ $\{k \text{ and } A \text{ uninitialized}\}$
 ℓ : \perp $\{\ell = 1, \dots, 8, \text{infimum}\}$
1: ,2: ,3: ,6: $\llbracket k \{k = 0 < n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket$ $\{0: \text{ where } k = 0, k < n, k \leq 0\}$
7: $\llbracket k \{k = 1 \leq n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket$ $\{6: \text{ where } k = k + 1\}$
2: ,3: $\llbracket k \{0 \leq k \leq 1, k < n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket$ $\{\text{joining 1: and 7: , test } k < n\}$
4: $\llbracket k \{1 = k < n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket$ $\{3: \text{ with } k > 0\}$
5: $\llbracket k \{1 = k < n\} : \llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \rrbracket$
 $\{4: \text{ with } A[k] = 0 \text{ where } k = 1\}$
6: $\llbracket k \{0 \leq k \leq 1, k < n\} : \llbracket A1 \{k = 0\} : (Av : \top) \rrbracket \ 1$
 $\llbracket A1 \{k = 1\} : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \rrbracket$
 $\{\text{joining 3: and } k \leq 0 \text{ so } k = 0 \text{ together with 5: where } k = 1\}$
7: $\llbracket k \{1 \leq k \leq 2, k \leq n\} : \llbracket A1 \{k = 1\} : (Av : \top) \rrbracket \ 2$
 $\llbracket A1 \{k = 2\} : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \rrbracket$ $\{6: \text{ where } k = k + 1\}$
1: \sqcup_t 7: $\llbracket k \{0 \leq k \leq 2, k \leq n\} : \llbracket A1 \{0 \leq k \leq 1\} : (Av : \top) \rrbracket \ 2$
 $\llbracket A1 \{k = 2\} : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \rrbracket$ $\{\text{join of 1: and 7:}\}$
2: ,3: $\llbracket k \{0 \leq k < n\} : \llbracket A1 \{0 \leq k \leq 1\} : (Av : \top) \rrbracket \ 2$
 $\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k \ (Av : \top) \rrbracket \rrbracket$ $\{2: \nabla (1: \sqcup_t 7)^7, \text{ test } k < n\}$
4: $\llbracket k \{0 < k < n\} : \llbracket A1 \{k = 1\} : (Av : \top) \rrbracket \ 2$
 $\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k \ (Av : \top) \rrbracket \rrbracket$ $\{3: \text{ with } k > 0\}$
5: $\llbracket k \{0 < k < n\} : \llbracket A1 \{k = 1\} : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \ 2$
 $\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k + 1 \ (Av : \top) \rrbracket \rrbracket$ $\{4: \text{ with } A[k] = 0\}$
6: $\llbracket k \{0 \leq k < n\} : \llbracket A1 \{k = 0\} : (Av : \top) \rrbracket \ 1$
 $\llbracket A1 \{k = 1\} : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \ 2$
 $\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k + 1 \ (Av : \top) \rrbracket \rrbracket$
7: $\llbracket k \{0 < k \leq n\} : \llbracket A1 \{k = 1\} : (Av : \top) \rrbracket \ 2$ $\{\text{joining 3: and } k \leq 0 \text{ with 5:}\}$

$$\begin{array}{l}
\llbracket A1 \{k = 2\} : (Av : \top) \ 1 \ (Av : 0) \ 2 \ (Av : \top) \rrbracket \ 3 \\
\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k \ (Av : \top) \rrbracket \rrbracket \quad \{6: \text{ where } k = k + 1\} \\
1: \sqcup_t 7: \llbracket k \{0 \leq k \leq n\} : \llbracket A1 \{0 \leq k \leq 1\} : (Av : \top) \rrbracket \rrbracket \ 2 \\
\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k \ (Av : \top) \rrbracket \rrbracket \quad \{\text{join of } 1: \text{ and } 7:\} \\
2:, 3: \llbracket k \{0 \leq k < n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket \ 2 \\
\llbracket A1 : (Av : \top) \ 1 \ (Av : 0) \ k \ (Av : \top) \rrbracket \rrbracket \\
\{2: \nabla (1: \sqcup_t 7), \text{ test } k < n, \text{ convergence, } 3: \text{ is the abstract loop invariant}\} \\
8: \llbracket k \{0 \leq k = n, 0 \leq A1 < n\} : \llbracket A1 : (Av : \top) \rrbracket \rrbracket \ 2 \\
\llbracket A1 \{0 \leq k < n\} : (Av : \top) \ 1 \ (Av : 0) \rrbracket \rrbracket \\
\{2: \text{ and } k \geq n, \text{ program postcondition}\}
\end{array}$$

Observe that the segmented decision tree automatically discovers a partition of the loop body as given by the condition $k > 0$ while the segmented array partitions the values of the array elements according to variable k .

9.3 Multidimensional Arrays

The program below partially initializes a matrix M .

```

int m, n; /* m, n > 0 */
int i, j, M[m,n];
/* 0: */ i = 0;
/* 1: */ while /* 2: */ (i < m) {
/* 3: */     j = i+1;
/* 4: */     while /* 5: */ (j < n) {
/* 6: */         M[i,j] = 0;
/* 7: */         j = j+1;
/* 8: */     };
/* 9: */     i = i+1;
/* 10: */ };
/* 11: */

```

A global invariant is $0 \leq M1 < m$ and $0 \leq M2 < n$, so we keep it implicit in the following fixpoint iteration:

$$\begin{array}{l}
0: \llbracket M1 : \llbracket M2 : (Mv : \top) \rrbracket \rrbracket \quad \{\text{program precondition: } i, j, \text{ and } A \text{ uninitialized}\} \\
\ell: \perp \quad \{\ell = 1, \dots, 11, \text{ infimum}\} \\
1:, 2:, 3: \llbracket M1 \{i = 0\} : \llbracket M2 : (Mv : \top) \rrbracket \rrbracket \quad \{0: \text{ with } i = 0, i < m^8\} \\
4:, 5:, 6: \llbracket M1 \{i = 0, j = i + 1 = 1 < n\} : \llbracket M2 : (Mv : \top) \rrbracket \rrbracket \\
\quad \{3: \text{ with } j = i+1; j < n\} \\
7: \llbracket M1 \{i = 0, j = i + 1 = 1 < n\} : \\
\quad \llbracket M2 : (Mv : \top) \ j \ (Mv : 0) \ j + 1 \ (Mv : \top) \rrbracket \ i + 1 \llbracket M2 : (Mv : \top) \rrbracket \\
\quad \rrbracket \quad \{6: \text{ with } M[i, j] = 0;\}
\end{array}$$

⁷ When a new branch is taken in a test within a loop the widening is usually delayed, which we avoid to shorten the example.

8: $\llbracket \text{M1 } \{i = 0, j = i + 2 = 2 \leq n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } j - 1 \text{ } (\text{Mv} : 0) \text{ } j \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{7: \text{ with } j = j + 1;\}$

4: \sqcup_t 8: $\llbracket \text{M1 } \{i = 0, i + 1 \leq j \leq i + 2 \leq n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \text{ } j \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{\text{join of 4: and 8:}\}$

5: $\llbracket \text{M1 } \{i = 0, i + 1 \leq j \leq n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \text{ } j \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{5: \nabla (4: \sqcup_t 8:)^9\}$

9: $\llbracket \text{M1 } \{i = 0, i + 1 \leq j = n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{5: \text{ and } j \geq n\}$

10: $\llbracket \text{M1 } \{i = 1, i \leq j = n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{9: \text{ and } i = i + 1;\}$

1: \sqcup_t 10: $\llbracket \text{M1 } \{i = 1, i \leq j = n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{\text{join of 1: and 10:}\}$

2: $\llbracket \text{M1 } \{0 \leq i\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{2: \nabla (1: \sqcup_t 10:)\}$

3: $\llbracket \text{M1 } \{0 \leq i < m\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{2: \text{ and } j < n\}$

4:, 5:, 6: $\llbracket \text{M1 } \{0 \leq i < m, j = i + 1 < n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{3: , j = i + 1; \text{ and } j < n\}$

7: $\llbracket \text{M1 } \{0 \leq i < m, j = i + 1 < n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \text{ } j \text{ } (\text{Mv} : 0) \text{ } j + 1 \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{6: \text{ and } M[i, j] = 0;\}$

8: $\llbracket \text{M1 } \{0 \leq i < m, j = i + 2 \leq n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \text{ } n \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \text{ } j - 1 \text{ } (\text{Mv} : 0) \text{ } j \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{7: \text{ with } j = j + 1;\}$

4: \sqcup_t 8: $\llbracket \text{M1 } \{0 \leq i < m, i + 1 \leq j \leq i + 2 \leq n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \text{ } i + 1 \text{ } (\text{Mv} : 0) \text{ } j \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{\text{join of 4: and 8:}\}$

5: $\llbracket \text{M1 } \{0 \leq i < m, i + 1 \leq j \leq n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \text{ } i + 1 \text{ } (\text{Mv} : 0) \text{ } j \text{ } (\text{Mv} : \top) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{5: \nabla (4: \sqcup_t 8:)\}$

9: $\llbracket \text{M1 } \{0 \leq i < m, i + 1 \leq j = n\} : \llbracket \text{M2} : (\text{Mv} : \top) \text{ } 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i \llbracket \text{M2} : (\text{Mv} : \top) \text{ } i + 1 \text{ } (\text{Mv} : 0) \rrbracket \text{ } i + 1 \llbracket \text{M2} : (\text{Mv} : \top) \rrbracket \rrbracket$
 $\} \quad \{5: \text{ and } j \geq n\}$

- 10: $\llbracket \mathbf{M1} \{0 < i \leq m, i \leq j = n\} : \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \ 1 \ (\mathbf{Mv} : 0) \rrbracket i - 1$
 $\llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \ i \ (\mathbf{Mv} : 0) \rrbracket i \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \rrbracket$
 \rrbracket {9: and $i = i+1$;}
- 1: \sqcup_t 10: $\llbracket \mathbf{M1} \{0 \leq i \leq m\} : \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \ \mathbf{M1} + 1 \ (\mathbf{Mv} : 0) \rrbracket i \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \rrbracket \rrbracket$
{join of 1: and 10:(segments unification yields $1 \leq \mathbf{M1} + 1 \leq i$ for subtree merges)}
- 2: $\llbracket \mathbf{M1} \{0 \leq i\} : \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \ \mathbf{M1} + 1 \ (\mathbf{Mv} : 0) \rrbracket i \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \rrbracket \rrbracket$
{2: \sqsubseteq (1: \sqcup_t 10:), stabilization at a fixpoint}
- 11: $\llbracket \mathbf{M1} \{0 < m = i\} : \llbracket \mathbf{M2} : (\mathbf{Mv} : \top) \ \mathbf{M1} + 1 \ (\mathbf{Mv} : 0) \rrbracket \rrbracket$
{2: and $i \geq m$, program postcondition.}

10 Conclusion

Many static analyses are very impressive on small examples but fail to scale up. The problem mainly originates from the explosion of possible cases in handling disjunctions. Mastering the exponential growth is the key to scalability, while enabling weak forms of disjunction is essential to the precision which is necessary to avoid false alarms. Based on two abstract domain functors that have shown experimentally to scale up, we have proposed a new combination which expressivity is better than each of them taken separately and which complexity can be mastered by imposing both static restrictions (like maximal depth or variable packing) and dynamic restrictions (by widening to control the breath of the tree).

References

1. Randal E. Bryant. Graph based algorithms for boolean function manipulation. *IEEE Transactions on Computers*, C-35:677–691, August 1986.
2. P. Cousot. *Méthodes itératives de construction et d'approximation de points fixes d'opérateurs monotones sur un treillis, analyse sémantique de programmes (in French)*. Thèse d'État ès sciences mathématiques, Université scientifique et médicale de Grenoble, Grenoble, 21 March 1978.
3. 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. LNCS 2772, Springer, Berlin, Taormina, 29 June – 4 July 2003.
4. P. Cousot and R. Cousot. Static determination of dynamic properties of programs. In *Proc. 2nd Int. Symp. on Programming*, pages 106–130, Paris, 1976. Dunod, Paris.
5. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *4th POPL*, pages 238–252, Los Angeles, 1977. ACM Press.
6. P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *6th POPL*, pages 269–282, San Antonio, 1979. ACM Press.

⁹ When a new branch is taken in a test within a loop the widening is usually delayed, which we avoid to shorten the example.

7. P. Cousot and R. Cousot. Abstract interpretation frameworks. *Journal of Logic and Computation*, 2(4):511–547, August 1992.
8. P. Cousot and R. Cousot. Comparing the Galois connection and widening/narrowing approaches to abstract interpretation. In M. Bruynooghe and M. Wirsing, editors, *Proc. 4th Int. Symp. Programming Language Implementation and Logic Programming, PLILP '92*, Leuven, 26–28 August 1992, LNCS 631, pages 269–295. Springer, Berlin, 1992.
9. P. Cousot and R. Cousot. Formal language, grammar and set-constraint-based program analysis by abstract interpretation. In *Proc. of the Seventh ACM Conf. on Functional Programming Languages and Computer Architecture*, pages 170–181, La Jolla, 25–28 June 1995. ACM Press.
10. P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, and X. Rival. The ASTRÉE static analyzer. www.astree.ens.fr and www.absint.com/astree/.
11. P. Cousot, R. Cousot, and F. Logozzo. A parametric segmentation abstract domain function for the fully automatic and scalable inference of array properties. Research report, Microsoft Research, Remond, September 2009.
12. P. Cousot and N. Halbwegs. Automatic discovery of linear restraints among variables of a program. In *5th POPL*, pages 84–97, Tucson, 1978. ACM Press.
13. S. Graf and H. Saïdi. Verifying invariants using theorem proving. In R. Alur and T.A. Henzinger, editors, *Proc. 8th Int. Conf. on Computer Aided Verification, CAV '97*, New Brunswick, LNCS 1102, pages 196–207. Springer, Berlin, July 31 – August 3 1996.
14. P. Granger. Static analysis of arithmetical congruences. *Int. J. Comput. Math.*, 30(3 & 4):165–190, 1989.
15. G. Kildall. A unified approach to global program optimization. In *Conf. Record of the First Annual ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages*, pages 194–206, Boston, October 1973. ACMpress.
16. L. Mauborgne. Analyse statique et domaines abstraits symboliques. Mémoire d'habilitation à diriger les recherches en informatique, Université de Paris Dauphine, 12 February 2007.
17. L. Mauborgne. Abstract interpretation using TDGs. In B. Le Charlier, editor, *Proc. 1st Int. Symp. on Static Analysis, SAS '94*, Namur, 20–22 september 1994, LNCS 864, pages 363–379. Springer, Berlin, 1994.
18. L. Mauborgne. Abstract interpretation using typed decision graphs. *Science of Computer Programming*, 31(1):91–112, May 1998.
19. L. Mauborgne. Binary decision graphs. In A. Cortesi and G. Filé, editors, *Proc. 6th Int. Symp. on Static Analysis, SAS '99*, Venice, 22–24 september 1999, LNCS 1694, pages 101–116. Springer, Berlin, 1999.
20. A. Miné. The octagon abstract domain. *Higher-Order and Symbolic Computation*, 19:31–100, 2006.
21. A. Tarski. A lattice theoretical fixpoint theorem and its applications. *Pacific Journal of Mathematics*, 5:285–310, 1955.
22. Arnaud Venet. Abstract cofibered domains: Application to the alias analysis of untyped programs. In *Static Analysis Symp. (SAS'96)*, LNCS 1145. Springer, Berlin, 1996.