

DelExp: a Relational Container Abstraction with Applications to Compositional Analysis

Milla Valnet  

LIP6, Sorbonne Université, F-75005, Paris, France

Raphaël Monat  

Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 CRISTAL, F-59000 Lille, France

Antoine Miné  

LIP6, Sorbonne Université, F-75005, Paris, France

Abstract

Data containers, such as lists, arrays, trees, etc, raise challenges for program verification. In static analysis by abstract interpretation, one popular approach is summarization: multiple elements of a data structure are abstracted into a single one, favoring performance over precision. This technique is at the core of most container abstractions – from smashing to segmentation – of arrays, lists or algebraic data types.

However, summarization approaches are unable to express relations between containers, even when relational numerical abstract domains are used. Our work introduces DelExp, a new domain able to express relations between summarized variables. DelExp can state that the content of a data structure is included in the content of another data structure, up to a given transformation. DelExp is language-agnostic, modular in the abstraction chosen for any other types (integers, strings, functions, etc.), and can be seamlessly combined with existing container abstractions. We show how DelExp allows us to infer precise summaries for compositional analyses of container-manipulating functions in a pure functional language. We present extensions to DelExp supporting polymorphism and higher-order transformations.

Our implementation of DelExp within the MOPSA static analysis platform confirms that DelExp works out of the box with pre-existing container abstractions. Our evaluation targets both Python programs manipulating lists and relational summary generation for OCaml functions handling algebraic data types.

2012 ACM Subject Classification Software and its engineering → Automated static analysis; Theory of computation → Program analysis; Software and its engineering → Functional languages

Keywords and phrases Static Value Analysis, Abstract Interpretation, Functional Programming

Digital Object Identifier [10.4230/LIPIcs...12](https://doi.org/10.4230/LIPIcs...12)

1 Introduction

Most programs use various kinds of containers such as lists, maps, dictionaries or sets, i.e. data structures of possibly unbounded size storing a collection of objects. Collecting information on those collections can be crucial to verify some properties on the program behaviours, e.g. assertions over elements inside the containers. In this work, we aim at discovering *relational properties* over containers in order to improve the precision of static value analysis. Our approach is rooted in the abstract interpretation framework: we compute over-approximations of the possible behaviour of a program. The resulting analysis is fully automated and does not require annotations.

Value analysis on containers raises specific challenges: standard abstract domains manipulate a fixed number of objects whereas a container can be of unbounded size. A standard approach from Blanchet et al. [3] is to fold, or “smash”, together a set of objects into one *summarized* object. A single abstract object then represents multiple concrete objects. Gopan et al. [13] proposed a systematic approach to transform standard numerical



© Milla Valnet Raphaël Monat Antoine Miné;
licensed under Creative Commons License CC-BY 4.0

Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

```

# assumes l:List[int]
r=[]
for x in l:
    r.append(2*x)

```

■ **Figure 1** Motivating list-manipulating function, written in Python.

domains, including relational domains, into summarizing domains, manipulating summarized variables – that we will call weak variables. However, this approach cannot express relations between the set of values inside containers – e.g. inclusion between one set and another.

Cox et al. [11] defined a domain called QUIC (Quantified Union/Intersection Constraint) graphs. It can express constraints over numerical sets, under the form $\bigcap X_i \subseteq \{\nu \in \bigcup Y_j \mid B[\nu]\}$, B being an abstract domain acting as a predicate on the elements ν of the collection. When abstracting containers as sets, this domain can express some relations between containers. However, this domain is limited to inclusions between union/intersection of sets refined by a predicate. As is, this domain cannot represent an element-wise transformation on a given set. By contrast, the domain presented in our article can handle transformations similar to map.

Consider the Python program in Fig. 1. Given a list l of integer, it builds the list r by multiplying every element in l by 2. Neither standard relational domains on integers, such as the polyhedra domain, nor the QUIC graphs domain can express that the elements of the list r are elements of the list l multiplied by two. We develop here a family of domains, called DelExp, which, given a transformation θ , can state relations of the form $x \prec \theta(y)$, meaning that x abstracts a set of elements included in $\theta(y) = \{\theta(e) \mid e \in y\}$. In particular, by introducing the weak variables l_c and r_c to represent all elements of l and r , our derived analysis can infer that r contains only elements that are the double of elements from l , i.e. $l_c \prec 2 \times r_c$. This family of domains is *language-agnostic*: it only operates on weak variables (here, l_c and r_c), and therefore, can be used in product with any analysis manipulating weak variables. Indeed, our implementation has been exploited within both a Python and an OCaml analysis.

An important use case of our relational domains is the definition of compositional (or function-modular) analyses. These analyses consist in analyzing functions once, at their definition site, to infer a summary, i.e. a contract, valid for every input, that over-approximates the function behaviour. The summary is then applied to analyze each function call. One difficulty faced by compositional analyses is precision loss: analyses need to produce a unique summary, correct for every possible input. A key ingredient to achieve interesting precision levels is to rely on relational abstract domains which are able to express relations between variables. Relational domains have been used extensively in previous works to create input-output relational summaries [8, 12, 20, 21].

Recently, a compositional analysis was developed for higher-order languages such as Haskell or OCaml [31]. However, since this analysis is based on summarization [13] to abstract user-defined Algebraic Data Types (ADTs), it fails to express input-output relations between two ADTs. Consider the OCaml program in Fig. 2. The previous analysis [31] creates two weak variables, namely `l.Cons.1` and `r.Cons.1`, abstracting the content of the input list l and the output list r (contained in the 1st field of `Cons`). It only infers \top as summary for `mult2`, as it is unable to express a relation between `l.Cons.1` and `r.Cons.1`. However, with our DelExp domains, the analysis is able to infer: $\forall l, r. r = \text{mult2}(l) \implies r.\text{Cons.1} \prec 2 \times l.\text{Cons.1}$.

```

type list = Cons of int * list | Nil
let rec mult2 l =
  match l with
  | Cons(h,q) -> Cons(2*h, mult2 q)
  | Nil -> Nil

```

■ **Figure 2** Motivating list-manipulating function, written in OCaml

86 It means that the integer `Cons` fields of variable r (represented by `r.Cons.1`) are included in
 87 the values of the integer `Cons` fields of variable l (represented by `l.Cons.1`) multiplied by two.

88 Note that the variables that are used in the generated summary for `mult2` are not only
 89 weak, i.e., they represent an unbounded set of concrete memory locations, but are also
 90 *optional*, i.e., this set of location may be empty. Indeed, if l is `Nil`, then `l.Cons.1` is not
 91 defined, neither is `r.Cons.1`. Environments with optional variables, called *heterogeneous*
 92 environments [19], require careful handling, as the standard join and meet operators are
 93 unsound and need to be adapted. Our domains support them.

94 The OCaml list type is a user-defined Algebraic Data Type whose abstraction was
 95 delegated to the weak variable `l.Cons.1`. Our DelExp domain then manipulates those
 96 variables. This works similarly with tree data-structures. Consider the following ADT:

```

type tree = Node of tree * int * tree | Leaf

```

97 The ADT domain delegates the abstraction of all node labels of a variable x to `x.Node.2`,
 98 on which our DelExp domain can operate. Consequently, it can infer a similar summary for
 99 a function multiplying by 2 every label of a tree.

100 Finally, a relation $x < y$, meaning that the elements of x form a subset of the elements of
 101 y , does not need hypotheses on the type and underlying abstract domain for x and y . They
 102 can abstract containers of integers, strings, etc. In fact, they can even abstract polymorphic
 103 containers – e.g. a polymorphic list `'a list`.

104 **Contributions.** This paper presents the following contributions:

- 105 ■ We define the *DelExp* abstract domain, stating simple constraints of the form $x < y$. This
 106 formalization supports optional weak variables – i.e. possibly empty containers. It is
 107 language-agnostic and domain-modular – i.e. no assumption is made on the underlying
 108 domain and type of weak variables. Therefore, it can be used on polymorphic lists,
 109 algebraic data types or array segments.
- 110 ■ We generalize it to a family of abstract domains stating constraints “up to transformation”,
 111 of the form $x < \theta(y)$. Depending on the chosen transformation, those domains can also
 112 express relations between polymorphic variables.
- 113 ■ As a case study, we develop a compositional analysis for a pure polymorphic functional
 114 language leveraging those domains. It can generate summaries for higher-order polymorphic
 115 recursive functions manipulating user-defined algebraic data types.
- 116 ■ We have implemented the resulting analysis in the MOPSA platform [18] and assessed
 117 its efficiency and precision both on list-manipulating Python programs and on ADTs-
 118 manipulating functions from the OCaml standard library. DelExp was integrated
 119 seamlessly into MOPSA: it did not require any change to the other abstract domains for
 120 the Python and OCaml analyses to work.

$$\begin{aligned}
e_1, e_2, \dots \in \mathcal{E} &::= x \in \mathbb{V} \mid n \in \mathbb{Z} \mid \text{rand}(1, n) \mid e_1 \oplus e_2 \mid t[n] \mid [e_1; \dots; e_n] \\
s_1, s_2, \dots \in \mathcal{S} &::= s_1; s_2 \\
&\mid x \leftarrow e \quad \mid \text{while } e \text{ do } s \\
&\mid t[e] \leftarrow e \quad \mid \text{if } (e)\{s_1\} \text{ else } \{s_2\}
\end{aligned}$$

■ **Figure 3** Syntax of the Universal toy language

Set of program variables	\mathbb{V}
Values	$\mathcal{V} ::= \mathbb{Z} \cup \{[v_1; \dots; v_n] \mid n \in \mathbb{N}, v_i \in \mathcal{V}\}$
Concrete environments	$\Sigma = \mathbb{V} \rightarrow \mathcal{V}$
Concrete semantics of expr. e	$\mathbb{E}[[e]] : \mathcal{P}(\Sigma) \rightarrow \mathcal{P}(\mathcal{V})$
Concrete semantics of statement. s	$\mathbb{S}[[s]] : \mathcal{P}(\Sigma) \rightarrow \mathcal{P}(\Sigma^\perp)$

■ **Figure 4** Formalization of the concrete semantics

121 **Limitations.** As it is based on weak variables, our container analysis, as a first step, over-
122 approximates containers as sets. The abstraction is content-based: it expresses content
123 inclusion between containers – possibly up to transformation. Therefore, it cannot express
124 order-dependent specification – e.g. list sortedness – nor multiplicity-sensitive properties
125 – e.g. equality between two lists. It cannot express relations between specific elements
126 of the container – e.g. the n -th element of a container relates to the m -th element of
127 another. Inferring such fine-grained relational properties on generic containers would require
128 developing new domains and is considered future work.

129 **Outline.** Section 2 presents a relational analysis on a simple toy language with weak variables.
130 Section 3 proposes a simple abstract domain expressing relations on summarization variables.
131 Section 4 shows how we can derive more expressive, "up to transformation" relations from
132 this domain. Section 5 defines a compositional analysis on a pure functional language and
133 motivates the use of the DelExp domain in this context. Section 6 adds the support of
134 heterogeneous environment to the DelExp domain, to support possibly empty containers.
135 Section 7 derives a sound and automatic analysis for our polymorphic functional language.
136 Section 8 presents our implementation and experimental results for both OCaml and Python
137 analysis. Section 9 presents related works and Sec. 10 concludes.

138 2 Relational Analysis with Weak Variables

139 This section introduces backgrounds on relational analyses for containers in a toy imperative
140 language. Section 2.1 introduces a toy imperative language. Section 2.2 motivates and defines
141 relational analyses. Section 2.3 introduces weak variables and how they can be abstracted
142 in relational domains. These weak variables are leveraged in Sec. 2.4 to showcase different
143 container abstractions, on Python lists and on OCaml Algebraic Data Types.

144 2.1 An Imperative Toy Language

145 We define in Fig. 3 a simple imperative toy language with integers and arrays. \mathbb{V} is the set of
146 variables, \oplus any operator of the language ($+$, $-$, \times , $/$, $==$, $<$, $>$, etc.). Figure 4 describes

```

1  x <- rand(1,10);
2  y <- 2*x ;
3  if {x=10} z <- y else z <- 20 ;

```

■ **Figure 5** Example program motivating relational analysis.

Abstract environments	$\Sigma^\# = (\mathbb{V} \rightarrow \mathcal{D}) \times \mathcal{D}_r$
Abstract semantics of expr. e	$\mathbb{E}^\#[e] : \Sigma^\# \rightarrow \mathbb{V} \times \Sigma^\#$
Abstract semantics of stmt. s	$\mathbb{S}^\#[s] : \Sigma^\# \rightarrow \Sigma^\#$

■ **Figure 6** Formalization of the abstract semantics.

147 the semantic domain of the language. The semantic values of the language \mathcal{V} are integers
 148 and arrays. Σ is the set of concrete environments, associating a value to each variable. The
 149 concrete semantics of an expression e is $\mathbb{E}[e] : \mathcal{P}(\Sigma) \rightarrow \mathcal{P}(\mathcal{V})$. It computes all possible
 150 semantic values of an expression in a set of possible environments (e.g. $\text{rand}(1, n)$ can yield n
 151 different values). The concrete semantics of a statement s is $\mathbb{S}[s] \in \mathcal{P}(\Sigma) \rightarrow \mathcal{P}(\Sigma^\perp)$, where
 152 \perp denotes non-termination. It updates input environments by executing the statement. Its
 153 semantics is standard, and not detailed here.

154 2.2 Relational Analysis

155 Given a language, we want to define a computable abstract semantics $\mathbb{E}^\#[e]$ of an
 156 expression e , based on abstract domains. It delegates the abstraction of an object of type τ
 157 (e.g. integers, pairs, ADTs) to the relevant abstract domain \mathcal{D}_τ . Usually, such a semantics
 158 takes as parameter an expression e and an abstract environment (abstracting the possible
 159 values of all the variables), and returns an abstract value representing the set of possible
 160 values of the expression. However, this formulation cannot easily express a relation between
 161 the value of the expression and those of the variables. Indeed, let us consider the program in
 162 Fig. 5. A numerical domain such as the interval domain can abstract x and y independently,
 163 as $x \in [1, 10]$ and $y \in [2, 20]$, but cannot infer the relation $y = 2*x$. It then fails to infer
 164 that $z=20$. To overcome this limitation, we need an abstract numerical domain able to store
 165 relations between variables – e.g. the polyhedra domain [10].

166 In our solution to this problem (see Fig. 6), we will separate non-relational domains \mathcal{D} and
 167 relational domains \mathcal{D}_r . We will maintain as abstract environments a pair, containing a map
 168 $m \in \mathbb{V} \rightarrow \mathcal{D}$ from variables to abstract values (i.e. elements of non-relational domains) and a
 169 tuple $d \in \mathcal{D}_r$ of relations from the relational domains. Thus, $\sigma^\# = (m, d) \in \Sigma^\#$. Consequently,
 170 the abstract semantics of an expression e , $\mathbb{E}^\#[e]$, returns a variable aside with an abstract
 171 environment storing the relations over this variable. The abstract semantics of a statement
 172 $\mathbb{S}^\#[s]$ updates an abstract environment with statement s . For example, let p be the program
 173 in Fig. 5, and $[\cdot]$ a starting environment with no information on program variables. The
 174 analysis of p is:

175 $\mathbb{S}^\#[p][\cdot] = (x \rightarrow [1, 10], y \rightarrow [2, 20], z \rightarrow [20, 20]; y = 2x)$

176 2.3 Content Summarization with Weak Variables

```

1 x <- [a-1;a;a+1]; (* array definition *)
2 y <- x[2]; (* get array cell *)
3 x[0] <- a+2; (* modify array cell *)

```

■ **Figure 7** Small example to illustrate weak variables

177 Abstracting containers implies representing an unknown number of objects. However,
 178 standard abstract domains such as the interval domain or the polyhedra domain can only
 179 manipulate a known, finite number of variables. To re-use those abstractions, Blanchet
 180 et al. [3] proposed to *summarize* a (possibly infinite) set of objects into one unique variable.
 181 Therefore, we can delegate the analysis of containers' content to standard domains. This
 182 approach was formalized by Gopan et al. [13] and extended to relational domains. Those
 183 summarization variables, called *weak variables*, represent multiple concrete elements and
 184 consequently concretize to a set of values. Standard variables are then called *strong variables*.
 185 We define \mathbb{V}^{weak} (resp. $\mathbb{V}^{\text{strong}}$) the sets of weak (resp. strong) variables. We add those
 186 variables to our language from Fig. 3.

187 Consider the array-manipulating program in Fig. 7. The array \mathbf{x} is composed of 3 concrete
 188 cells (i.e. 3 objects). Those concrete elements can be summarized together into one single
 189 weak variable $x_c \in \mathbb{V}^{\text{weak}}$ abstracting the content of \mathbf{x} . Weak variables then concretize to a
 190 set of values. We update our definition of concrete environments:

191 ► **Definition 2.1** (Environments). *The set of concrete environments is defined as $\Sigma =$*
 192 *$(\mathbb{V}^{\text{weak}} \rightarrow \mathcal{P}(\mathcal{V})) \times (\mathbb{V}^{\text{strong}} \rightarrow \mathcal{V})$.*

193 We define the lift $\uparrow \cdot : \mathcal{P}(\mathcal{V}) \cup \mathcal{V} \rightarrow \mathcal{P}(\mathcal{V})$ as $\uparrow X = X$ if $x \in \mathcal{P}(\mathcal{V})$, or $\{X\}$ if $x \in \mathcal{V}$:
 194 it extracts the set of possible values of a variable, whether it is weak or strong. A sound
 195 abstraction for $x_c \in \mathbb{V}^{\text{weak}}$ in the polyhedra domain is $a - 1 \leq x_c \leq a + 1$: all concrete cells
 196 represented by x_c are between $a - 1$ and $a + 1$. Note that those variables must be handled
 197 with care: e.g. modifying any of the concrete cells a summary variable represents must affect
 198 its corresponding weak variable. Their manipulation then requires specific operators.

199 ► **Definition 2.2** (Expand). *Given two variables v and v' , $\text{expand}(v, v')$ adds to the environment*
 200 *a variable v' with the same potential values as v :*

$$201 \quad \mathbb{S}[\text{expand}(v, v')]\mathbb{S} = \{\sigma[v' \rightarrow Z] \mid \sigma \in S \wedge (v' \in \mathbb{V}^{\text{weak}} \implies Z \subseteq \uparrow \sigma(v)) \\ 202 \quad \wedge (v' \in \mathbb{V}^{\text{strong}} \implies Z \in \uparrow \sigma(v))\}$$

203 It corresponds to copying every constraint verified by v to v' . This is the operation used
 204 when reading a cell of a container (cf. line 2 from Fig. 7): if the summarization variable
 205 respects a constraint, so does the value contained in the cell. Knowing $a - 1 \leq x_c \leq a + 1$,
 206 we cannot derive that $x_c = y$. Indeed, if we discover that $y = a + 1$ it doesn't mean that
 207 $x_c = a + 1$, i.e. that every element in \mathbf{x} is equal to $a + 1$. However, we can safely derive that
 208 y verifies the same constraints as x_c , i.e. $a - 1 \leq y \leq a + 1$.

209 ► **Definition 2.3** (fold). *Given two variables v and v' , $\text{fold}(v, v')$ adds the values of v' as*
 210 *possible values for v :*

$$211 \quad \mathbb{S}[\text{fold}(v, v')]\mathbb{S} = \{\sigma[v \rightarrow Z] \mid \sigma \in S \wedge Z \subseteq \uparrow \sigma(v) \cup \uparrow \sigma(v')\}$$

212 This is the operation used when adding or modifying the cell of a container (cf. line 3):
 213 given a cell from the modified container, it is either a cell from the former container or a

214 new cell. In the abstract, it either respects constraints from the former container or from the
 215 new one. Consequently, after line 3, $a - 1 \leq x_c \leq a + 1 \cup x_c = a + 2$, i.e. $a - 1 \leq x_c \leq a + 2$.

216 Abstraction of those operators are provided by each domain implementing the support
 217 for weak variables. Gopan et al. [13] defines computable abstractions for classic numerical
 218 domains (e.g. polyhedra, octagons). Once those operators are defined, given an abstract
 219 semantics $\mathbb{S}^\sharp[\cdot]$ on a language without weak variables, we can derive an abstract semantics
 220 of assignments for weak variables:

► **Definition 2.4** (Assign by weak update).

$$221 \quad \mathbb{S}^\sharp[x^{\text{strong}} \leftarrow e]\sigma^\sharp = \mathbb{S}^\sharp[x \leftarrow \text{expanded}(e)]$$

$$222 \quad \mathbb{S}^\sharp[x^{\text{weak}} \leftarrow e]\sigma^\sharp = \mathbb{S}^\sharp[\text{fold}(x, tmp)] \circ \mathbb{S}^\sharp[tmp \leftarrow \text{expanded}(e)]\sigma^\sharp$$

223 $\mathbb{S}^\sharp[x^{\text{strong}} \leftarrow e]$ means that the value of x is in the set of values that expression e evaluates
 224 into. When the right-hand side e contains weak variables, the $\text{expanded}(e)$ models reading
 225 from a collection by replacing every occurrence of a weak variable y in expression e with a
 226 temporary strong variable tmp , initialized with $\text{expand}(y, tmp)$. Given $x \in \mathbb{V}^{\text{strong}}$, $\sigma^\sharp \in \Sigma^\sharp$,
 227 $\mathbb{S}^\sharp[\text{remove}(x)]\sigma^\sharp$ removes the variable x from the abstract states Σ^\sharp . This statement is
 228 used to remove temporary variables after the assignment. This prevents inferring unsound
 229 relations: the assignment $x^{\text{strong}} = y^{\text{weak}}$ means that the value of x is the value of one of the
 230 concrete elements represented by y (i.e. $\sigma(x) \subseteq \sigma(y)$) whereas storing $x^{\text{strong}} = y^{\text{weak}}$ in a
 231 relational domain (e.g. the polyhedra domain) would mean that x and y have the same
 232 concretization (i.e. $\sigma(x) = \sigma(y)$).

233 Finally, if the left-hand side is a weak variable x^{weak} , then the assignment means that we
 234 add a possible value for x : we fold the right-hand side value into x .

235 ► **Example 2.5.** Consider the abstract environment $\sigma^\sharp = (y^{\text{weak}} : [1, 12])$. Then:

$$236 \quad \mathbb{S}^\sharp[x^{\text{strong}} \leftarrow 2 * y^{\text{weak}}]\sigma^\sharp = \mathbb{S}^\sharp[x \leftarrow \text{expanded}(2 * y^{\text{weak}})]\sigma$$

237 $\text{expanded}(2 * y^{\text{weak}})$ creates a strong temporary variable tmp , initialized with the same values
 238 as y^{weak} , which is removed after the assignment, i.e.:

$$239 \quad \mathbb{S}^\sharp[x^{\text{strong}} \leftarrow y^{\text{weak}}]\sigma^\sharp = \mathbb{S}^\sharp[\text{remove}(tmp)] \circ \mathbb{S}^\sharp[x^{\text{strong}} \leftarrow 2 * tmp](\sigma^\sharp[tmp : [1, 12]])$$

$$240 \quad = (x^{\text{strong}} : [2, 24], y^{\text{weak}} : [1, 12])$$

241 Those summary variables are used in the literature, to abstract strings [26], numerical
 242 trees [19] or algebraic data types [31]. They can be used in combination with a relational
 243 domain. For example, the polyhedra domain can express that a weak variable x is lower
 244 than some strong variable y , meaning that any value represented by x is lower than y (e.g.
 245 $a - 1 \leq x_c \leq a + 2$ in the example). However, it cannot infer relations between the content
 246 of two weak variables. Indeed, consider the semantics of the assign: the weak variables from
 247 the right-hand side are expanded to temporary fresh strong variables before performing an
 248 assign. Ultimately, as a limitation, it cannot express how the content of a container relates to
 249 the content of another container. The difficulty to express relations between weak variables
 250 was underlined by Siegel and Simon [30]. This very problem will be overcome by our DelExp
 251 domain (Section 3).

252 2.4 Container Abstraction

253 Weak variables can be leveraged to abstract containers in various contexts.

254 **Arrays.** Weak variables can be used to abstract arrays or lists in any language. To any
 255 array a , we associate a weak variable a^{weak} which summarizes every object contained in a .
 256 As seen in the previous section, the operations of array manipulation derive from `fold` and
 257 `expand` operators. When updating the cell of an array a with the expression e , we compute
 258 $a^{\text{weak}} \leftarrow e$, which folds inside t^{weak} the value of e . When reading the n -th cell $a[n]$ of a inside a
 259 variable x , we perform $x \leftarrow t^{\text{weak}}$, which expands the values of t^{weak} in x . This works for static
 260 arrays of fixed size (as in C), as well as dynamic arrays with changing, unbounded size, such
 261 as Python lists.

262 In particular, weak variables can be used to abstract Python lists, as highlighted by this
 263 example (abstraction in comment):

```
a = [12,3,7] #t_c: [3,12]
a.append(15) #t_c: [3,15]
h = a[0] #h: [3,15]
```

264 **OCaml Algebraic Data Types.** Valnet et al. [31] defines an abstraction for algebraic data
 265 types based on weak variables. Given any *user-defined* algebraic data type, they associate
 266 one weak variable per non-recursive construct field:

```
type t = C1 of t1,1 * ... t1,n | ... | Cm of tm,1 * ... * tm,n
```

267 The field (i,j) of type $t_{i,j}$ of a variable x is then abstracted by a weak variable called
 268 $x.Ci.j$. Therefore, all the elements of the field (i,j) recursively nested inside the variable x
 269 are smashed (i.e. folded) together inside the variable $x.Ci.j$. They also record the possible
 270 variants Ci used by each variable x , which thus limits the set of weak variables $x.Ci.j$ used
 271 to abstract x . Consider the following tree:

```
type tree = Node of tree * int * tree | Leaf
let t0 = Node(Node(Leaf,12,Leaf)3,Node(Leaf,7,Leaf)) (* t0.Node.1:[3,12], t0:{Node} *)
let t = Node(Leaf,15,t0) (* t.Node.1:[3,15], t:{Node} *)
let h = match t with Node(_,h,_) -> h (* h:[3,15] *)
```

272 All integer fields of `Node` recursively defined in `t0` are represented by `t0.Node.1`. The
 273 abstraction of this weak variable in the interval domain is $[3,12]$. $t0:\{Node\}$ means that `t0`
 274 starts with the variant `Node`.

275 Note that OCaml and Python examples share a similar precision loss. Indeed, when h is
 276 defined by `expand` of the weak variable, the abstract constraints on variable t are propagated
 277 to h , but we forget that h relates to t . With the polyhedra domain, we cannot store how
 278 h and t relate together, i.e. that h belongs to the set of elements that t represents (or,
 279 equivalently, that $\sigma(h) \subseteq \sigma(t)$). If we discover information about t afterwards, we will not be
 280 able to propagate them to h .

281 To overcome this limitation, the next section will define a domain able to store such a
 282 relation. The domain we will present is language-agnostic and manipulates weak variables
 283 independently of which set of objects they abstract. In the following, we will show examples
 284 both on Python lists and on OCaml ADTs.

285 3 The DelExp Domain

286 In this section, we define a domain able to express *relations* between weak variables. As a
 287 start, we want to express that a variable x verifies the same constraints as another variable
 288 y . Our domain is defined in order to be combined with existing analysis to improve their
 289 precision. Therefore, we suppose that an abstract semantics computing abstract environments
 290 in Σ^\sharp is already defined, of concretization $\gamma_\Sigma^\sharp : \Sigma^\sharp \rightarrow \mathcal{P}(\Sigma)$.

3.1 DelExp Domain

We define $\Sigma_{\mathbb{C}}^{\sharp} = \mathbb{V} \rightarrow \mathcal{P}(\mathbb{V}^{\text{weak}})$, a map from variables to a set of weak variables. Given $\mathbb{C} \in \Sigma_{\mathbb{C}}^{\sharp}$, if $x \in \mathbb{C}(v)$, it means that it concretizes to an environment σ where $\uparrow \sigma(v) \subseteq \uparrow \sigma(x)$, i.e. that v verifies all the constraints that x does, or, equivalently, that the elements represented by v form a subset of the elements represented by x . This is expressed by the following concretization:

► **Definition 3.1** (Concretization). *We define the concretization $\gamma : \Sigma^{\sharp} \times \Sigma_{\mathbb{C}}^{\sharp} \rightarrow \mathcal{P}(\Sigma)$. For $(\sigma^{\sharp}, \mathbb{C}) \in \Sigma^{\sharp} \times \Sigma_{\mathbb{C}}^{\sharp}$:*

$$\gamma((\sigma^{\sharp}, \mathbb{C})) = \{\sigma \in \gamma_{\Sigma}^{\sharp}(\sigma^{\sharp}) \mid \forall v \in \mathbb{V}, \forall x \in \mathbb{C}(v), \uparrow \sigma(v) \subseteq \uparrow \sigma(x)\}$$

We will write $x \in \mathbb{C}(v)$ as $v \prec x$. Considering the definition of **expand**, $\uparrow \sigma(v) \subseteq \uparrow \sigma(x)$ also means that $\sigma = \mathbb{S}[\mathbf{expand}(v, x)]\sigma$. This domain states the possibility to perform some expansions on the environment: therefore, we will call it the DelExp (Delay Expand) domain.

► **Example 3.2.** This domain is language agnostic: it operates directly on weak variables, independently of what they abstract. Suppose l is non-empty. Consider those OCaml and Python programs, both computing the tail of the list l :

```
let x = match l with | Cons(h, q) -> q | Nil -> Nil
```

```
x=l[1:]
```

In OCaml, we can express $\mathbf{x.Cons.1} \prec \mathbf{q.Cons.1}$. Similarly, in Python, we can express $x_0 \prec l_0$. Those are sound over-approximations of the final environments: the elements of the tail x of the list l satisfy every constraint that the elements of the list l do.

► **Definition 3.3** (Partial order). *We define the abstract partial order on \mathbb{C} :*

$$\mathbb{C} \subseteq_c \mathbb{C}' \iff \forall v \in \mathbb{V}, \mathbb{C}'(v) \subseteq \mathbb{C}(v)$$

That is, \mathbb{C} is more precise than \mathbb{C}' if it contains more constraints.

► **Example 3.4.** Given $\mathbb{V}^{\text{weak}} = \{x, y\}$ and $\mathbb{V}^{\text{strong}} = \{a\}$, we define $\mathbb{C} : \mathbb{V} \rightarrow \mathcal{P}(\mathbb{V}^{\text{weak}}) = x \rightarrow \{y\}; y \rightarrow \emptyset; a \rightarrow \{y\}$, meaning that x and a verify all the constraints that y does, and $\mathbb{C}' =: \mathbb{V} \rightarrow \mathcal{P}(\mathbb{V}^{\text{weak}}) = x \rightarrow \{y\}; y \rightarrow \emptyset; a \rightarrow \emptyset$, meaning that x verifies all constraints that y does. Then $\mathbb{C} \subseteq_c \mathbb{C}'$: \mathbb{C} expresses strictly more constraints on the environments than \mathbb{C}' , therefore it represents a smaller set of environments.

3.2 Operators

We now define the operators of our abstract domain. Note that before performing those operations, as is standard in relational domains (e.g. polyhedra domain), we perform the transitive closure of constraints under the rule: $x \prec y \wedge y \prec z \implies x \prec z$.

► **Definition 3.5** (Join). *We define the join \sqcup_c on $\Sigma_{\mathbb{C}}^{\sharp}$: $\mathbb{C} \sqcup_c \mathbb{C}' = v \mapsto \mathbb{C}(v) \cap \mathbb{C}'(v)$*

When performing a join, we need to keep constraints present in both environments.

► **Property 3.1.** \sqcup_c is a sound over-approximation of \cup and is a widening (because \mathbb{V} and $\mathcal{P}(\mathbb{V}^{\text{weak}})$ are finite).

325 ▶ **Definition 3.6** (Meet). We define the meet \sqcap_c on Σ_C^\sharp : $\mathbb{C} \sqcap_c \mathbb{C}' = v \mapsto \mathbb{C}(v) \cup \mathbb{C}'(v)$

326 Similarly to the join, when performing the meet, we keep all the constraints from both
327 environments: they are satisfied in the result environment.

328 ▶ **Property 3.2.** \sqcap_c is a sound over-approximation of \cap .

329 3.3 Transfer Functions

330 We update our abstract semantics with the transfer functions for this domain.

331 **Discovering constraints.** First, we define the `add_expand` function (Algorithm 1) which,
332 given a variable x , an expression e , and a DelExp environment \mathbb{C} , returns an environment in
333 which the constraint $x \prec y$ is added if e is a variable y .

■ Algorithm 1 `add_expand`

```

Input :  $x, e, \mathbb{C} \in \mathbb{V} \times \mathcal{E} \times \mathbb{C}$ 
if  $e = y \in \mathbb{V}$  then
  | if  $x \in \mathbb{V}^{strong}$  then
  | | return  $\mathbb{C}[x \rightarrow \{y\}]$ ;
  | else
  | | return  $\mathbb{C}[x \rightarrow \mathbb{C}(x) \cap \{y\}]$ ;
  | end
else
  | return  $\mathbb{C}$ ;
end

```

334 We then update the semantics of `expand`:

335 $S^\sharp \llbracket \text{expand}(x, e) \rrbracket (\sigma^\sharp, \mathbb{C}) = S^\sharp \llbracket \text{expand}(x, e) \rrbracket \sigma^\sharp, \text{add_expand}(x, e, \mathbb{C})$

336 When we expand a variable x with an expression $e = y \in \mathbb{V}$, then x verifies the same
337 constraints as y (Definition 2.2). Consequently, we add $x \prec y$ in the DelExp domain.

338 **Reduction.** We now define a reduction function to propagate information from our DelExp
339 domain to underlying domains. First, we define the `apply_constraint` function (Algorithm 2)
340 which, given a constraint $x \prec y$, applies it to the environment: to do so, it enforces that x is
341 the expansion of y in the current environment. Then, the `reduce` function (Algorithm 3)
342 takes as input an abstract environment σ^\sharp and a DelExp "to-expand" environment \mathbb{C} . It
343 returns a refined version of σ^\sharp . To do so, it applies sequentially each constraint in \mathbb{C} to
344 update σ^\sharp . Thanks to this reduction, other domains can benefit from information stored in
345 the DelExp domain. Note that this reduction is agnostic with respect to those underlying
346 domains.

■ Algorithm 2 `apply_constraint`

```

Input :  $\sigma^\sharp, x \prec y$ 
 $\sigma^\sharp \leftarrow E^\sharp \llbracket \text{expand}(x, y) \rrbracket \sigma^\sharp$ ;
return  $\sigma^\sharp$ 

```

Algorithm 3 reduce

```

Input :  $\sigma^\#, \mathbb{C}$ 
for  $x \prec y \in \mathbb{C}$  do
  |  $\sigma^\# \leftarrow \text{apply\_constraint}(\sigma^\#, x \prec y)$ 
end
return  $\sigma^\#$ 

```

347 A possible interpretation for this analysis is that it delays the application of the `expand`
 348 operator. When we uncover information about variable y , knowing that $x \prec y$, we can
 349 propagate them with the `expand` operator.

350 ► **Example 3.7.** Given programs from Example 3.2 computing the tail of a non-empty list
 351 l , the OCaml analysis will perform `q.Cons.1` \leftarrow `l.Cons.1` then `x.Cons.1` \leftarrow `q.Cons.1` using
 352 `expand`, then remove intermediate variable q , therefore yielding `x.Cons.1` \prec `l.Cons.1`. Similarly,
 353 the Python analysis will perform $x_0 \leftarrow l_0$, therefore yielding $x_0 \prec l_0$. The inference of those
 354 constraints works at the level of weak variables, independently of the underlying containers.

4 Inclusions up to Transformations

356 The DelExp domain presented so far yields precise results when a variable verifies the same
 357 constraints as another variable (e.g. when a list is a sub-list of another). Therefore, the
 358 DelExp domain can then be seen as a simplification of the QUIC graphs domain [11], which
 359 can express constraints over numerical sets, under the form $\bigcap X_i \subseteq \{\nu \in \bigcup Y_j \mid B[\nu]\}$, B
 360 being an abstract domain acting as a predicate on ν . However, many functions on containers
 361 apply a *transformation* to its elements. As is, our DelExp domain, as well as QUIC graphs,
 362 fail to capture a relation “up to transformation” between two variables. This section proposes
 363 a systematic methodology to extend the expressivity of the DelExp domain in this direction.

4.1 AffineDelExp Domain

365 We recall the Python program from the introduction:

```

# assumes l:List[int]
for x in l:
  r.append(2*x)

```

366 Given a non-empty list l , it builds r as the list l where every element was multiplied by
 367 2. The standard polyhedra domain on weak variables, the QUIC graphs domain [11] or our
 368 current DelExp domain from Section 3 all fail to capture a relation between l and r the input
 369 and the output lists. We would like to express constraints under the form $x \prec 2y$ to express
 370 that x respects every constraint that $2y$ does, i.e. that the set of values represented by x
 371 is included in the set of values in y multiplied by 2. Given a set X , we denote as $P_f(X)$ the set
 372 of finite subsets of X . We define an *affine* DelExp domain able to express such constraints:

373 ► **Definition 4.1** (Abstract domain). We define $\Sigma_A^\# = \mathbb{V} \rightarrow \mathcal{P}_f(\{aX + b \mid a, b \in \mathbb{Z}, X \in \mathbb{V}\})$.

374 This domain, that we will name AffineDelExp, stores constraints under the form $x \prec ay + b$.
 375 Note that, like the standard polyhedra domain, this domain ignores possible overflows (there
 376 are none in Python).

377 ► **Definition 4.2** (Concretization). *We define the concretization $\gamma : \Sigma^\sharp \times \Sigma_A^\sharp \rightarrow \Sigma$. For*
 378 *$(\sigma^\sharp, \mathbb{C}) \in \Sigma^\sharp \times \Sigma_A^\sharp$:*

$$379 \quad \gamma((\sigma^\sharp, \mathbb{C})) = \{\sigma \in \gamma_\Sigma^\sharp(\sigma^\sharp) \mid \forall x \in \mathbb{V}, x \prec ay + b \in \mathbb{C} \implies \uparrow \sigma(x) \subseteq a \times \uparrow \sigma(y) + b\}$$

380 Consequently, if $x \prec 2y \in \mathbb{C}$, it means that in each environment, the set of elements in the
 381 container summarized by x is a subset of the set of elements in the container summarized
 382 by y multiplied by 2. The order relation, join, meet and widening defined in Section 3
 383 can be left unchanged, up to the update of the transitive closure. In this new domain,
 384 the transitive closure, needed for the precision of the operators, is computed with the rule:
 385 $x \prec ay + b \wedge y \prec cz + d \implies x \prec acz + bc + d$.

386 We then need to update how we discover constraints and how we perform the reduction
 387 applying them. This is done by trivial modifications on `add_expand` and `apply_constraint`:
 388 we store $x \prec ay + b$ when we encounter $\mathbb{E}^\sharp[\text{expand}(x, ay + b)]\sigma^\sharp$ and we apply $\mathbb{E}^\sharp[\text{expand}(x, ay +$
 389 $b)]\sigma^\sharp$ to perform a reduction. Consequently, our analysis is able to automatically discover as
 390 final state for the Python programs $r_0 \prec 2 \times l_0$, i.e. the list r is included in a transformation
 391 of the list l where each element is multiplied by 2.

392 4.2 General Methodology

393 Section 4.1 demonstrates that only a few modifications need to be applied to DelExp
 394 to express relations up to affine transformations. This section aims at generalizing this
 395 methodology. Given a certain set of transformations Θ (e.g. affine transformations,
 396 polynomial transformation, etc.), we want to be able to infer constraints under the form
 397 $x \prec \theta(y)$ for θ in Θ , meaning that x verifies the same constraints that the set $\theta(y)$, i.e. that
 398 the elements of y to which we applied the transformation θ . To do so, we build an abstract
 399 domain of the form $\Sigma_\Theta^\sharp = \mathbb{V} \rightarrow \mathcal{P}_f(\{\theta(X) \mid X : \tau_1 \in \mathbb{V}, \theta : \tau_1 \rightarrow \tau_2 \in \Theta\})$.

400 ► **Definition 4.3** (General concretization). *The concretization is of the form:*

$$401 \quad \gamma((\sigma^\sharp, \mathbb{C})) = \{\sigma \in \gamma_\Sigma^\sharp(\sigma^\sharp) \mid \forall \theta : \tau_1 \rightarrow \tau_2 \in \Theta, \forall x : \tau_2, y : \tau_1 \in \mathbb{V},$$

$$402 \quad x \prec \theta(y) \in \mathbb{C} \implies \sigma(x) \subseteq \theta(\sigma(y))\}$$

403 As for affine transformations, we need to update how we discover constraints and how we
 404 perform the reduction applying them. We only need to perform those three steps:

- 405 ■ Soundly update `add_expand` definition to decide when to add a constraint $x \prec \theta(y)$.
- 406 ■ Soundly update `reduce` definition to decide how to enforce this constraint in the abstract
 407 environment.
- 408 ■ Soundly modify how the transitive closure is computed (depending on Θ).

409 Defining a new domain on the set of transformations Θ (e.g. polynomial transformations)
 410 then permits to precisely express an element-wise transformation from this set. For instance,
 411 Θ can be the set of polynomial transformation. This methodology is at the heart of our
 412 domain expressivity, and will be further developed in the context of OCaml compositional
 413 analysis (see Section 7.1). Θ will then be the set of pure OCaml functions from the analyzed
 414 code, enabling us to express higher-order properties on programs.

415 5 Compositional Analysis for a Pure Functional Language

416 This section provides background on a compositional analysis for a pure monomorphic
 417 functional language introduced by Valnet et al. [31]. Then, it shows that the precision

$$\begin{aligned}
e_1, e_2, \dots \in \mathcal{E} ::= & x \in \mathbb{V} \\
& | n \in \mathbb{Z} & | e_1 \oplus e_2 \\
& | e_0 \ e_1 \ \dots \ e_n & | \text{fun } x_1 \dots x_n \rightarrow e \\
& | \text{let } x = e_1 \text{ in } e_2 & | \text{let rec } x = e_1 \text{ in } e_2 \\
& | C(e_1, \dots, e_n) & | \text{match } e_0 \text{ with } p_1 \rightarrow e_1 \mid \dots \mid p_m \rightarrow e_m \\
& | \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \\
p_1, p_2, \dots \in \Pi ::= & _ \mid x \in \mathbb{V} \mid n \in \mathbb{Z} \mid C(p_1, \dots, p_n) \\
& \text{where } \forall i \neq j, fv(p_i) \cap fv(p_j) = \emptyset
\end{aligned}$$

■ **Figure 8** Syntax of the functional language.

$$\begin{aligned}
\mathcal{V} ::= & \mathbb{Z} \cup \{ C(v_1, \dots, v_n) \mid C \in \mathbb{C}, v_i \in \mathcal{V} \} \cup \bigcup_{n \in \mathbb{N}} [\mathcal{V}^n \rightarrow \mathcal{V}] \\
\Sigma = & \mathbb{V} \rightarrow \mathcal{V} \\
\mathbb{E}[\cdot] : & \mathcal{E} \times \Sigma \rightarrow \mathcal{V}^\perp
\end{aligned}$$

■ **Figure 9** Semantic domain of the functional language.

418 issue of relational domains on weak variables prevents this analysis from generating precise
419 summaries. We will extend the DelExp domain in Sec. 6 to combine it with this analysis
420 and generate precise summaries for this language in Sec. 7

421 5.1 Syntax and Semantics

422 Figure 8 details a pure monomorphic functional language without side-effects. Its key features
423 are recursion, algebraic data types and higher-order. \mathcal{E} is the set of expressions. $fv(e)$ is the
424 set of free variables of expression e . Π , the set of patterns, contains the wildcard $_$, variables,
425 integers, and type variants. Free variables from the same pattern are required to be disjoint,
426 as is the case in OCaml: the same variable cannot be bound twice.

427 Figure 9 describes the semantic domain of the language. \mathbb{C} is the set of all possible variants.
428 Semantic values of the language \mathcal{V} are integers, type variants containing values, and continuous
429 functions from values to values. Finally, Σ is the set of concrete environments, associating a
430 value to each variable. The concrete semantics of an expression e is $\mathbb{E}[e] \in \Sigma \rightarrow \mathcal{V}^\perp$, where \perp
431 denotes non-termination. It associates a semantic value to an expression in an environment.
432 This semantics follows the standard for eagerly evaluated functional languages (such as
433 OCaml). We suppose that programs are well-typed according to usual type inference on
434 functional languages [16, 24]. Note that this language is monomorphic so far. This limitation
435 will be overcome later on thanks to our DelExp domain.

436 5.2 Compositional Analysis

437 Compositional analysis aims at over-approximating at definition site all the possible
438 behaviours of a function. The function is analyzed once and for all: we then store the result
439 of the analysis, its *summary*. Compositional analyses significantly improve scalability when
440 a function is called multiple times, as summary applications are in general computationally

```

let a = random 5 10 (* environment before *)
let mult2 x = 2*x   (* function definition *)
let b = mult2 a     (* function application *)

```

■ **Figure 10** Small example to illustrate compositional analysis.

441 cheaper than reanalyzing a function at every calling context [5, 8]. Therefore, it can help the
 442 analysis to scale up. Besides, compared to non-compositional analyses, they can be used to
 443 infer functions specifications.

444 Since our language (Figure 8) is pure (no side effects), the function does not modify the
 445 environment: its behavior can be represented by a relation between its inputs and its output.
 446 Bautista et al. [2] formalize a so-called folk technique to generalize this relation with side
 447 effects at first-order.

448 We recall the notations of Valnet et al. [31], which we will extend later on. Given a
 449 function type $\tau = \tau_1 \rightarrow \dots \rightarrow \tau_{n+1}$, we denote as \mathcal{F}^τ the set of functions of type τ . We
 450 can then *choose* an abstract domain $\mathcal{R}(V)$ of concretization $\gamma_{\mathcal{R}} : \mathcal{R}(V) \rightarrow \Sigma|_V$ to represent
 451 relations between input and output variables V : it abstracts $\Sigma|_V$, i.e. concrete environments
 452 restricted to V . The domain for functions of type τ is then $\mathcal{D}_\tau = \mathbb{V}^n \times \mathbb{V} \times \mathcal{R}(V)$. A function
 453 is abstracted as: the names x_1, \dots, x_n of its formal inputs, a result variable, and a relation
 454 between those variables. A key advantage of this formalization is that it is *domain-modular*.
 455 This framework can then be re-used as is in our analysis with our DelExp abstract domain
 456 as relational domain.

457 ► **Definition 5.1** (Functions concretization). *The concretization $\gamma_\tau : \mathcal{D}_\tau \rightarrow \mathcal{F}^\tau$ is:*

$$458 \quad \gamma_\tau(\langle (x_1, \dots, x_n), r, p \rangle) = \{f : \tau \mid \forall (v_1, \dots, v_n) : \tau_1 \times \dots \times \tau_n, \\ 459 \quad \quad \quad [x_1 = v_1] \dots [x_n = v_n][r = f(v_1, \dots, v_n)] \in \gamma_{\mathcal{R}}(p)\}$$

460 A function f is abstracted as $\langle (x_1, \dots, x_n), r, p \rangle$ if the relation between its inputs and its
 461 output can be described by p . We call this abstraction the *function summary* and denote it
 462 as $\lambda x_1 \dots x_n. p(x_1, \dots, x_n, r)$.

463 ► **Example 5.2** (Concretization). Given `mult2` of type $\tau = \text{int} \rightarrow \text{int}$ and $V = \{x, r\}$, we
 464 choose $\mathcal{R}(V)$ as the polyhedra domain on V . A sound over-approximation for `mult2` is then
 465 $\langle x, r, r = 2x \rangle$, i.e. $\lambda x. r = 2x$.

466 Sound operations on this domain are defined by Valnet et al. [31], together with a domain
 467 for *disjunctive* summaries (allowing partitioning on the inputs) and transfer functions to
 468 analyze (possibly recursive) functions and apply those summaries.

469 ► **Example 5.3** (Analysis of `mult2`). As an example of function analysis, we analyze `mult2` :

$$\begin{array}{l}
 \text{S}^\# \llbracket \text{let } mult2 = \text{fun } x \rightarrow 2 * x \rrbracket [.] \\
 \left\{ \begin{array}{l}
 \rightarrow \text{E}^\# \llbracket \text{fun } x \rightarrow 2 * x \rrbracket [.] \\
 \quad \left\{ \begin{array}{l}
 \rightarrow \text{E}^\# \llbracket 2 * x \rrbracket \sigma^\# [x = \top] \\
 \quad \leftarrow 2 * x, [x = \top] \\
 \leftarrow \langle x, r, r = 2 * x \rangle, [.] \\
 \leftarrow [mult2 \rightarrow \langle x, r, r = 2 * x \rangle]
 \end{array} \right.
 \end{array} \right.
 \end{array}$$

470

471 This analysis trace shows the analysis steps. The right arrow means that we analyze a
 472 sub-expression. The left arrow means that a result is returned. `mult2` is analyzed once and
 473 for all at *definition site*. Consequently, we compute the abstract semantics of the body with
 474 no hypothesis on input values: x is initialized at \top . The generated relation is then valid for
 475 all inputs. Note that relationality is critical there. It enables us to discover information on
 476 the function behavior despite making no hypothesis on the values of the input x .

477 ► Remark 5.4 (Recursive functions). Note that the analysis of recursive functions is done via
 478 fixpoint iteration: starting from a function with bottom input-output relation, we iteratively
 479 complete the function summary by reanalysing the function, until reaching a fixpoint, and use
 480 a widening to converge in finite time.

481 ► Example 5.5 (Application of `mult2`). As an example of function application, we analyze:

$$\begin{array}{l}
 \mathbb{E}^\# \llbracket \text{mult2 } (7 + 1) \rrbracket \sigma^\# \\
 \begin{array}{l}
 \longrightarrow \mathbb{E}^\# \llbracket \text{mult2} \rrbracket \sigma^\# \\
 \longleftarrow \langle x, r, r = 2 * x \rangle, \sigma^\# \\
 \longrightarrow \mathbb{S}^\# \llbracket x = 7 + 1 \rrbracket \sigma^\# \\
 \longleftarrow \sigma^\# [x \rightarrow 8] \\
 \longleftarrow r, \sigma^\# [x \rightarrow 8, r \rightarrow \top] \sqcap \sigma^\# [x \rightarrow \top, r = 2x] = \\
 \longleftarrow r, \sigma^\# [x \rightarrow 8, r \rightarrow 16]
 \end{array}
 \end{array}$$

482

483 First, we compute the summary of the applied function. Since `mult2` was already analyzed,
 484 its summary is stored in the environment. We only need to instantiate this summary with
 485 input values: we assign formal parameters to the input values in the current environment
 486 and perform its intersection with the summary to get the result $r = 16$.

487 6 DelExp with Heterogeneous Environments

488 Our analysis manipulates summarization variables to abstract contents from containers, e.g.
 489 to abstract the content of a list. However, those containers may be empty: as a consequence,
 490 the summarization variables may correspond to no concrete objects. Such variables are called
 491 *optional* and must be handled with extra care. In particular, they require reformulating
 492 standard concretizations.

493 6.1 Heterogeneous Environments

494 We introduce here the formalization from Journault et al. [19]. Together with a standard
 495 abstract environment $\sigma^\#$, they define two sets of variables $l, u \in \mathcal{P}(\mathbb{V})$, $l \subseteq u$. l represents the
 496 set of non-optional variables, i.e. variables existing for sure, in every concretize environments.
 497 u represents the set of all possibly existing variables in the program. Therefore, variables
 498 in $u \setminus l$ are optional: they may not exist in some concretized environments. Those variable
 499 sets are soundly inferred throughout the analysis of the programs. We denote as $dom(\sigma)$
 500 the definition domain of the environment σ . Given $\sigma' \in \Sigma$, $\sigma'_{|dom(\sigma)}$ restricts the definition
 501 domain of σ' to $dom(\sigma)$. Given a set of environments S , $S_{|dom(\sigma)}$ lifts this restriction to
 502 every environment of the set. The new concretization of abstract environments is:

503 ► Definition 6.1 (Heterogeneous concretization). We define the concretization:

$$504 \quad \gamma_{\mathcal{H}}((\sigma^\#, l, u)) = \{ \sigma \mid l \subseteq dom(\sigma) \subseteq u \wedge \sigma \in \gamma_{\mathbb{Z}}(\sigma^\#)_{|dom(\sigma)} \}$$

505 We denote as $\sqcup_{\mathcal{H}}$ (resp. $\sqcap_{\mathcal{H}}$) the join (resp. the meet) on heterogeneous environments,
 506 defined in [19]. As shown in the introduction, summarization variables of ADTs fields may
 507 not exist in some executions. Therefore, their abstraction, from [31], is based on careful
 508 manipulation of those optional variables.

509 ► **Example 6.2** (Heterogeneous summaries). We illustrate the difficulty of handling heterogeneous
 510 environments:

```

type either = Left of int | Right of int
let to_int a =
  match a with
  | Left n -> n
  | Right n -> n
let b = if rand () then Left 5 else Right 6
let x = to_int b
```

511 A sound heterogeneous summary is $\langle a, r, a : \{\text{Left}, \text{Right}\}, r = \text{a.Left} = \text{a.Right} \rangle$ with
 512 $l = \{r\}$ and $u = \{r, \text{a.Left}, \text{a.Right}\}$. The heterogeneous relational constraint $r = \text{a.Left} =$
 513 a.Right where r is always defined and a.Left and a.Right are optional variables represents
 514 multiple cases of homogeneous relational constraints. It includes the case where only r and
 515 a.Left exists and $r = \text{a.Left}$, and the case where only r and a.Right exists and $r = \text{a.Right}$.
 516 Note that however, it loses the information that exactly one variable among a.Left and
 517 a.Right exists: it also abstracts a concrete environment where a.Left and a.Right both exist
 518 and $r = \text{a.Left} = \text{a.Right}$. This can be recovered by the fact that **Either** is a non-recursive
 519 ADT.

520 A sound abstract value for b is $\{\text{Left}, \text{Right}\}, \text{b.Left} = 5, \text{b.Right} = 6$ with $l = \emptyset$ and
 521 $u = \{\text{b.Left}, \text{b.Right}\}$. If we naively compute the intersection when applying the function,
 522 we end up with $x = \text{b.Left} = 5 = \text{b.Right} = 6$, i.e. \perp . This is unsound: we propagate
 523 information from values that may not exist as if their existence was certain. A sound
 524 intersection was proposed in Journault et al. [19].

525 6.2 Heterogeneous Version of DelExp

526 The DelExp domains (Secs. 3 and 4) presented so far in this section are only defined
 527 on *homogeneous* environments: the operators (join, meet, widening) are defined between
 528 environments on the same set of non-optional variables. To use it in combination with the ADT
 529 domain from Valnet et al. [31], we need to provide support for heterogeneous environments.
 530 To express relations between optional variables, we have to adapt its concretization by
 531 providing it with the lower and upper sets of defined variables l and u .

532 ► **Definition 6.3** (Heterogeneous concretization of DelExp). *We define the concretization*
 533 $\gamma : \Sigma^{\#} \times \mathcal{P}(\mathbb{V}) \times \mathcal{P}(\mathbb{V}) \times \mathbb{C} \rightarrow \Sigma$ *on the standard DelExp domain (Sec. 3):*

$$534 \quad \gamma((\sigma^{\#}, l, u, \mathbb{C})) = \{\sigma \in \Sigma \mid \sigma \in \gamma_{\mathcal{H}}(\sigma^{\#}, l, u) \quad (1)$$

$$535 \quad \wedge \forall v \prec x \in \mathbb{C}, v \in \text{dom}(\sigma) \implies x \in \text{dom}(\sigma) \quad (2)$$

$$536 \quad \wedge v \in \text{dom}(\sigma) \implies \sigma(v) \subseteq \sigma(x) \quad (3)$$

537 The concretization requires that if a constraint $v \prec x$ exists then it only represents
 538 environments σ where $v \in \text{dom}(\sigma) \implies x \in \text{dom}(\sigma)$. Adding this precision is a *choice* that
 539 modifies the semantics of the domain. Without this restriction, we could only propagate

540 information from the constraint $v \prec x$ if we are sure that v and x exists. With it, we only
 541 need to make sure that x exist. Note that to be sound, an analysis may only add a constraint
 542 $v \prec x$ in the environment if the existence of v indeed guarantees that x exists.

543 ► **Property 6.1.** *The unchanged meet $\sqcap_c^{\mathcal{H}} = \sqcap_c$ is sound.*

544 ► **Definition 6.4** (Heterogeneous join). *We define the join $\sqcup_c^{\mathcal{H}}$ on \mathbb{C} :*

$$\begin{aligned}
 545 \quad (f_1, l_1, u_1, \mathbb{C}_1) \sqcup_c^{\mathcal{H}} (f_2, l_2, u_2, \mathbb{C}_2) &= ((f_1, l_1, u_1) \sqcup_{\mathcal{H}} (f_2, l_2, u_2), \\
 546 \quad &\{x \prec y \in \mathbb{C}_1 \mid x \notin u_2\} \cup \\
 547 \quad &\{x \prec y \in \mathbb{C}_2 \mid x \notin u_1\} \cup (\mathbb{C}_1 \cap \mathbb{C}_2))
 \end{aligned}$$

548 ► **Property 6.2.** $\sqcup_c^{\mathcal{H}}$ is sound.

549 ► **Example 6.1.** *We compute $(f_1, \{a, x_w, y_w\}, \{a, x_w, y_w\}, \{x_w \prec y_w\}) \sqcup_c^{\mathcal{H}} (f_2, \{a, y_w\}, \{a, y_w\}, \top)$.
 550 Since x_w only exists in the left component, we can keep $x_w \prec y_w$ in the final result. This is
 551 due to the semantics of this constraint: it only applies to environments where x_w exists. We
 552 can safely keep it: it will not propagate information to environments where x_w does not exist.
 553 Therefore, the formula yields $(f_1 \sqcup_{\mathcal{H}} f_2, \{a, y_w\}, \{a, x_w, y_w\}, \{x_w \prec y_w\})$.*

■ **Algorithm 4** `apply_constraintH`

```

Input:  $(f, l, u, \mathbb{C}), x \prec y$ 
 $\sigma^{\#} \leftarrow (f, l, u, \mathbb{C})$ 
if  $x \in l$  then
  |  $\sigma^{\#} \leftarrow \mathbb{S}^{\#}[\text{expand}(x, y)](f, l, u)(f, l, u), \mathbb{C}$ 
end
else
  |  $\sigma^{\#} \leftarrow \mathbb{S}^{\#}[\text{remove}(x)](f, l, u) \sqcup \mathbb{S}^{\#}[\text{expand}(x, y)](f, l \cup \{x, y\}, u), \mathbb{C}$ 
end
return  $\sigma^{\#}$ 

```

554 The heterogeneous reduction, `reduceH`, is an update of `reduce` (Algorithm 3) where
 555 `apply_constraint` is replaced with `apply_constraintH`. `apply_constraintH` behaves like
 556 `apply_constraint` when the variables are non-optional. Otherwise, if x is optional, we
 557 perform a join between environments where x does not exist (in which case we perform
 558 `remove(x)`) and environments where it is non-optional (in which case we perform `expand(x, y)`).

559 ► **Property 6.3.** `reduceH` is sound.

560 Those operators are proven sound with regard to the heterogeneous concretization. This
 561 work on heterogeneous environments allows the DelExp domain to work with optional
 562 variables, i.e. variables that may not exist. A similar formalization can adapt AffineDelExp
 563 to heterogeneous environments. In particular, it makes this domain compatible with the state-
 564 of-the-art relational abstractions for numerical trees [19] and for Algebraic Data Types [31].

565 ► **Example 6.5** (OCaml Algebraic Data Types). Given the OCaml function:

```

let rec filter_le l inf =
  match l with
  | Cons(h,q) ->
    if h > inf then filter_le q inf

```

```

else Cons(h,filter_le q inf)
| Nil -> Nil

```

566 The analysis can infer the summary:

567 $\lambda l \text{ inf}. r : \{\text{Nil}, \text{Cons}\} \wedge l = \{\text{Nil}, \text{Cons}\} \wedge (r_0 < \text{inf}, \{\}, \{r_0, l_0\}) \wedge r_0 < l_0$

568 The concretization of the summary means: “the result of `filter_le` may be empty; if it
569 is not, then the input list is not empty either, and the content of the output list is lower than
570 `inf` and included in the input list”. This is more precise than the summary for this function
571 in Valnet et al. [31], i.e. this summary without the constraint $r_0 < l_0$.

572 7 Case Study: Functional Compositional Analysis Over ADTs

573 This section shows how the domains introduced in Secs. 3 and 4 can be used to perform
574 a precise analysis in the presence of containers. We extend the analysis from Valnet et al.
575 [31] on pure monomorphic ADT-manipulating OCaml programs described in Sec. 5.1. We
576 highlight key design choices that make this combination possible.

577 **Heterogeneity.** The reduction of DelExp domains is sound when in product with a domain
578 that implements sound operations for heterogeneous environments. Non-relational domains
579 can be lifted without efforts to heterogeneous environments [19]. Defining heterogeneous
580 operations for relational domains requires more care. The OCaml analysis we built upon only
581 supports two relational domains: the polyhedra domain and the equality domain (equalities
582 are only stated between variables that are sure to exist together). Both of them implement
583 sound heterogeneous operators.

584 **Modularity.** Second, note that those domains are agnostic in the domain chosen to abstract
585 underlying variables: consequently, they can work in combination with any domain. Valnet
586 et al. [31] use summarization variables to represent Algebraic Data Types fields in a similar
587 modular fashion: the ADT domain is agnostic in the domain chosen to abstract them. Thanks
588 to those parametric definitions, they can be combined, and used in a higher-order settings.

589 **Polymorphism.** In fact, the default DelExp domain do not make any assumption on the
590 very type of the variables they manipulate. In particular, those variables can be *polymorphic*.
591 We store the constraints between polymorphic variables during the analysis of a polymorphic
592 function. When the function is applied to a known ground type τ , we can then perform the
593 reduction with the domain \mathcal{D}_τ . We therefore extend the existing OCaml analysis to support
594 polymorphism (e.g. the `map` function from Fig. 11).

595 7.1 ApplyDelExp: a Higher-Order Transformation

596 As a showcase of the methodology from Sec. 4.2, we will expand our DelExp domain with
597 pure higher-order transformation (without side-effect). We want to express constraints under
598 the form $x < f(y)$ with x , y , and f being variables of the language. We define as abstract
599 domain $\Sigma_{\text{fun}}^\# = \mathbb{V} \rightarrow \mathcal{P}_f(\{f(X) \mid \exists \tau_0, \tau_1, f \in \mathbb{V}_{\tau_0 \rightarrow \tau_1}, X \in \mathbb{V}_{\tau_1}\})$.

600 ■ We update `add_expand`: if the abstract expression $e^\#$ is (syntactically) $f(y)$ then we add
601 the constraint $x < f(y)$.

```

let rec map f l =
  match l with
  | Cons(h,q) -> Cons(f h, map f q)
  | Nil -> Nil

```

■ **Figure 11** Motivating example for our higher-order polymorphic transformation

- 602 ■ We update **reduce**: if the environment contains $x \prec f(y)$, we execute **expand**($x, f(y)$) in
603 the environment.
- 604 ■ We compute the transitive closure under the rule that $x \prec f(y)$ and $y \prec z$ implies that
605 $x \prec f(z)$.

606 Note that the reduction is sound in a *pure* fragment of the language, i.e. where functions
607 have no side-effects. We name this domain ApplyDelExp. It enables us to derive a precise
608 summary for **map**.

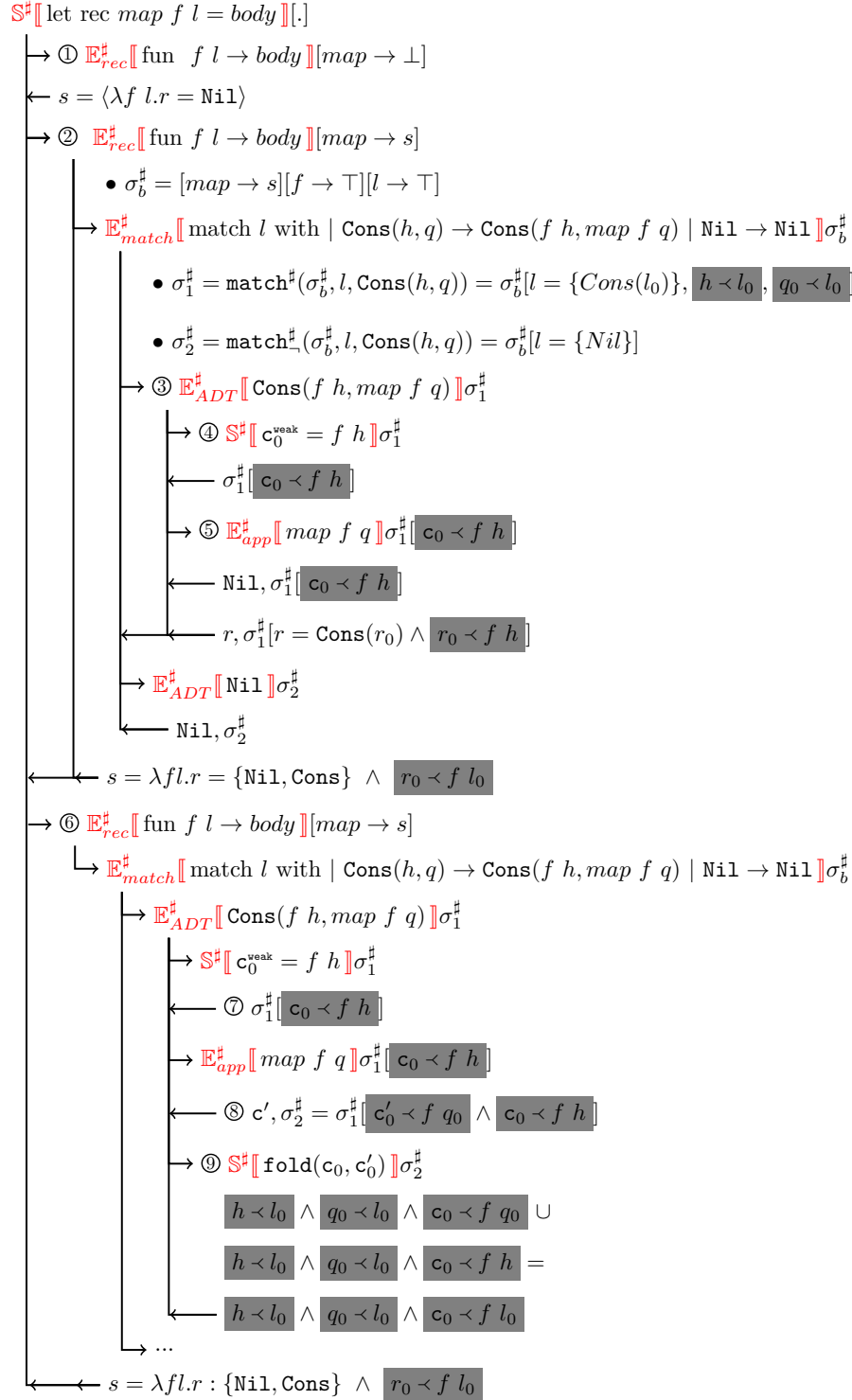
609 7.2 Analysis of map

610 The trace in Fig. 12 details the analysis of **map** from Fig. 11. Note that some details are
611 skipped for pedagogical purposes: not all steps are described and the sets l and u from
612 heterogeneous environments are not represented.

613 Given an OCaml list l , $l.\text{Cons}.1$ summarizes the elements of the integer fields of
614 the variant **Cons** nested inside l (see Sec. 2.4). For readability purposes, $l.\text{Cons}.1$ (resp.
615 $r.\text{Cons}.1$ and $q.\text{Cons}.1$) will be denoted as l_0 (resp. r_0 and q_0). c is a temporary variable
616 for **Cons**($f\ h, \text{map}\ f\ q$) and c' a temporary variable for its queue **map** $f\ q$. Constraints
617 from DelExp domains are highlighted in gray. The variable r corresponds to the result of
618 **map**. Depending on the expression e , $\mathbb{E}^\sharp[e]$ delegates the abstraction of e to the relevant
619 abstract domain (e.g. $\mathbb{E}_{ADT}^\sharp[e]$).

620 We start with an empty environment $[\cdot]$. As **map** is recursive, we need to perform fixpoint
621 iterations to infer its summary, as described in Remark 5.4. We first start by assuming
622 that **map**'s summary is \perp and analyze the body of **map** to get a new summary s . We then
623 reanalyze the body of **map** by assuming that s is the new summary of **map**, etc. until a
624 fixpoint is reached:

- 625 1. After the first iteration ①, the inferred behavior for **map** is that the input can only be
626 **Nil** and the output is then **Nil**.
- 627 2. We then start the second iteration with this hypothesis on **map**'s behaviour ②. **match**[♯]
628 computes in σ_1^\sharp an over-approximation of the environments in which l matches the first
629 pattern, whereas **match**_♯ computes in σ_2^\sharp an over-approximation of the environments in
630 which they do not. We then analyze **Cons**($f\ h, \text{map}\ f\ q$) in σ_1^\sharp ③: the head is $f\ h$
631 ④, which yields $c_0 \prec fh$. The tail ⑤ is empty (by application of the summary inferred
632 at the first iteration), so we do not fold any other constraint inside c_0 . When removing h
633 in the final summary and assigning $r = c$, we deduce $r_0 \prec fl_0$.
- 634 3. We start the third iteration with this new hypothesis on **map**'s behaviour ⑥. We focus
635 on the analysis of **Cons**($f\ h, \text{map}\ f\ q$). First, the analysis of the head derives $c_0 \prec fh$
636 ⑦. Then the queue c' (summarized by c'_0) of the list is **map** $f\ q$, which yields $c'_0 \prec fq_0$
637 ⑧. We then fold c'_0 inside c_0 ⑨: we join $c_0 \prec fh$ with $c_0 \prec fq_0$ knowing that $h \prec l_0$
638 and $h \prec q_0$, i.e. we only keep the constraints that are true in both DelExp environments.
639 We deduce $r_0 \prec fl_0$.



■ **Figure 12** Analysis of the `map` function with the DelExp domain. DelExp constraints are highlighted in gray. l_0 (resp. r_0, q_0, c_0, c_0') summarizes the elements of l (resp. r, q, c, c'), r is the result of `map` and s its summary.

640 4. Iterations 2 and 3 ended up with the same summary: we reached a fixpoint.

641 As a conclusion, the generated summary states: if the input list is empty, the output list
642 is empty. Otherwise, the output list content is included in the input list content to which
643 we applied the (pure) input function. The ApplyDelExp domain we created following the
644 methodology from Sec. 4.2 enhances the precision of the inferred summary. The ability to
645 derive “up-to-transformation” DelExp domains can help significantly improve compositional
646 analysis in the presence of containers. Note however that we did not prove the equality
647 between those two sets, nor that the order of elements is maintained in the resulting list. To
648 this end, we would need to express that h and q cover all elements in l : this cannot be done
649 in our DelExp domains.

650 **8 Experimental Evaluation**

651 **8.1 Mopsa**

652 The analysis described in this article is implemented in MOPSA [18]. MOPSA (Modular
653 Open Platform for Static Analysis) is a platform aiming at easing the development of new
654 static analyzers [25, 29]. This platform, written in OCaml, is open-source and multi-language.
655 It features language-agnostic abstractions such as integer domains (polyhedra, octagons,
656 intervals) and implements the transfer functions for standard constructs such as conditional
657 branching, `while` loops or assignments. The core of the platform is composed of 22 000
658 LoC of OCaml. Existing analyzers for C, Python, and OCaml base their implementation on
659 those components by adding additional syntactic structures, iterators, and language-specific
660 domains that complete the language-agnostic abstractions.

661 MOPSA traditionally encourages both relational abstractions and information sharing
662 between domains through powerful communication mechanisms (such as reduced products
663 or expression rewriting). Those design choices perfectly align with our domain-agnostic
664 relational domain for weak variables: DelExp was integrated seamlessly into MOPSA. We
665 implemented the basic DelExp domain from Sec. 3, the AffineDelExp from Sec. 4 and the
666 ApplyDelExp domain from Sec. 7.1 as a separate module of around 1000 LoC of OCaml. We
667 simply added it both to the OCaml compositional analysis from Valnet et al. [31] (restricted
668 to a pure and monomorphic fragment) and to the Python analysis from Monat et al. [27],
669 which uses weak variables to summarize lists and sets. It did not require any change to the
670 existing abstract domains and transfer function. Section 4.2 mentioned that `add_expand`
671 and `reduce` functions need to be redefined when formalizing a new transformation: the same
672 applies for the implementation. Once the basic DelExp domain was implemented, we simply
673 extended the `add_expand` and `reduce` functions to handle those affine and apply constraints.

674 The resulting analyses are fully automatic. They compute over-approximations of the
675 values of the program variables. In particular, the OCaml analysis infers functions summaries.

676 **8.2 Benchmarks**

677 All our tests were performed on a Intel(R) Core(TM) Ultra 7 165H CPU with 32 GB of RAM.
678 Execution time were computed as an average over 10 runs of the analysis. All programs
679 are between 6 and 20 lines and consist of a unique container-manipulating function. The
680 summaries have been visually checked in the execution trace of the analysis.

681 **OCaml.** As our OCaml analysis is compositional, we focused on assessing the precision of
682 the generated summaries on ADT-manipulating functions. Figure 13 presents our results.

Function	Delexp (ms)	Overhead	Summary
hd.ml	8	2.6	$r \prec l.Cons.1$
tl.ml (Ex. 3.2)	8	2.5	$r.Cons.1 \prec l.Cons.1$
filter.ml	46	2.3	$r.Node.1 \prec l.Node.1$
map.ml (Fig. 11)	48	3.4	$r.Cons.1 \prec f(l.Cons.1)$
copy.ml	56	3.0	$r.Cons.1 \prec l.Cons.1$
filter_le.ml (Ex. 6.5)	160	3.7	$r.Cons.1 \prec l.Cons.1 \wedge r.Cons.1 \prec inf$
mult2.ml (Fig. 2)	65	3.4	$r.Cons.1 \prec 2 \times l.Cons.1$
mult2tree.ml	99	3.2	$r.Node.1 \prec 2 \times l.Node.1$
mult3plus4.ml	67	3.6	$r.Cons.1 \prec 3 \times l.Cons.1 + 4$
maptree.ml	61	2.8	$r.Node.1 \prec f(l.Node.1)$
listtotree.ml	62	2.9	$r.Cons.1 \prec l.Node.1$

■ **Figure 13** Analysis of ADT-manipulating OCaml functions

Function	Delexp (ms)	Overhead	Summary
copy.list.py	69	1.3	$l2 \prec l1$
copy.list.comprehension.py	69	1.4	$l2 \prec l1$
copy.set.py	56	1.03	$s2 \prec s1$
copy.set.comprehension.py	59	1.03	$s2 \prec s1$
filter.list.py	92	1.5	$l2 \prec l1$
filter.list.comprehension.py	91	1.5	$l2 \prec l1$
filter.set.py	57	1.04	$s2 \prec s1$
filter.set.comprehension.py	65	1.08	$s2 \prec s1$
mult2.list.py (Fig. 1)	70	1.4	$l2 \prec 2 \times l1$
mult2.list.comprehension.py	50	2	$l2 \prec 2 \times l1$
mult2.set.py	57	1.1	$s2 \prec 2 \times s1$
mult2.set.comprehension.py	62	1.06	$s2 \prec 2 \times s1$
partition.list.py	124	1.7	$l2 \prec l1 \wedge l3 \prec l1$
partition.set.py	76	1.1	$s2 \prec s1 \wedge s3 \prec s1$
partition.set-list.py	104	1.4	$l2 \prec s1 \wedge l3 \prec s1$
partition.list-set.py	96	1.2	$s2 \prec l1 \wedge s3 \prec l1$
tail.list.py (Ex. 3.2)	76	1.3	$l2 \prec s1$

■ **Figure 14** Analysis of list-manipulating Python programs

683 The first four programs are inspired from the OCaml standard library over lists. For
684 `filter_le.ml`, the OCaml analysis without DelExp could infer that `l.Cons.1 < inf`, but
685 for every other functions, the generated summary without DelExp is \top , i.e. we have no
686 information about the function’s behaviour. The DelExp domains computed the intended
687 summaries for the ADT-manipulating functions. In particular, notice that it also inferred
688 precise summaries for the versions of `mult2` and `map` that operate on the labels of a tree
689 rather than the elements of a list. Note that given the `filter` function, which, given a
690 functional predicate as input, filters the list with the elements respecting this predicate, we
691 can compute that the output list is included in the input list. We can also infer a specification
692 for `listtotree.ml`, that translates a list to a tree: we can infer that the elements of the
693 output tree are included in the elements of the input list. It highlights that our domain is
694 agnostic in the underlying structure and type of the objects abstracted by the weak variables.
695 We also computed the overhead of our DelExp domain by comparing the execution time
696 with and without our domain. Our overhead is between 2 and 4. Part of this cost is due to
697 the computation of our transitive closure. As an optimization, our transitive closure could
698 be computed more efficiently using labeled union-find data structure from Lesbre et al. [22].

699 **Python.** The Python static analysis we build upon [27] is not compositional, as this is
700 made very difficult by Python’s dynamic typing and flexible semantics. However, our Python
701 programs define lists and sets (l_1 and s_1) initialized with a random size and content and give
702 them as input to functions, which compute new lists or sets (l_2, l_3 or s_2, s_3) by applying a
703 certain transformation. Operations on lists and sets can be defined with `for` loops or with
704 comprehensions. Figure 14 shows our results. `partition.x.py` splits a container into two
705 containers with respect to an integer pivot. This is the intermediate operation of quicksort.
706 Note that `partition.set-list.py` builds two lists from a set and `partition.list-set.py`
707 two sets from a list. The overhead of the DelExp domains in Python is between 1.1 and 1.7.
708 Since the `reduce` function is called at every function application in our implementation, the
709 OCaml compositional analysis of recursive functions performs more reduction than a similar
710 program in Python. This may explain the overhead difference between the two languages.

711 **Comparison with QUIC graphs.** Compared to QUIC graphs, we have a greater expressivity
712 on map-like functions, i.e. functions that apply element-wise transformations. `copy.set.py`,
713 `partition.set.py` and `filter.set.py` from Fig. 14 are examples from QUIC’s Table 1 [11],
714 chosen for their interest and translated to Python from their toy language. Note that we
715 cannot experimentally compare both approaches: no artifact was archived at the time, and
716 we could not compile the available source code.¹ However, Cox et al. [11] reported an
717 overhead of a factor of 9 on similar standard set-manipulating functions. Note that they
718 operate on a toy set language, whereas DelExp was run on top of existing analyses, for real
719 languages of different paradigms (statically typed and functional for OCaml, dynamically
720 typed and object-oriented for Python).

721 An artifact reproduces those experimental claims².

¹ <https://plv.colorado.edu/projects/quicgraphs/>

² <https://zenodo.org/records/18498584>

722 **9 Related Work**

723 **Weak variables.** Programs often manipulate an unbounded number of objects (e.g. heap-
 724 allocated data or containers) whereas most abstract domains manipulate fixed and finite
 725 numbers of objects. Gopan et al. [13] overcame this issue by grouping a possibly unbounded
 726 collection of concrete objects into one abstract *summarization* variable. They defined
 727 operators to soundly manipulate them. This generic framework was used to abstract
 728 different kind of containers. For numerical trees, Journault et al. [19] use tree automata
 729 (a generalization of regular expressions) to represent a set of paths from root to node – or
 730 equivalently, a set of nodes. To each abstract path, they associate a weak variable collecting
 731 the integer label of those nodes. Monat [26] used them to summarize all possible characters,
 732 represented as their ASCII code, whereas Valnet et al. [31] used them to summarize ADTs’
 733 fields.

734 **Arrays.** The abstraction of fixed-size containers such as arrays also leverages those summariz-
 735 ation variables [3]. Gopan et al. [14] partitioned arrays with respect to the value of an integer
 736 variable i (provided by the user or by a pre-analysis) into three index-bounded segments,
 737 $a_{<i}$ (resp. $a_{>i}$) for array slots of index lower (resp. greater) than i and a_i for slot i . For
 738 both segments $a_{<i}$ and $a_{>i}$, a weak variable abstracts the indices of the segment and another
 739 abstracts its content. With a similar approach, Halbwachs and Péron [15] define for each
 740 segment *slice variables*, ranging over arrays indices. They allow expressing relations between
 741 segments from two different arrays ($a = b$ on a set of indices p means that $\forall i \in p, a[i] = b[i]$).
 742 *Segmentation* is introduced by Cousot et al. [9] as a functor both deciding how to partition
 743 arrays and how to summarize their contents. Partitioning is then semantically-guided and
 744 can dynamically change while analyzing the program. However, it cannot express content
 745 relations between summarized segments. Those methods were reviewed by Bautista et al [1]
 746 and extended [2] so that segment summaries can contain non-recursive algebraic data types
 747 and refer to arbitrary program variables. Non-contiguous partitions are supported by Liu
 748 and Rival [23]. Segmentation allows very precise abstraction. However, since this approach
 749 is index-based, it is not suited for general containers (sets, maps, linked lists or user-defined
 750 algebraic data types). As segmentation defines when to create weak variables whereas DelExp
 751 increases precision in the presence of those variables, they could be used in combination.
 752 Braine et al. [6], Monniaux and Gonnord [28] handle array analysis by translation into Horn
 753 Clauses, then delegated to solvers: they show the approach can prove array sortedness.

754 **Compositional analysis.** Compositional analysis is a long-known technique to improve
 755 analysis scalability [8]. First applied to `while` programs [21], it has been used in various
 756 settings, to analyze independently decomposed parts of a program [12] or for inter-procedural
 757 analysis [20]. Bautista et al. [1, 2] define an ADT domain in an input-output analysis for
 758 non-recursive imperative programs. The precision of function summaries can be improved
 759 with partitioning [4, 5, 31].

760 **Relational analysis on containers.** A key ingredient to achieve interesting precision levels
 761 with compositionality is to rely on relational abstract domains. However, retaining relational
 762 information between weak variables is a non-trivial concern [30]. Cox et al. [11] defined the
 763 domain QUIC graphs, presented for sets of scalars. Our domain differs from QUIC graphs in
 764 three main aspects: whereas QUIC graphs can express union and intersection of sets, our
 765 domain can instead express transformations (Sec. 4) on weak variables. Besides, our domain

766 supports polymorphic variables (Sec. 7) and heterogeneous environments (Sec. 6). Another
767 relevant line of work is shape analysis, which aims at abstracting memory regions (e.g. heap
768 allocated data structures), partitioning them according to high-level predicates. They can
769 express relations with sophisticated folding and unfolding operators [7]. For example, Illous
770 et al. [17] can express in-place modification of data structures ("transform-into" relation).
771 Techniques based on separation logic sometimes use fold and expand operators to abstract
772 memory regions: they may leverage the DelExp domains to improve their precision in this
773 context.

774 10 Conclusion

775 This article presents a new family of relational domains for containers and shows how to
776 derive precise analyses for container-manipulating programs. Built upon the concept of weak
777 variables [3, 13], those domains can express relations between them. Compared to the state
778 of the art, our domains yield a higher expressivity for containers-transforming programs, e.g.
779 it can precisely abstract map-like functions.

780 Those domains are highly parametric. First, they are language-agnostic: they are
781 formalized on a very generic language manipulating weak variables. They are also modular
782 in the domain chosen to represent weak variables: those variables can abstract a set of
783 objects of any type, even a polymorphic one. Since they support heterogeneous environments,
784 they can manipulate *optional* weak variables: therefore, they can express properties on the
785 content of possibly empty containers. Finally, new domains expressing element-wise content
786 transformation can be derived at low cost from the DelExp domain – e.g. the ApplyDelExp
787 domain can express higher-order transformations. As a consequence, they can be used in
788 product with any existing analysis manipulating weak variable, for applications as different
789 as Python lists and user-defined OCaml Algebraic Data Types.

790 Since they are *relational*, those domains make it possible to define precise compositional
791 analyses for container-manipulating functions. In particular, we used it to extend the OCaml
792 compositional analysis from Valnet et al. [31] on algebraic data types, so far less expressive and
793 restricted to monomorphism. Those domains were implemented into the MOPSA platform
794 in product with the existing OCaml analysis [31]. Experimental results show that they enable
795 to automatically discover precise invariants on ADT-manipulating, possibly higher-order,
796 functions from the OCaml standard library. The same domains were used in product with
797 an existing Python analysis [26], proving the genericity of the approach. The overhead of
798 the DelExp domain was around 1.5 for Python and 3 for OCaml, lower than the overhead of
799 the QUIC graphs approach [11] (around 9).

800 For future work, we consider using the DelExp domain in combination with segmentation
801 techniques or shape analysis abstractions to evaluate its use in different contexts. We also
802 aim at defining more sophisticated domains on containers able to express more fine-grained
803 properties. As a foreseen next step, we want to design a relational domain able to express
804 properties on balanced binary search trees to automatically prove correctness of functional
805 algorithms manipulating them.

806 Acknowledgements

807 This work was supported by project ANR-22-PECY0005 "Secureval" managed by the French
808 National Research Agency for France 2030. This work was supported by France's Agence
809 Nationale de la Recherche (ANR), program France 2030, reference ANR-22-PTCC-0001.

810 **References**

- 811 1 Bautista, S.: Static Analysis of Algebraic Data Types and Arrays. Ph.D. thesis, ENS
812 Rennes (2023)
- 813 2 Bautista, S., Jensen, T., Montagu, B.: An input–output relational domain for algebraic
814 data types and functional arrays. *Formal Methods in System Design* pp. 1–74 (2024)
- 815 3 Blanchet, B., Cousot, P., Cousot, R., Feret, J., Mauborgne, L., Miné, A., Monniaux, D.,
816 Rival, X.: Design and implementation of a special-purpose static program analyzer for
817 safety-critical real-time embedded software. In: *The essence of computation: complexity,*
818 *analysis, transformation*, pp. 85–108, Springer (2002)
- 819 4 Bourdoncle, F.: Abstract interpretation by dynamic partitioning. *Journal of Functional*
820 *Programming* **2**(4), 407–435 (1992)
- 821 5 Boutonnet, R., Halbwachs, N.: Disjunctive relational abstract interpretation for
822 interprocedural program analysis. In: *Verification, Model Checking, and Abstract*
823 *Interpretation: 20th International Conference, VMCAI 2019, Cascais, Portugal, January*
824 *13–15, 2019, Proceedings 20*, pp. 136–159, Springer (2019)
- 825 6 Braine, J., Gonnord, L., Monniaux, D.: Data Abstraction: A General Framework to
826 Handle Program Verification of Data Structures. In: *Lecture notes in computer science,*
827 *Lecture notes in computer science*, vol. 12913, pp. 215–235, Chicago, United States
828 (Oct 2021), doi:10.1007/978-3-030-88806-0_11, URL [https://inria.hal.science/
829 hal-03321868](https://inria.hal.science/hal-03321868)
- 830 7 Chang, B.Y.E., Rival, X.: Relational inductive shape analysis. *ACM SIGPLAN Notices*
831 **43**(1), 247–260 (2008)
- 832 8 Cousot, P., Cousot, R.: Modular static program analysis. In: *International Conference*
833 *on Compiler Construction*, pp. 159–179, Springer (2002)
- 834 9 Cousot, P., Cousot, R., Logozzo, F.: A parametric segmentation functor for fully
835 automatic and scalable array content analysis. *ACM SIGPLAN Notices* **46**(1), 105–
836 118 (2011)
- 837 10 Cousot, P., Halbwachs, N.: Automatic discovery of linear restraints among variables of a
838 program. In: *Proceedings of the 5th ACM SIGACT-SIGPLAN symposium on Principles*
839 *of programming languages*, pp. 84–96 (1978)
- 840 11 Cox, A., Chang, B.Y.E., Sankaranarayanan, S.: Quic graphs: Relational invariant
841 generation for containers. In: *European Conference on Object-Oriented Programming*,
842 pp. 401–425, Springer (2013)
- 843 12 Farzan, A., Kincaid, Z.: Compositional recurrence analysis. In: *2015 Formal Methods in*
844 *Computer-Aided Design (FMCAD)*, pp. 57–64, IEEE (2015)
- 845 13 Gopan, D., DiMaio, F., Dor, N., Reps, T., Sagiv, M.: Numeric domains with summarized
846 dimensions. In: *Tools and Algorithms for the Construction and Analysis of Systems: 10th*
847 *International Conference, TACAS 2004, Held as Part of the Joint European Conferences*
848 *on Theory and Practice of Software, ETAPS 2004, Barcelona, Spain, March 29–April 2,*
849 *2004. Proceedings 10*, pp. 512–529, Springer (2004)
- 850 14 Gopan, D., Reps, T., Sagiv, M.: A framework for numeric analysis of array operations.
851 In: *Proceedings of the 32nd ACM SIGPLAN-SIGACT symposium on Principles of*
852 *programming languages*, pp. 338–350 (2005)
- 853 15 Halbwachs, N., Péron, M.: Discovering properties about arrays in simple programs. In:
854 *Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design*
855 *and Implementation*, p. 339–348, PLDI '08, ACM (2008)
- 856 16 Hindley, R.: The principal type-scheme of an object in combinatory logic. *Transactions*
857 *of the american mathematical society* **146**, 29–60 (1969)

- 858 **17** Illous, H., Lemerre, M., Rival, X.: A relational shape abstract domain. Formal methods
859 in system design **57**(3), 343–400 (2021)
- 860 **18** Journault, M., Miné, A., Monat, R., Ouadjaout, A.: Combinations of reusable abstract
861 domains for a multilingual static analyzer. In: Verified Software. Theories, Tools, and
862 Experiments: 11th International Conference, VSTTE 2019, New York City, NY, USA,
863 July 13–14, 2019, Revised Selected Papers 11, pp. 1–18, Springer (2020)
- 864 **19** Journault, M., Miné, A., Ouadjaout, A.: An abstract domain for trees with numeric
865 relations. In: European Symposium on Programming, pp. 724–751, Springer (2019)
- 866 **20** Kincaid, Z., Breck, J., Boroujeni, A.F., Reps, T.: Compositional recurrence analysis
867 revisited. ACM SIGPLAN Notices **52**(6), 248–262 (2017)
- 868 **21** Kozen, D.: Kleene algebra with tests. ACM Transactions on Programming Languages
869 and Systems (TOPLAS) **19**(3), 427–443 (1997)
- 870 **22** Lesbre, D., Lemerre, M., Ait-El-Hara, H.R., Bobot, F.: Relational abstractions based
871 on labeled union-find. Proceedings of the ACM on Programming Languages **9**(PLDI),
872 1194–1219 (2025)
- 873 **23** Liu, J., Rival, X.: Abstraction of arrays based on non contiguous partitions. In:
874 International Workshop on Verification, Model Checking, and Abstract Interpretation,
875 pp. 282–299, Springer (2015)
- 876 **24** Milner, R.: A theory of type polymorphism in programming. Journal of computer and
877 system sciences **17**(3), 348–375 (1978)
- 878 **25** Miné, A.: Mopsa (2025), URL <https://mopsa.lip6.fr>
- 879 **26** Monat, R.: Static type and value analysis by abstract interpretation of Python programs
880 with native C libraries. Ph.D. thesis, Sorbonne Université (2021)
- 881 **27** Monat, R., Ouadjaout, A., Miné, A.: Static type analysis by abstract interpretation
882 of Python programs. In: 34th European Conference on Object-Oriented Programming
883 (ECOOP 2020), pp. 17–1, Schloss Dagstuhl–Leibniz-Zentrum für Informatik (2020)
- 884 **28** Monniaux, D., Gonnord, L.: Cell morphing: From array programs to array-free horn
885 clauses. In: International Static Analysis Symposium, pp. 361–382, Springer (2016)
- 886 **29** Ouadjaout, A., Monat, R., Miné, A., Journault, M., Parolini, F., Milanese, M., Boillot,
887 J.: Mopsa (2025), URL <https://gitlab.com/mopsa/mopsa-analyzer>
- 888 **30** Siegel, H., Simon, A.: Summarized dimensions revisited. Electronic Notes in Theoretical
889 Computer Science **288**, 75–86 (2012)
- 890 **31** Valnet, M., Monat, R., Miné, A.: Compositional static value analysis for higher-order
891 numerical programs. In: 39th European Conference on Object-Oriented Programming
892 (ECOOP 2025), vol. 333, p. 15, Dagstuhl Publishing (2025)