

Shape Analysis in a Functional Language by Using Regular Languages *

Manuel Montenegro
Dpto. de Sistemas
Informáticos y Computación
Univ. Complutense de Madrid
C/ Prof. José García
Santesmases s/n
28040, Madrid, Spain
montenegro@fdi.ucm.es

Ricardo Peña
Dpto. de Sistemas
Informáticos y Computación
Univ. Complutense de Madrid
C/ Prof. José García
Santesmases s/n
28040, Madrid, Spain
ricardo@sip.ucm.es

Clara Segura
Dpto. de Sistemas
Informáticos y Computación
Univ. Complutense de Madrid
C/ Prof. José García
Santesmases s/n
28040, Madrid, Spain
csegura@sip.ucm.es

ABSTRACT

Shape analysis is concerned with the compile-time determination of the ‘shape’ the heap may take at runtime, meaning by this the pointer chains that may happen within, and between, the data structures built by the program. This includes detecting alias and sharing between the program variables.

Functional languages facilitate somehow this task due to the absence of variable updating. Even though, sharing and aliasing are still possible. We present an abstract interpretation-based analysis computing precise information about these relations. In fact, the analysis gives an information more precise than just the existence of sharing. It informs about the paths through which this sharing takes place. This information is critical in order to get a modular analysis and not to lose precision when calling an already analysed function.

The main innovation with respect to the literature is the use of regular languages to specify the possible pointer paths from a variable to its descendants. This additional information makes the analysis much more precise while still being affordable in terms of efficiency. We have implemented it and give convincing examples of its precision.

Keywords: functional languages, abstract interpretation, shape analysis, points-to analysis, regular languages.

1. MOTIVATION

Shape analysis is concerned with statically determining the connections between program variables through pointers in the heap that may occur at runtime. As particular cases it includes sharing and alias between variables. To know the shape of the heap for every possible program execution

is undecidable in general, but the analysis computes an over-approximation of this shape. This means that it may include relations that will never happen at runtime.

Much work has been done in imperative languages (see Sec. 7), specially for C. There, the sharing detection is aggravated by the fact that variables are mutable, and they may point to different places at different times. We have addressed the problem for a first order functional language. This simplifies some of the difficulties since variables do not mutate. A consequence is that the inferred relations are immutable considering different parts of the program text. Another consequence is that the heap is never updated. It can only be increased with new data structures, or decreased by the garbage collector. But the latter cannot produce effects in its live part.

Our analysis puts the emphasis on three properties: (1) modularity; (2) precision; and (3) efficiency. For the sake of scalability, it is important for the analysis to be modular. The results obtained for a function should summarize the shape information so that the user functions should be able to compute all the sharing produced when calling it. Looked at from outside, and given that the language is functional, a function may only create sharing between its result and its arguments, or between the results themselves, but it can never create new sharing between the arguments. The internal variables become dead after the call, so the result of analysing a function only contains its input-output sharing behaviour. Differently from previous works, we compute the *paths* through which this sharing may occur in a precise way. This information is used to propagate to the caller the sharing created by a call. In this way, large programs can be analysed with a cost linear in the number of functions.

The motivation for our analysis is a type system we have developed for a functional language with explicit memory disposal [10]. This feature may create dangling pointers at runtime. The language also provides automatically allocated and deallocated heap regions, instead of having a runtime garbage collector. This feature can never create dangling pointers, so it plays no role in the current work and we will not mention it anymore. We have proved that passing successfully the type inference phase gives total guarantee that there will not be such dangling pointers. For typechecking a function, it is critical to know at compile time which variables may point to the disposed data structures, and for this a precise sharing analysis was needed. Nevertheless,

*Work partially funded by the projects TIN2008-06622-C03-01/TIN (STAMP), and S2009/TIC-1465 (PROMETIDOS).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.

PPDP '13, September 16 - 18 2013, Madrid, Spain

Copyright 2013 ACM 978-1-4503-2154-9/13/09 ... \$15.00.

<http://dx.doi.org/10.1145/2505879.2505893>.

```

unshuffle [] = ([],[])
unshuffle (x:xs) = (x:ys2, ys1)
                    where (ys1,ys2) = unshuffle xs
merge [] ys = ys
merge (x:xs) [] = x:xs
merge (x:xs) (y:ys) | x <= y = x : merge xs (y:ys)
                    | otherwise = y : merge (x:xs) ys
msort [] = []
msort [x] = [x]
msort xs = merge (msort xs1) (msort xs2)
            where (xs1, xs2) = unshuffle xs

```

Figure 1: *mergesort* algorithm in *Full-Safe*

<i>prog</i>	→	$\overline{data_i}; \overline{dec_j}; e$	{ <i>Core-Safe</i> program }
<i>dec</i>	→	$f \overline{x_i} = e$	{ recursive, polymorphic }
<i>e</i>	→	c	{ literal of a basic type }
		x	{ variable }
		$f \overline{a_i}$	{ function application }
		$C \overline{a_i}$	{ constructor application }
		let $x_1 = e_1$ in e_2	{ non-rec., monomorphic }
		case x of $\overline{alt_i}$	{ case expression }
<i>alt</i>	→	$C \overline{x_i} \rightarrow e$	{ case alternative }

Figure 2: Simplified *Core-Safe* syntax

we believe that the sharing analysis presented here could be equally useful for other purposes, since it provides precise information about the heap shape. Note that some shapes, such as cyclic or doubly chained lists, cannot be created by a functional language, so they are out of the scope of our analysis. But, in some cases, the analysis is capable of asserting that a given structure is a tree, i.e. it does not have internal sharing.

Our prior prototype shape analysis done in [12] was correct but imprecise, specially in function applications and **case** expressions. The reason for this is that it does not suffice knowing that two variables share a common descendant. We should more precisely know through which paths this sharing occurs.

The main contribution of this paper with respect to [12] is the incorporation of regular languages to our abstract domain. Each word of the language defines a pointer path within a data structure. Having regular languages introduces additional problems such as how to combine them during the analysis, how to compare them, and specially how to guarantee that a fixpoint will be reached after a finite number of iterations. We show that we have increased the precision of our prior analysis, and that the new problems can be tackled with a reasonable efficiency.

The plan of the paper is as follows: Sec. 2 provides a mild introduction to the analysis via a small example. Then, Sections 3, 4 and 5 contain all the technical material about the abstract domain, abstract interpretation rules, correctness, widening, decidability, and cost of the operations done on regular expressions. Sec. 6 presents our implementation and gives more examples. Finally, Sec. 7 concludes and discuss some related work.

2. SHAPE ANALYSIS BY EXAMPLE

Our reference language *Safe* is a first-order eager language with a syntax similar to Haskell's. Fig. 1 shows a mergesort algorithm written in *Full-Safe*. The compiler's front-end processes *Full-Safe* and produces a bare-bones functional

language called *Core-Safe*. This transformation desugars pattern matching into **case** expressions, transforms **where** clauses into **let** expressions, collapses several function-defining equations into a single one, and ensures unique names for the variables. In Fig. 2 we show a simplified *Core-Safe*'s syntax. A program *prog* is a sequence of possibly recursive polymorphic data and function definitions followed by a main expression e whose value is the program result. The abbreviation $\overline{x_i}$ stands for $x_1 \cdots x_n$, for some n . In Fig. 3 we show the translation to *Core-Safe* of the *msort* function of Fig. 1.

Our shape analysis infers the following sharing information for the functions *unshuffle* and *merge*:

$$\begin{aligned} \Sigma(\text{unshuffle}) &= \{ res \xrightarrow{12^*1+22^*1} \bullet \xleftarrow{2^*1} xs \} \\ \Sigma(\text{merge}) &= \{ res \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} xs, res \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} ys, \\ &\quad res \xrightarrow{2^*} \bullet \xleftarrow{2^*} xs, res \xrightarrow{2^*} \bullet \xleftarrow{2^*} ys \} \end{aligned}$$

The meaning for *unshuffle* is the following: the resulting tuple *res* of calling the function with an input list *xs*, may share the elements of this list. Moreover, the path reaching a common descendant, in the case of *res*, begins either with a 1 or a 2 (this should be understood as descending to the left or to the right element of the tuple), and then follows by the path 2^*1 , by this meaning that we should take the tail of the (left or right) list a number of times, and then take the head. From *xs*'s point of view of, the common descendant can be reached by a similar path 2^*1 .

The meaning for *merge* is quite precise: the resulting list *res* may share its elements with any of the input lists *xs* and *ys*, and additionally one or more tails of *res* may be shared with one or more tails of both *xs* and *ys*. This is what the path 2^* means.

When analysing *msort*'s code of Fig. 3, we have the information about *unshuffle* and *merge* available. By substituting the actual arguments for the formal ones, we get the following relations:

$$\begin{aligned} R_1 &= \{ p \xrightarrow{12^*1+22^*1} \bullet \xleftarrow{2^*1} xs \} \\ R_2 &= \{ res \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} z_1, res \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} z_2, \\ &\quad res \xrightarrow{2^*} \bullet \xleftarrow{2^*} z_1, res \xrightarrow{2^*} \bullet \xleftarrow{2^*} z_2 \} \end{aligned}$$

The **case** and **let** expressions in *msort* introduce more relations:

$$R_3 = \{ x \xrightarrow{\epsilon} \bullet \xleftarrow{1} xs, y_1 \xrightarrow{\epsilon} \bullet \xleftarrow{1} p, y_2 \xrightarrow{\epsilon} \bullet \xleftarrow{2} p \}$$

From these relations, we can derive other by reflexivity, symmetry and transitivity, such as:

$$R_4 = \{ y_1 \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} xs, y_2 \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} xs, y_1 \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} y_2 \}$$

In the first iteration of *msort*'s analysis, the only relation that can be inferred between its result and its argument is $xs \xrightarrow{1} \bullet \xleftarrow{1} res$. This is due to the third line. The rest of the code gives us sharing information between the internal variables, but this cannot be propagated to the arguments, because in the internal recursive calls to *msort* we have nothing to start with. But, by interpreting the internal calls with the sharing information $\Sigma_1(\text{msort}) = \{ xs \xrightarrow{1} \bullet \xleftarrow{1} res \}$, we get the following bigger result: $\Sigma_2(\text{msort}) = \{ xs \xrightarrow{2^*1+1} \bullet \xleftarrow{2^*1+1} res \}$. If we interpret the code a third time by using this information when interpreting the internal calls, we get

```

msort xs = case xs of
  []   -> []
  x:xx -> case xx of [] -> x:[]
          _:_ -> let p = unshuffle xs           in
                  let y1 = case p of (s1,s2) -> s1 in
                  let y2 = case p of (w1,w2) -> w2 in
                  let z1 = msort y1           in
                  let z2 = msort y2           in
                  merge z1 z2

```

Figure 3: Function *msort* in *Core-Safe*

again the same result. So, a fixed point has been reached, and we consider this result as a correct approximation of the sharing created by `msort`.

It is worthwhile to remark that our prior analysis [12] of the same program gave us the additional spurious sharing information $\{z_2 \rightarrow \bullet \leftarrow z_1\}$, meaning that a descendant of z_2 is shared by z_1 (the regular languages were absent in that analysis). Having spurious relations is not incorrect, but just imprecise. Since we used this analysis to type a destructive version of `msort`, using in turn a destructive version of `merge`, our type system rejected the function because of this additional sharing. The cause of this imprecision was a worse analysis of `case` expressions and function applications due to the absence of the paths represented by the regular languages.

3. THE ANALYSIS

We formally define here the analysis approximating the runtime sharing relations between the program variables. At this point, types have already been inferred, so the analysis can ask for type-related issues, such as the positions of constructor descendants, their types, and the like.

3.1 Sharing relation

In order to capture sharing, we define a binary relation between variables:

DEFINITION 1. *Given two variables x and y , in scope in an expression, a sharing relation is a set of two pairs $\{(x, p_1), (y, p_2)\}$ specifying that x and y share a common descendant. Moreover, the regular languages denoted by p_1 and p_2 respectively define the possible pointer chains through which x and y reach their common descendant. We shall denote this sharing relation either by $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ or $y \xrightarrow{p_2} \bullet \xleftarrow{p_1} x$.*

For the sake of readability, we shall assume in the following p_1 and p_2 to be regular expressions that denote regular languages, but the actual implementation does not use them, though. Notice that, if $p_1 = \epsilon$, then x is a descendant of y , and symmetrically for p_2 .

The regular languages have pairs i_C as alphabet symbols, where i is a natural number starting at 1, and C is a data constructor. The symbol i_C denotes a singleton pointer path in the heap passing through the i -th argument of constructor C . For instance, $x \xrightarrow{2^*} \bullet \xleftarrow{1_{(,)}} y$ indicates that a tail of the list x is pointed-to by the first element of the tuple y . In the examples, we shall usually omit the constructor.

The relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2}$ is symmetric by definition and reflexive by writing $p_1 = p_2 = \epsilon$. But the transitivity does not

hold, i.e. $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ and $y \xrightarrow{p_3} \bullet \xleftarrow{p_4} z$, with $p_2 \neq \epsilon$, does not necessarily imply $x \xrightarrow{p_1} \bullet \xleftarrow{p_4} z$. However, the transitivity holds in some cases, for example when y reaches its common descendant with x through *the same path* as it reaches its common descendant with z , as shown in Figure 4a.

More generally, we can investigate the languages denoted by p_2 and p_3 , and decide whether a path in p_2 coincides with, or is a prefix of, a path in p_3 (as shown in Figure 4b), or the other way around. In these cases, there may exist a sharing path through y between x and z . Notice that both p_2 and p_3 are upper approximations to the actual runtime paths, so the risk of imprecision is still there, but if there are no such paths we are certain that there will not be paths at runtime either, and we can safely omit a tuple relating x and z from the sharing relation. The rules computing the sharing derived by transitivity are explained in detail in Section 3.4.

3.2 The abstract interpretation

Based on the above considerations, we define an abstract interpretation S (meaning *sharing*) which, given an expression e and a set R containing an upper approximation to the sharing relations between the variables in scope in e , delivers another set R_{res} (*res* stands for result) containing (an upper approximation to) all the relations between the result of evaluating e , named *res*, and its variables in scope. To be precise, R and R_{res} must record at least the *minimum* information needed in order to compute all possible sharing, i.e. if we have $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ in R or R_{res} , and p_3 denotes all possible paths inside the data structure pointed-to by x and y , then we understand that $x \xrightarrow{p_1 \cdot p_3} \bullet \xleftarrow{p_2 \cdot p_3} y$ is implicitly included in the relation.

Notice that this means that:

- If two variables x and y share a substructure in the heap as in Figure 5a, there must exist a sharing relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ containing at least the paths w_1 and w_2 , leading to the first point of confluence. Their extensions with a common path w need not.
- In case a variable x has internal sharing, as shown in Figure 5b, there must exist a sharing relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} x$ containing at least the paths w_1 and w_2 leading to the first point of confluence.

In order to achieve a modular analysis, it is very important to reflect the result of the analysis of a function f in a *function signature environment*, so that when the analysis finds calls to f in the body of another function g , it uses

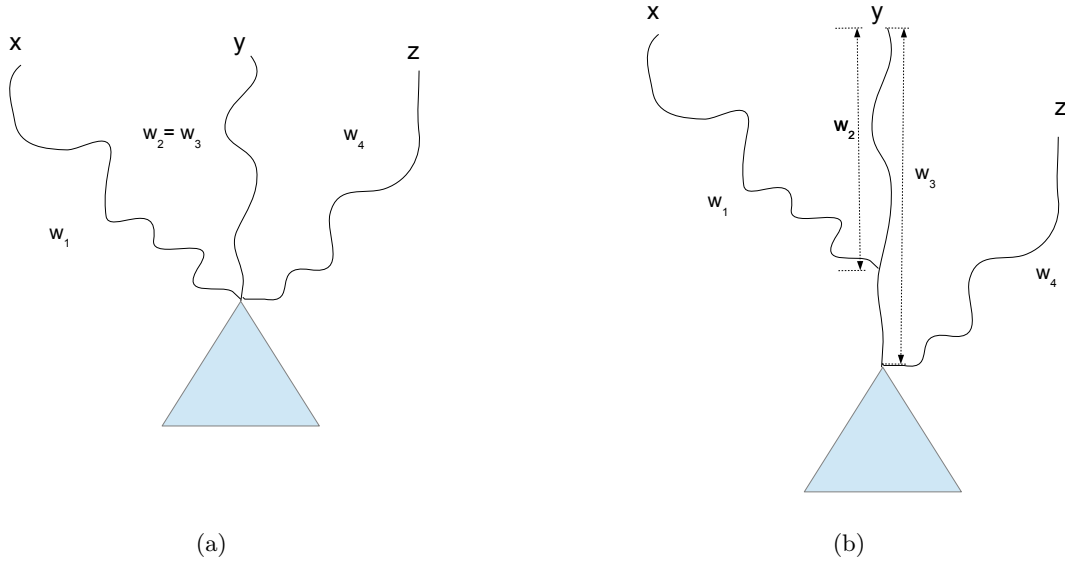


Figure 4: Particular cases of transitivity

this knowledge to compute the sharing relations for g . We keep function signatures in a global environment Σ , so that $\Sigma(f)$ is a set R_{res} containing the sharing relations between the result of calling f and its arguments. The interpretation $S[e] R \Sigma$ gives us the relations between (the normal form of) e and its variables in scope, provided Σ gives us correct approximations to the sharing relations of the functions called from e .

The rules for expressions are explained in detail in Section 3.3. The interpretation S_d of a function definition $f x_1 \dots x_n = e_f$ begins with the interpretation of its body. It is straightforward to extract the signature of the function, which just describes the relations between the result of e_f and its formal arguments \bar{x}_f^n , which are the only variables in scope. In case f is recursive, the interpretation is run several times, by starting with an empty signature for f and then computing the least fixpoint. Each iteration updates f 's signature in the signature environment:

$$S_d[f x_1 \dots x_n = e_f] \Sigma = fix (\lambda \Sigma. \Sigma[f \rightarrow S[e_f] R_0 \Sigma]) \Sigma_0$$

where $\Sigma_0 = \Sigma[f \rightarrow \emptyset]$
 $R_0 = \{x_i \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x_i \mid i \in \{1..n\}\}$

where $\Sigma[f \rightarrow R]$ either adds signature R for f or replaces it in case there was already one for it. Notice that the right hand side of the function definition is analysed starting with a neutral initial relation R_0 in which each argument is only related to itself. This means that the signatures are computed assuming that all the parameters are disjoint and they do not present internal sharing in addition to the trivial sharing relation given by R_0 . When they are not, the function caller knows the additional sharing of the actual arguments and the rule for application merges both information, as we will see in Section 3.3.

As function S is monotonic over a lattice, the least fixpoint exists and could be computed using Kleene's ascending chain if the chain were finite. We come back to this issue in Section 5.

3.3 Interpretation of expressions

The interpretation defined in Figure 6 does a top-down traversal of a function definition, accumulating these relations as soon as bound variables become free variables.

The notation $R[y/x]$ means the substitution of the variable y for the variable x in the relation R . In order to avoid name capture, y must be fresh in R . The operator $R \setminus \{x\}$ removes from R any tuple containing the variable x . The union operator \cup is the usual set union. The closure operation $R_1 \uplus_x^* R_2$ takes a relation R_1 and completes it by adding R_2 and the tuples involving x that can be derived by transitivity. This operation also generates the reflexive relation $x \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x$. We explain this operator in detail in Section 3.4.

An important invariant of the rules presented in Figure 6 is that, in each occurrence of $S[e] R \Sigma$, the set R contains an upper approximation of *all* the sharing relations that at runtime may happen between the variables in scope in e . Also, the set R_{res} returned by $S[e] R \Sigma$ enjoys the same property. It is easy to check that if the property holds for the original call $S[e_f] R_0 \Sigma$, then the rules preserve it.

The rule for a constant c introduces no new sharing. The rule of a variable x specifies that the result is an alias of x , and \uplus_{res}^* propagates to the result the variables to which x is related.

When a constructor application $C \bar{a}_i^m$ is returned as a result, parent-child sharing relations are created with the constructor's children. These are added to the current set R , and then the closure computes all the derived sharing.

When a function application $g \bar{a}_i^m$ is returned as a result, first we get from g 's signature the sharing relations between g 's result and its formal arguments. These are copied by replacing the formal arguments by the actual ones, and then added to the current set. As before, the closure computation does the rest.

The **let** rule is almost self-explanatory: first e_1 is analysed and the sharing computed for e_1 's result is assigned to the new variable in scope x_1 . Using this enriched set R_1 as as-

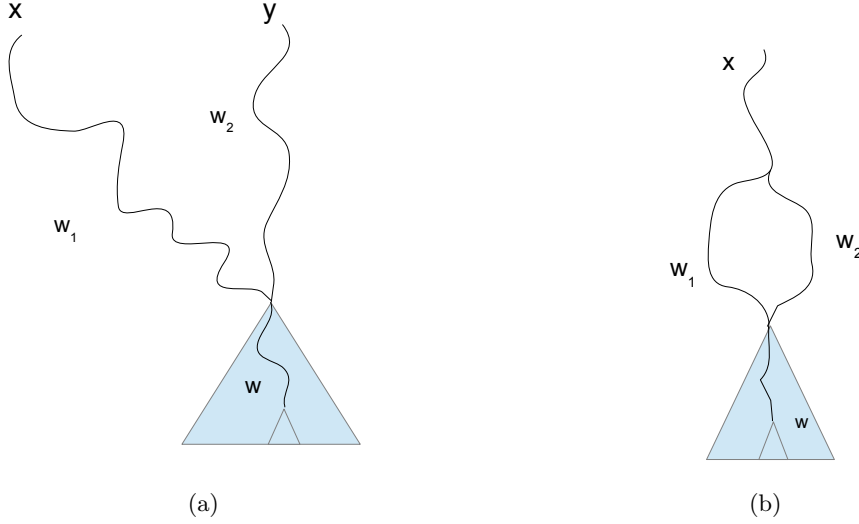


Figure 5: At least paths w_1 and w_2 must be recorded in a sharing relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ (a) or $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} x$ (b)

$$\begin{aligned}
S \llbracket c \rrbracket R \Sigma &= R \\
S \llbracket x \rrbracket R \Sigma &= R \uplus_{res}^* \{res \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x\} \\
S \llbracket C \bar{a}_i^m \rrbracket R \Sigma &= R \uplus_{res}^* \{res \xrightarrow{jC} \bullet \xleftarrow{\epsilon} a_j \mid j \in \{1..m\}, var(a_j)\} \\
S \llbracket g \bar{a}_i^m \rrbracket R \Sigma &= R \uplus_{res}^* \Sigma(g) \llbracket \bar{a}_j / x_j^m \rrbracket \\
S \llbracket \text{let } x_1 = e_1 \text{ in } e_2 \rrbracket R \Sigma &= (S \llbracket e_2 \rrbracket R_1 \Sigma) \setminus \{x_1\} \\
&\quad \text{where } R_1 = (S \llbracket e_1 \rrbracket R \Sigma) \llbracket x_1 / res \rrbracket \\
S \llbracket \text{case } x \text{ of } C_i \bar{x}_{ij}^{n_i} \rightarrow e_i \rrbracket R \Sigma &= \bigcup_i (S \llbracket e_i \rrbracket R_i \Sigma) \setminus \{\bar{x}_{ij}^{n_i}\} \\
&\quad \text{where } R_i = R \uplus_{x_{ij}}^* \{x \xrightarrow{jC_i} \bullet \xleftarrow{\epsilon} x_{ij} \mid j \in \{1..n_i\}\}
\end{aligned}$$

Figure 6: Definition of the abstract interpretation S

sumption, the main expression e_2 is analysed, and its result is the result of the whole **let** expression.

Finally, a **case** expression introduces the pattern variables $\bar{x}_{ij}^{n_i}$ in the scope of a branch e_i . Their sharing relations are derived from the parent x 's ones by first adding the child-parent relation between each x_{ij} and x , and then computing the closure. After analysing the branches, the least upper bound of all the analyses must be computed, expressing the fact that at compile time it is not known which branch will be taken at runtime.

It is important to see whether the relations inferred by the analysis are well-typed. For instance, we could have a relation $x \xrightarrow{p_x} \bullet \xleftarrow{p_y} y$ in which the descendant reached from x and p_1 had a type t , while the descendant reached from y and p_y had a different type t' . This would obviously be a spurious relation since in well-typed programs, an ill-typed sharing may not occur at runtime.

The expression $type(t, p)$ returns the type computed starting at the type t , and then descending through the constructors of the words in p according to its type and to the child chosen at each step. In our language this type can be statically computed. Let t_x be the type computed by the compiler for the variable x .

DEFINITION 2. We say that the relation $x \xrightarrow{p_x} \bullet \xleftarrow{p_y} y$ is well-typed if $type(t_x, p_x) = type(t_y, p_y)$.

LEMMA 3. If the relations in R and Σ are well-typed, then for every expression e , the relations in $S \llbracket e \rrbracket R \Sigma$ are well-typed.

3.4 The closure of a relation

The closure operation \uplus_x^* is defined in terms of the simpler one \uplus_x , which completes a relation set R with a new relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$, where $y \neq x$, by adding the relations that bind x to the variables contained in R , and are derived by transitivity. Both operators are defined in Figure 7.

The inclusion of R and the relations $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$, $x \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x$ is self-explanatory. We shall concentrate on the remaining relations shown in the last lines of the definition.

The second line corresponds to the case illustrated in Figure 4b, while the third one corresponds to the symmetric case. These relations involve the derivative operator $_|_$ whose meaning is:

$$p_1 |_{p_2} = \{w_3 \mid \exists w_2 \in L(p_2). w_2 w_3 \in L(p_1)\}.$$

If p_1 and p_2 denote regular languages so do $p_1 |_{p_2}$, and in

$$\begin{aligned}
R \uplus_x \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y\} &\stackrel{\text{def}}{=} \\
R \cup \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y\} \cup \{x \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x\} \\
\cup \{x \xrightarrow{p_1 \cdot p_3 | p_2} \bullet \xleftarrow{p_4} z \mid y \xrightarrow{p_3} \bullet \xleftarrow{p_4} z \in R\} \\
\cup \{x \xrightarrow{p_1} \bullet \xleftarrow{p_4 \cdot p_2 | p_3} z \mid y \xrightarrow{p_3} \bullet \xleftarrow{p_4} z \in R\} \\
\cup \{x \xrightarrow{p_1 \cdot p_3 | p_2} \bullet \xleftarrow{p_1 \cdot p_4 | p_2} x \mid y \xrightarrow{p_3} \bullet \xleftarrow{p_4} y \in R\} \\
R \uplus_x^* R' &\stackrel{\text{def}}{=} \\
\frac{R \uplus_x \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y\} \mid x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y \in R', y \neq x\}}{\cup \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} x \mid x \xrightarrow{p_1} \bullet \xleftarrow{p_2} x \in R'\}}
\end{aligned}$$

Figure 7: Definition of the closure operation

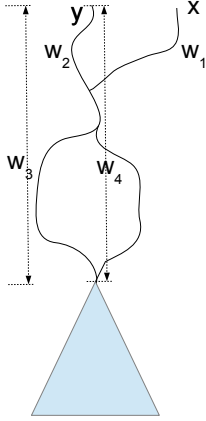


Figure 8: Transitivity with internal sharing.

Section 5 we explain how to compute it¹. In the second line of the definition of \uplus_x the language describing $p_1 \cdot p_3 | p_2$ might be empty. In this case we can discard the corresponding sharing relation from the result of the closure operation. If it is not empty then there exists a word $w_2 \in L(p_2)$ such that it is a prefix of another word $w_3 \in L(p_3)$, so we can start from x , follow a path $w_1 \in L(p_1)$, and then follow the path w_2 without the prefix w_3 (which results in a path of $L(p_3 | p_2)$) in order to reach the common descendant of x and z . The third line of \uplus_x is applicable when a path of p_3 is a prefix of a path of p_2 , and works similarly.

The fourth line deals with the case in which variable x gets internal sharing through variable y , shown in Figure 8. This happens when the path w_2 through which y reaches its common descendant with x is a prefix of both paths w_3 and w_4 representing the internal sharing of y . Then $p_3 | p_2$ and $p_4 | p_2$ are not empty, and contain respectively the paths w_3 and w_4 without the prefix w_2 , which prepended with $w_1 \in L(p_1)$ represent two paths of internal sharing from variable x .

In spite of the restrictions of the \uplus operator, we could replace the \uplus^* operator by a sequence of \uplus operations in

¹When $p_2 = \{a\}$, the language $p_1 | a$ is sometimes called the *derivative* of $L(p_1)$ with respect to a , and it is denoted $a \setminus L$, being $L = L(p_1)$.

all the rules but in function application, because only in the application non-trivial reflexive relations may be added.

In fact, operation $R \uplus_x^* R'$ is used to define the confluence of information happening in a function call. R represents the context of the call, while R' represents the sharing generated by the function between the result and the arguments.

Its definition is divided into two parts:

1. First, we take each relation in R' of the form $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ such that $y \neq x$ and apply the previous transitivity operator incrementally. This is well defined because operator \uplus_x is in a sense *commutative*, as we will prove in Section 4. So the order in which we add the relations of R' is not relevant: the final result may be different but *equivalent*, in the sense that it records the same information.
2. Second, we just add those reflexive relations $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} x \in R'$. In the abstract interpretation, this kind of relations only appear in the application of a function: it may happen that the result of a function f has internal sharing, so a relation $res \xrightarrow{p_1} \bullet \xleftarrow{p_2} res \in \Sigma(f)$. It is not necessary to apply transitivity here because the internal sharing of res either comes from the function itself (i.e. is reflected in R') or through a real argument which already has internal sharing (i.e. is reflected in R). Transitivity, as we will prove in Section 4, would only add redundant information.

4. CORRECTNESS

In this section we provide the main results needed to prove the analysis is well-defined and correct. Full proofs and auxiliary lemmas can be found in [11].

4.1 Properties of the abstract interpretation

First, we prove that operator \uplus_x^* is well defined. As we said in Section 3, the order in which we apply transitivity to the rules belonging to R' may lead to different but equivalent results. First, we give the notion that two sets of relations contain the same information in terms of the sharing paths. Then, we prove well-definedness of the operator.

DEFINITION 4. *A set of sharing relations R is included in R' (written $R \sqsubseteq R'$) if for every sharing relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y \in R$ and every pair of words $w_1 \in L(p_1)$, $w_2 \in L(p_2)$ there exists a sharing relation $x \xrightarrow{p'_1} \bullet \xleftarrow{p'_2} y \in R'$ such that $w_1 \in L(p'_1)$ and $w_2 \in L(p'_2)$. Two sets of relations R and R' are said to be equivalent (written $R \equiv R'$) if both $R \sqsubseteq R'$ and $R' \sqsubseteq R$ hold.*

LEMMA 5 (COMMUTATIVITY OF CLOSURE OPERATION). *Let R be a set of sharing relations and $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$, $x' \xrightarrow{p'_1} \bullet \xleftarrow{p'_2} y'$ a pair of sharing relations such that $y \neq x$ and $y' \neq x'$. Let us define $R_{x,x'}$ and $R_{x',x}$ as follows:*

$$\begin{aligned}
R_{x,x'} &\stackrel{\text{def}}{=} (R \uplus_x \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y\}) \uplus_{x'} \{x' \xrightarrow{p'_1} \bullet \xleftarrow{p'_2} y'\} \\
R_{x',x} &\stackrel{\text{def}}{=} (R \uplus_{x'} \{x' \xrightarrow{p'_1} \bullet \xleftarrow{p'_2} y'\}) \uplus_x \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y\}
\end{aligned}$$

If $y' \neq x$ and $x' \neq y$ then $R_{x,x'} \equiv R_{x',x}$.

Consequently, following different orders in adding the relations of R' lead to equivalent sets of relations.

4.2 Notion of correct approximation

Now we define when a set of relations correctly approximates the real sharing in a heap and the notion of correct signature. The first definition reflects the fact that at least the minimum sharing must be recorded in the relations, i.e. the paths leading to the first point of confluence must be recorded, while their extensions with a common path need not. Notice that this means that in case of internal sharing, each point of internal confluence must also be recorded.

A correct function signature must record enough sharing information to be able to approximate each possible call to that function, i.e. each possible execution of the body. The operational semantics of Core-Safe can be found at [10]. It is a standard big-step operational eager semantics: judgment $E \vdash h, e \Downarrow h', v$ means that expression e in a variable environment E and initial heap h evaluates to value v and the heap changes to h' . If in a heap h there exists an actual sharing between two variables x and y through respective pointer paths w_1 and w_2 , we say that there exists a *sharing condition* in h and denote it by $E(x) \xrightarrow{w_1} \bullet \xleftarrow{w_2} E(y)$ (in h).

DEFINITION 6. A sharing relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ is said to approximate a sharing condition $E(x) \xrightarrow{w_1} \bullet \xleftarrow{w_2} E(y)$ (in h) iff there exists a word w such that $w_1 \in L(p_1 w)$ and $w_2 \in L(p_2 w)$.

DEFINITION 7. Let R be a set of sharing relations, E a runtime environment, and h a heap. We say that R is a correct approximation of E and h , denoted $R \succeq (E, h)$, iff for every pair of variables $x, y \in \text{dom } E$, and pair of words $w_1, w_2 \in \mathcal{V}^*$ if the condition $E(x) \xrightarrow{w_1} \bullet \xleftarrow{w_2} E(y)$ (in h) holds, it is approximated by a sharing relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ in R .

DEFINITION 8 (CORRECT SIGNATURE). A set R of relations is a correct signature for a function definition $f \bar{x}_i^n = e_f$ iff for each execution $E_f \vdash h, e_f \Downarrow h', v$ of the body of the function and every set of relations R' such that $R' \succeq (E_f, h)$ it holds that $R' \uplus_{\text{res}}^* R \succeq (E_f \uplus [\text{res} \mapsto v], h')$. A signature environment Σ is said to be correct iff every signature it contains is correct.

4.3 Correctness

Correctness of the analysis is divided into two steps. First, we prove that given correct signatures of the functions which are called from an expression, the interpretation of the expression is correct. Then, we prove that the interpretation of a function generates a correct signature. For both theorems we need to prove that the transitive closure operator is correct, which we show in the following two lemmas. The second lemma concerns the case in which a variable gets internal sharing through another variable having internal sharing, i.e. the fourth line of operator \uplus_x definition.

LEMMA 9 (TRANSITIVE CLOSURE LEMMA). Let us assume a runtime environment E , a heap h , a set of sharing relations R , some variables x, y, z , (with $y \neq z$) words w_x, w_y, w_z , and paths $p_{xy}, p_{yx}, p_{yz}, p_{zy}$ such that the following holds:

$$\begin{aligned} E(x) &\xrightarrow{w_x} \bullet \xleftarrow{w_y} E(y) \text{ (in } h), \text{ approximated by } x \xrightarrow{p_{xy}} \bullet \xleftarrow{p_{yx}} y \in R \\ E(z) &\xrightarrow{w_z} \bullet \xleftarrow{w_y} E(y) \text{ (in } h), \text{ approximated by } z \xrightarrow{p_{zy}} \bullet \xleftarrow{p_{yz}} y \end{aligned}$$

Then there exists a sharing relation $x \xrightarrow{p_{xz}} \bullet \xleftarrow{p_{zx}} z \in R \uplus_x \{y \xrightarrow{p_{yz}} \bullet \xleftarrow{p_{zy}} z\}$ approximating $E(x) \xrightarrow{w_x} \bullet \xleftarrow{w_z} E(z)$ (in h).

LEMMA 10 (TRANSITIVE SELF-CLOSURE LEMMA). Let us assume a runtime environment E , a heap h , a set of sharing relations R , some variables x, y (with $x \neq y$), words w_x, w_y, w_1, w_2 and paths $p_{x1}, p_{x2}, p_{xy}, p_{yx}$ such that the following holds:

$$\begin{aligned} E(x) &\xrightarrow{w_x w_1} \bullet \xleftarrow{w_x w_2} E(x) \text{ (in } h), \text{ approx. by } x \xrightarrow{p_{x1}} \bullet \xleftarrow{p_{x2}} x \in R \\ E(x) &\xrightarrow{w_x} \bullet \xleftarrow{w_y} E(y) \text{ (in } h), \text{ approximated by } x \xrightarrow{p_{xy}} \bullet \xleftarrow{p_{yx}} y \end{aligned}$$

Then there exists a sharing relation $y \xrightarrow{p_{y1}} \bullet \xleftarrow{p_{y2}} y \in R \uplus_y \{x \xrightarrow{p_{xy}} \bullet \xleftarrow{p_{yx}} y\}$ approximating $E(y) \xrightarrow{w_y w_1} \bullet \xleftarrow{w_y w_2} E(y)$ (in h).

The following theorem establishes the correctness of the abstract interpretation modulo the correctness of function signatures.

THEOREM 11. Assume an expression e , a set of sharing relations R and a correct signature environment Σ . If $S \llbracket e \rrbracket R \Sigma = R'$, then for every execution $E \vdash h, e \Downarrow h', v$ in which $R \succeq (E, h)$, it holds that $R' \succeq (E \uplus [\text{res} \mapsto v], h')$.

Now we prove that the interpretation of a function returns a correct signature. A signature records the sharing between the result and the arguments of the function assuming these are disjoint and without internal sharing. However, a real call to the function may not satisfy such assumption. Given the real configuration (E, h) , we define an hypothetical execution where both the environment \hat{E} and the heap \hat{h} contain the same information as (E, h) but meeting the separation property. The signature of the function captures the sharing information corresponding to this hypothetical execution.

DEFINITION 12. Let (E, h) and (\hat{E}, \hat{h}) be two configurations such that $\text{dom } E = \text{dom } \hat{E}$. A mapping $\gamma : \text{dom } \hat{h} \rightarrow \text{dom } h$ is said to be an entanglement from (\hat{E}, \hat{h}) to (E, h) , iff:

1. For every pointer $\hat{p} \in \text{dom } \hat{h}$, if $\hat{h}(\hat{p}) = C \hat{v}_1 \dots \hat{v}_n$, then $h(\gamma(\hat{p})) = C \gamma(\hat{v}_1) \dots \gamma(\hat{v}_n)$.
2. For every variable $x \in \text{dom } \hat{E}$, $\gamma(\hat{E}(x)) = E(x)$.

As an example, assume a function definition $f x y = C y$. Its signature consists of the following relations: $\{ \text{res} \xrightarrow{c} \bullet \xleftarrow{c} \text{res}, \text{res} \xrightarrow{1c} \bullet \xleftarrow{c} y \}$. Assume we execute a call $f z z$ where $E(z) = p$, $h(p) = C' p' p'$, $h(p') = C'' 3$, i.e. $E_f = [x \mapsto p, y \mapsto p]$. In this case x and y are not disjoint and also contain internal sharing. We can define (\hat{E}_f, \hat{h}) such that $\hat{E}_f(x) = p_1$, $\hat{E}_f(y) = p_2$, $\hat{h}(p_1) = C' p'_1 p'_1$, $\hat{h}(p_2) = C' p'_2 p'_2$ and $\hat{h}(p'_1) = \hat{h}(p'_1) = \hat{h}(p'_2) = \hat{h}(p'_2) = C'' 3$. Then $\gamma(p_1) = \gamma(p_2) = p$, $\gamma(p'_1) = \gamma(p'_1) = \gamma(p'_2) = \gamma(p'_2) = p'$ is an entanglement from (\hat{E}_f, \hat{h}) to (E_f, h) .

The following lemma proves that both the hypothetical and the real execution proceed in parallel and that the information inside the heap is the same although with a different shape. In Figure 9 we show the final heaps of the executions corresponding to the previous example.

LEMMA 13. Assume an execution $E \vdash h, e \Downarrow h', v$ and a configuration (\hat{E}, \hat{h}) . For every entanglement γ from (\hat{E}, \hat{h}) to (E, h) there exist some \hat{h}' , \hat{v}' and γ' such that:

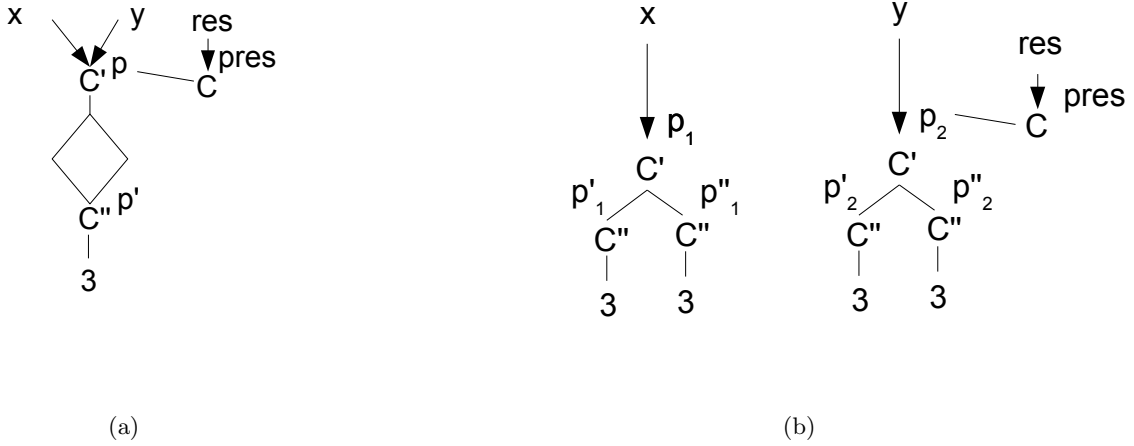


Figure 9: Final heaps in the real execution (a), and the untangled one (b).

1. $\widehat{E} \vdash \widehat{h}, e \Downarrow \widehat{h}', \widehat{v}$.
2. γ' is a conservative extension of γ . That is, $\gamma \subseteq \gamma'$.
3. γ' is an entanglement from $(\widehat{E}, \widehat{h}')$ to (E, h') .
4. $\gamma'(\widehat{v}) = v$.

For the same heap several entanglements may be defined, but we are interested in a configuration $(\widehat{E}, \widehat{h})$, where everything is untangled, as shown in the previous example. This is because, then $R_0 = \{x_i \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x_i \mid i = 1..n\}$ correctly approximates its sharing.

LEMMA 14. *For any configuration (E, h) there exists another configuration $(\widehat{E}, \widehat{h})$ and an entanglement γ from $(\widehat{E}, \widehat{h})$ to (E, h) such that the set $\{x \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x \mid x \in \text{dom } E\}$ is a correct approximation of $(\widehat{E}, \widehat{h})$.*

In the example above, signature is $R' = S[e_f] R_0 \Sigma = \{res \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} res, res \xrightarrow{1C} \bullet \xleftarrow{\epsilon} y\}$. The environment of the call is approximated by $R = \{x \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} y, x \xrightarrow{1C'} \bullet \xleftarrow{2C'} x, y \xrightarrow{1C'} \bullet \xleftarrow{2C'} y\}$. So the final sharing is approximated by $R \uplus_{res}^* R'$ which merges the context of the call with the signature, and contains $\{res \xrightarrow{1C} \bullet \xleftarrow{\epsilon} y, res \xrightarrow{1C} \bullet \xleftarrow{\epsilon} x, res \xrightarrow{1C \cdot 1C'} \bullet \xleftarrow{1C \cdot 2C'} res, res \xrightarrow{1C \cdot 1C'} \bullet \xleftarrow{2C'} x, res \xrightarrow{1C \cdot 1C'} \bullet \xleftarrow{2C'} y, res \xrightarrow{1C \cdot 2C'} \bullet \xleftarrow{1C'} x, res \xrightarrow{1C \cdot 2C'} \bullet \xleftarrow{1C'} y\}$. This happens for each R approximating a context call, so R' is a correct signature for f . We prove this in the following theorem.

THEOREM 15. *Assume a function definition $f \overline{x_i^n} = e_f$, a set of relations $R_0 = \{x_i \xrightarrow{\epsilon} \bullet \xleftarrow{\epsilon} x_i \mid i = 1..n\}$, and an environment Σ with correct signatures. If $R' = S[e_f] R_0 \Sigma$, then R' is a correct signature for f .*

PROOF. (Sketch) Assume a configuration (E, h) with $\text{dom } E = \{\overline{x_i^n}\}$ and a set of relations R such that $R \succeq (E, h)$. If we execute e_f under the configuration (E, h) we get $E \vdash h, e_f \Downarrow h', v$ for some h', v . By Lemma 14 there exists

a mapping γ which entangles a configuration $(\widehat{E}, \widehat{h})$ into (E, h) , where $(\widehat{E}, \widehat{h})$ is correctly approximated by R_0 . Assume we execute e_f under the untangled configuration so as to get $\widehat{E} \vdash \widehat{h}, e_f \Downarrow \widehat{h}', \widehat{v}$ for some \widehat{h}' and \widehat{v} . By correctness theorem (Theorem 11) we know that $R' \succeq (\widehat{E}, \widehat{h}')$, where $\widehat{E}' \stackrel{\text{def}}{=} \widehat{E} \uplus [res \mapsto \widehat{v}]$. Then, we prove that $R \uplus_{res}^* R' \succeq (E', h')$. In order to prove this we need two auxiliary properties:

1. For every variable $z \in \text{dom } E$ such that $E(z) \xrightarrow{w_z} \bullet \xleftarrow{w_v} v$ (in h') there exists a variable $y \in \text{dom } E$ and a word w_y such that $\widehat{E}(y) \xrightarrow{w_y} \bullet \xleftarrow{w_v} \widehat{v}$ (in \widehat{h}') and $E(z) \xrightarrow{w_z} \bullet \xleftarrow{w_y} E(y)$ (in h). This means, by Lemma 9, that the sharing between the result and a variable is captured by $R \uplus_{res}^* R'$.
2. For every w_1, w_2 such that $v \xrightarrow{w_2} \bullet \xleftarrow{w_1} v$ (in h') holds, but $\widehat{v} \xrightarrow{w_2} \bullet \xleftarrow{w_1} \widehat{v}$ (in \widehat{h}') does not, either
 - there exist two variables $y \neq z \in \text{dom } E$ and two words w_y, w_z such that:
 - (a) $E(y) \xrightarrow{w_y} \bullet \xleftarrow{w_z} E(z)$ (in h).
 - (b) $\widehat{E}(y) \xrightarrow{w_y} \bullet \xleftarrow{w_1} \widehat{v}$ (in \widehat{h}').
 - (c) $\widehat{v} \xrightarrow{w_2} \bullet \xleftarrow{w_z} \widehat{E}(z)$ (in \widehat{h}').
 - or, there exists a variable $z \in \text{dom } E$ and words w_v, w_z, w'_1, w'_2 such that
 - (a) $E(z) \xrightarrow{w_z w'_1} \bullet \xleftarrow{w_z w'_2} E(z)$ (in h).
 - (b) $\widehat{v} \xrightarrow{w_v} \bullet \xleftarrow{w_z} \widehat{E}(z)$ (in \widehat{h}').
 - (c) $w_1 = w_v w'_1$ and $w_2 = w_v w'_2$

This means that the internal sharing of res which is not created inside the function, can only come from an argument with internal sharing or from two arguments sharing between them and with the result in the appropriate way. The definition of $R \uplus_{res}^* R'$ also covers this situations, as Lemmas 9 and 10 show.

□

5. IMPLEMENTATION ISSUES AND COST

The analysis presented in the previous section contains some tests and operations that deserve a detailed comment in order to see whether all of them are decidable, and what their costs are.

Since the number of bound variables in a function definition is finite, so is the number of tuples in R . A relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ may occur multiple times in R with different p_1 and p_2 but, as we will see, always with different types. Then, all the set union operations are decidable.

In the implementation, we represent regular languages by non-deterministic finite automata (NFA). We will denote them by $A = (\Sigma, Q, i, F, \delta)$. This facilitates some of the operations needed on regular languages. These are the following:

1. To test whether a regular language L is empty, i.e. $L = \{\}$.
2. Given regular languages L_1 and L_2 , to compute its concatenation $L_1.L_2$.
3. Given regular languages L_1 and L_2 , to compute $L_1|L_2$.
4. Given regular languages L_1 and L_2 , to test whether $L_1 \subseteq L_2$.

The emptiness test can be achieved [6] by looking for a final state that is reachable from the initial one. If $n = |Q|$ is the number of states of A , the algorithm costs $O(n^2)$.

Given NFA automata A_1 and A_2 , the automaton recognizing $L(A_1).L(A_2)$ can be constructed with a cost $O(n)$, just by connecting with ϵ -transitions the final states of A_1 to the initial one of A_2 .

Given NFA automata $A_1 = (\Sigma_1, Q_1, i_1, F_1, \delta_1)$ and $A_2 = (\Sigma_2, Q_2, i_2, F_2, \delta_2)$, the automaton recognizing $L(A_1)|L(A_2)$ is more involved. In fact, we have not found in the literature an algorithm to compute it, and have invented our own:

1. Compute the automaton A'_2 by adding to A_2 transitions with every symbol in $\Sigma = \Sigma_1 \cup \Sigma_2$, from every final state of A_2 to itself. It is clear that A'_2 recognizes $L(A_2).\Sigma^*$.
2. Compute $A_3 = A_1 \cap A'_2$. It recognizes the words of $L(A_1)$ beginning with a word of $L(A_2)$. The construction implies that the states of A_3 are all the pairs of $Q_1 \times Q_2$.
3. Build an automaton A_4 with a fresh state q_0 as the initial one. If $L(A_1) \cap L(A_2)$ is empty, then add ϵ -transitions from q_0 to every state $(q_i, p_j) \in A_3$ such that q_i is a non-final state of A_1 and p_j is a final state of A'_2 . If $L(A_1) \cap L(A_2)$ is not empty, then q_i in $(q_i, p_j) \in A_3$ can be any state of A_1 .
4. Remove from A_4 all the states non-reachable from q_0 . The resulting automaton exactly recognizes $L(A_1)|L(A_2)$.

The dominant costs of the algorithm are the cartesian product and the state reachability computation, both in $O(n^2)$.

Given NFA automata A_1 and A_2 , $L(A_1) \subseteq L(A_2)$ if and only if $L(A_1) \cap L(A_2) = L(A_1)$, so inclusion is a particular case of equality. Unfortunately, equality cannot be directly

computed on NFA's. They must be converted to deterministic finite automata (DFA), and then their equality tested with the well-known table-filling algorithm [6], which has a cost $O(n^2)$. But the conversion from NFA to DFA has a worst-case cost in $O(n^{3 \cdot 2^n})$. This is because the number of states of the DFA are subsets of the NFA set of states, and can in theory be up to 2^n . As we will see, the equality of languages must be tested once every fixpoint iteration.

When interpreting the body of a recursive function f , we start by setting an empty signature for f , i.e. $\Sigma(f) = \emptyset$. It is easy to show that the interpretation is monotonic in the lattice:

$$\langle \mathcal{M}(Var_f \times \mathcal{P}(\Sigma^*) \times \mathcal{P}(\Sigma^*) \times Var_f), \emptyset, \top, \sqsubseteq, \cup, \cap \rangle$$

where \mathcal{M} stands for ‘multiset of’, Var_f are the bound variables of f , Σ^* is the top regular language, and \top is the maximum relation. We need to ensure that no two tuples with the same type exist relating the same variables. So, at the end of each iteration, the following collapsing rule is used:

$$\frac{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y \in R \quad x \xrightarrow{p_3} \bullet \xleftarrow{p_4} y \in R}{type(x, p_1) = type(x, p_3)} \quad OR$$

replace in R the two tuples by $x \xrightarrow{p_1+p_3} \bullet \xleftarrow{p_2+p_4} y$

Should not we use this rule, the abstract domain, regarding only the relations between program variables, would be infinite. The order relation between two tuples relating the same pair of variables, and having the same type, is as follows:

$$x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y \sqsubseteq x \xrightarrow{p'_1} \bullet \xleftarrow{p'_2} y$$

if $L(p_1) \subseteq L(p'_1)$ and $L(p_2) \subseteq L(p'_2)$. Let us call $\mathcal{I}_f \Sigma$ to the interpretation of e_f with current signature environment Σ , returning Σ with f 's signature updated. By monotonicity, we have:

$$\emptyset \sqsubseteq \mathcal{I}_f \emptyset \sqsubseteq \mathcal{I}_f (\mathcal{I}_f \emptyset) \sqsubseteq \dots \sqsubseteq (\mathcal{I}_f)^i \emptyset \sqsubseteq \dots$$

Disregarding the regular languages, this chain is finite because so is Var_f , and the number of different types of the program. Then, the least fixpoint can be reached after a finite number of iterations. If n is the number of f 's formal arguments, then at most n iterations are needed. This is because functional languages have no variable updates, and then there never may arise sharing relations between the formal arguments as a consequence of the function body actions. The only possible relations will be between the function's result and its arguments.

Considering now the regular languages, infinite ascending chains are possible, i.e. one can obtain infinite chains $L_1 \subseteq L_2 \subseteq L_3 \subseteq \dots$.

The least upper bound of such a sequence of regular languages needs not to be a regular one. But, at least, there always exists the regular language Σ^* greater than any other one. In order to ensure termination of the fixpoint computation, we use the following *widening* technique [3]:

1. Based on the form of the automata denoting the increasing language sequence, and by using some heuristics, we guess an automaton A such that $\bigcup_i L_i \subseteq L(A)$. Then, we iterate the interpretation by using this automaton as an assumption in f 's signature.

2. If A is a fixpoint or a post-fixpoint, then we are done. Otherwise, we use Σ^* as the upper bound of the sequence. In terms of precision, $x \xrightarrow{\Sigma^*} \bullet \xleftarrow{\Sigma^*} y$ is completely uninformative about the paths through which x and y share their common descendant.

The heuristic consists of comparing the automata sequence obtained for a given relation $x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y$ in the successive iterations, and discovering growing sequences reaching three or more states related by the same alphabet symbol. For example q_1, q_2, q_3 , with $(q_1, a, q_2), (q_2, a, q_3) \in \delta$. These sequences are collapsed into a single state class q , with a single iterative transition $(q, a, q) \in \delta$. The resulting automata is compared with the non-widened one, to ensure that they are equivalent regarding the remaining transitions. In all the examples we have tried, this heuristic appears to be enough to reach a fixed point.

We pay now attention to the asymptotic cost of the whole interpretation. We choose the size n of a function to be its number of bound variables. This figure is linearly related to the size of its abstract syntax tree, and to the number of lines of its source code. How is n related to the size of the inferred automata in terms of their number of states? It is easy to check that every bound variable y introduces a relation $x \xrightarrow{j_C} \bullet \xleftarrow{\leftarrow} y$ with a prior bound variable x . This increases by one the number of states of the y relations with respect to those of the x relations. So, the automata number of states grow from one to the abstract syntax tree height, when going from the initial expression to the deepest ones. Assuming a reasonably balanced syntax tree, we consider $\log n$ to be an accurate bound to the automata size.

If a function definition has n bound variables, and considering as a constant the number of different types, in the worst case there can be up to $O(n^2)$ tuples in the current relation R . The computation of a single closure operation $R \uplus_x \{x \xrightarrow{p_1} \bullet \xleftarrow{p_2} y\}$ (see Fig. 7) introduces as many relations $x \xrightarrow{p_x} \bullet \xleftarrow{p_z} z$ as prior relations $y \xrightarrow{p_y} \bullet \xleftarrow{p_z} z$ are there in R , i.e. $O(n)$ in the worst case. A single iteration of the abstract interpretation will compute one such closure for every bound variable, giving an upper bound of $O(n^2)$ new relations per iteration. For each one, two languages $A_1|_{A_2}.A_3$ must be computed, giving a total cost of $O(n^2 \log^2 n)$ per iteration.

It has been said that the number of iterations is at most the function's number of arguments, which is usually small. Even if it is not, in practice it suffices to perform only three iterations of the analysis before applying the widening, and then an additional iteration in order to check that the fixpoint has been reached. This checking is the most expensive operation of the analysis. A maximum of $O(n^2)$ languages are tested for equality, giving a total theoretical cost of $O(n^2 2^{\log n} \log^3 n)$ in the worst case, i.e. $O(n^3 \log^3 n)$.

A worst-case cost of $O(n^3 \log^3 n)$ is by no means a low one, but we consider it to be rather pessimistic. We remark that we are assuming each variable to be related to each other, and all conversions from NFA to DFA to produce an exponential blow-up of states. This leads us to think that this theoretical cost is almost never reached. Also, in functional programming it is common to write small functions. So, the number n of bound variables can be expected to remain below 20 for most of the functions (the reader is invited to check this assertion for the functions presented in

```
last xs = case xs of
  x:xx -> case xx of
    [] -> { * R1 * } x
    y:yy -> { * R2 * } last xx
```

Figure 10: Definition of the function `last`

this paper).

In practice our analysis is affordable for medium-size functions. More importantly, it is modular, because once a function definition is analysed, all its relevant information is recorded in the signature environment. Hence, the compilation of a big program is still linear in the program size, even if analysing each individual function of size n takes a time in $O(n^3 \log^3 n)$.

We have implemented the analysis presented here, which has been integrated into our *Safe* compiler, written in Haskell. We have extended the HaLeX library [17], which manipulates regular languages, with new operations such as language intersection, derivation and equality. While the implementation of the abstract interpretation rules of Fig. 6 is straightforward, the closure operation defined in Fig. 7 is much more involved.

Even though the automata library is not particularly efficient and there is much space for optimization, our prototype implementation is able to analyse a file with 40 small functions similar to `msort`, in less than ten seconds in a standard laptop computer.

In order to illustrate the analysis, we present in Fig. 10 the code of a function `last` computing the last element of a non-empty list. By iterating once the interpretation, and in the places marked in the text, we get the following two sets:

$$\begin{aligned} R_1 &= \{xs \xrightarrow{\leftarrow} \bullet \xleftarrow{\leftarrow} xs\} \uplus_x \{xs \xrightarrow{1} \bullet \xleftarrow{\leftarrow} x\} \uplus_{xx} \\ &\quad \{xs \xrightarrow{2} \bullet \xleftarrow{\leftarrow} xx\} \\ R_2 &= R_1 \uplus_y \{xx \xrightarrow{1} \bullet \xleftarrow{\leftarrow} y\} \uplus_{yy} \{xx \xrightarrow{2} \bullet \xleftarrow{\leftarrow} yy\} \end{aligned}$$

Then $\Sigma_1 = \mathcal{I}_{last} \{last \mapsto \emptyset\} = \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{1} xs\}$, where we omit the reflexive relations. By applying again the interpretation, we get:

$$\begin{aligned} \Sigma_2 = \mathcal{I}_{last} \{last \mapsto \Sigma_1\} &= \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{21} xs\} \cup \\ &\quad \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{1} xs\} \\ &= \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{21+1} xs\} \end{aligned}$$

The language 21 is obtained by the closure $\{res \xrightarrow{\leftarrow} \bullet \xleftarrow{1} xx\} \uplus_{res} \{xs \xrightarrow{2} \bullet \xleftarrow{\leftarrow} xx\}$. In the next round, we get $\Sigma_3 = \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{2(21+1)+1} xs\}$. Applying now the widening step, we get $\Sigma_3 = \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{2^*1+21+1} xs\}$, and by applying the interpretation once more:

$$\begin{aligned} \mathcal{I}_{last} \{last \mapsto \Sigma_3\} &= \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{2(2^*1+21+1)} xs\} \cup \\ &\quad \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{1} xs\} \\ &= \{res \xrightarrow{\leftarrow} \bullet \xleftarrow{2(2^*1+21+1)+1} xs\} \end{aligned}$$

The final test is $2(2^*1 + 21 + 1) + 1 \subseteq 2^*1 + 21 + 1$ which returns *true* because all the words in the left language are also in the right one. Notice that the right expression could be further simplified to 2^*1 . This language clearly expresses that the result of `last` is a descendant of the argument list that can be reached by taking the tail of the list a number of times and then by taking the head.

6. CASE STUDIES

Besides the examples already shown in the paper, we have applied our analysis to some additional ones involving list and binary tree manipulations. The following functions show how our analysis can also detect internal sharing in the data structure given as a result. This is useful to know whether a given data structure is laid out in memory without overlapping.

```
buildTree x 0 = Empty
buildTree x n = Node (buildTree x (n-1)) x (buildTree x (n-1))

buildTreeSh x 0 = Empty
buildTreeSh x n = let t = buildTree x (n-1) in Node t x t
```

The shape analysis yields the results given below. We also include the inferred sharing relations of the `append`, `partition` and `qsort` functions that make up a typical *Quicksort* implementation:

$$\begin{aligned}
 \text{buildTree } x \ n & : \{ \text{res} \xrightarrow{(1+3)^*2} \bullet \leftarrow^\epsilon x \\
 & \quad , \text{res} \xrightarrow{(1+3)^*2} \bullet \xleftarrow{(1+3)^*2} \text{res} \} \\
 \text{buildTreeSh } x \ n & : \{ \text{res} \xrightarrow{(1+3)^*2} \bullet \leftarrow^\epsilon x \\
 & \quad , \text{res} \xrightarrow{(1+3)^*2} \bullet \xleftarrow{(1+3)^*2} \text{res} \\
 & \quad , \text{res} \xrightarrow{(1+3)^*} \bullet \xleftarrow{(1+3)^*} \text{res} \} \\
 \text{append } xs \ ys & : \{ \text{res} \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} xs, \text{res} \xrightarrow{2^*} \bullet \leftarrow^\epsilon ys \} \\
 \text{partition } p \ xs & : \{ \text{res} \xrightarrow{12^*1+22^*1} \bullet \xleftarrow{2^*1} xs \} \\
 \text{qsort } xs & : \{ \text{res} \xrightarrow{2^*1} \bullet \xleftarrow{2^*1} xs \}
 \end{aligned}$$

7. RELATED WORK AND CONCLUSIONS

There exist many different analyses dedicated to extracting information about the heap, mainly in imperative languages where pointers are explicitly used and may be reassigned. *Alias analysis* is one of the most studied. It tries to detect program variables that point to the same memory location. *Pointer analysis* aims at determining the storage locations a pointer can point to, so it may be also used to detect aliases in a program. These analyses are used in many different applications such as live variable analysis for register allocation and constant propagation. In [13, 5, 14] we can find surveys about pointer analysis applied to imperative languages from the 80's. Related to these analyses, an *Escape Analysis* tries to determine statically the dynamic scope of the data structures that will be created at runtime, whereas *Shape analysis* [16, 9, 15] tries to approximate the 'shape' of the heap-allocated structures. That information has been used, for example, for binding time optimizations.

The level of detail of all these analyses mainly depends on the user of the analysis. Our analysis tries to capture a kind of sharing information more refined than alias and pointer analysis may provide, and in fact both are subsumed in our relations: if $x \xrightarrow{\epsilon} \bullet \leftarrow^\epsilon y$, then, x and y are aliases; if $x \xrightarrow{j} \bullet \leftarrow^\epsilon y$, then x points to y (i.e. y is the j -th child of the data structure x). In the area of escape analysis, Blanchet [2] applies the concept of paths in order to determine which pointers in a data structure survive the current execution scope. The sets of paths are subsequently abstracted by integer numbers denoting escape contexts, whereas in this work we use regular expressions for abstracting those sets. Moreover, our analysis aims to infer sharing relations between our structures. That is why shape analysis is nearer to our needs.

Jones and Muchnick [9] associate sets of k -limited graphs to each program point in order to approximate the sharing relations between variables. The k limits the length of the paths in the graphs modeling the heap in order to make the domain finite and obtain the minimal fixpoint by iteration. The graphs obtained after the abstract execution of a program instruction must be transformed in order to maintain themselves k -limited. Our widening operator resembles this operation. Our path relations are in general uncomparable in precision to these sets of limited graphs. First, having sets of graphs may provide more precision because our union operation loses information: adding $x \xrightarrow{p_1+p_3} \bullet \xleftarrow{p_2+p_4} y$ introduces combinations of paths $x \xrightarrow{p_1} \bullet \xleftarrow{p_4} y$ and $x \xrightarrow{p_3} \bullet \xleftarrow{p_2} y$ which did not exist previously. Second, paths longer than k may be more precise than k -limited graphs: $x \xrightarrow{2221} \bullet \leftarrow^\epsilon y$ indicating that y is the fifth element of the list x is more precise than saying in a 2-limited graph that y shares in an *unknown* way with x after the path 22. Additionally, the cost of having sets of graphs is doubly exponential in the number of variables.

In order to reduce the cost to polynomial, Reps [15] formulated the analysis as a graph-reachability problem over the dependence graph generated from the program. The reachability is defined in terms of those (context-free) paths one is interested in. The fixpoint calculation in this case is also finite because he just records the information about the variables, not the exact paths. We need the paths in order to make the analysis more precise as shown in the *mergesort* example, that is why we need the widening. The use of context-free paths in our framework would make undecidable most of our tests.

Other related works are those devoted to compile-time garbage collection, such as [7, 8]. The first one tries to save creating a new array when updating an array that is only referenced once. The second one provides an analysis also detecting when a cell is referenced at most once by the subsequent computation. Its aim is to destroy the cell after its last use so that it can be reused by the runtime system. Both analyses are done on a first-order eager functional language. After these ones, there have been many similar analyses, usually known as *usage analyses* (e.g. [18, 1, 19, 4]) whose aim is to detect when a cell is used at most once and then, either to recover or to avoid to update it, when the language is lazy. These analyses do not try to know *which* other data structures points to a particular cell, but rather *how many* of them do it, and in this sense they are simpler. The nearer to our problem is [8] since it pursues an aim similar to that of *Safe*: to save memory. The main difference is that, in our case, it is the programmer who decides to destroy a cell and the compiler just analyses whether doing this is safe or not. So, the programmer may have destructive and non-destructive versions of the same function and uses the first one in contexts where it is safe to do it. In [8] it is the compiler who decides to destroy the cell, when it is safe to do it in *all* the contexts in which the function is called. A single unsafe context will avoid to recover the cell in all the safe ones. Another important difference is that our analysis is modular, while theirs need to analyse the program as a whole. This makes it unpractical for big programs.

8. REFERENCES

- [1] E. Barendsen and S. Smetsers. Uniqueness typing for functional languages with graph rewriting semantics.

- Mathematical Structures in Computer Science*, 6(6):579–612, Dec. 1996.
- [2] B. Blanchet. Escape analysis for JavaTM: Theory and practice. *ACM Transactions on Programming Languages and Systems*, 25(6):713–775, Nov. 2003.
- [3] P. Cousot and R. Cousot. Abstract interpretation: A unified lattice model for static analysis of programs by construction of approximation of fixed points. In *Proc. 4th ACM Symp. on Principles of Prog. Languages*, pages 238–252. ACM, 1977.
- [4] J. Gustavsson and J. Sveningsson. A Usage Analysis with Bounded Usage Polymorphism and Subtyping. In *Selected Papers of the 12th International Workshop on Implementation of Functional Languages, IFL'00*, volume 2011 of LNCS, pages 140–157. Springer-Verlag, 2001.
- [5] M. Hind. Pointer analysis: Haven't we solved this problem yet? In *Proceedings of the 2001 ACM SIGPLAN-SIGSOFT Workshop on Program Analysis For Software Tools and Engineering, PASTE'01*, pages 54–61. ACM Press, 2001.
- [6] J. E. Hopcroft, R. Motwani, and J. D. Ullman. *Introduction to Automata Theory, Languages and Computation*. Addison Wesley, 2nd edition, 2001.
- [7] P. Hudak. A Semantic Model of Reference Counting and its Abstraction (Detailed Summary). In *ACM Symposium on Lisp and Functional Programming*, pages 351–363. ACM, 1986.
- [8] T. P. Jensen and T. A. Mogensen. A Backwards Analysis for Compile-Time Garbage Collection. In *European Symposium on Programming*, pages 227–239. LNCS 432, Springer, 1990.
- [9] N. D. Jones and S. S. Muchnick. Flow analysis and optimization of lisp-like structures. In *Proceedings of the 6th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, POPL '79, pages 244–256. ACM, 1979.
- [10] M. Montenegro, R. Peña, and C. Segura. A Type System for Safe Memory Management and its Proof of Correctness. In *ACM Principles and Practice of Declarative Programming, PPDP'08, Valencia, Spain, July. 2008*, pages 152–162, 2008.
- [11] M. Montenegro, R. Peña, and C. Segura. Shape Analysis in a Functional Language by Using Regular Languages (Extended Version). Technical report, TR-8-13. Dpto. de Sistemas Informáticos y Computación. Universidad Complutense de Madrid, 2013. Available at: <http://federwin.sip.ucm.es/sic/investigacion/publicaciones/informes-tecnicos>.
- [12] R. Peña, C. Segura, and M. Montenegro. A Sharing Analysis for SAFE. In *Selected Papers Trends in Functional Programming, TFP'06*, pages 109–128. Intellect, 2007.
- [13] V. Raman. Pointer analysis – a survey. CS203 UC Santa Cruz, <http://www.soe.ucsc.edu/~vishwa/publications/Pointers.pdf>, 2004.
- [14] D. Rayside. Points-to analysis. <http://www.cs.utexas.edu/~pingali/CS395T/2012sp/lectures/points-to.pdf>, 2005.
- [15] T. Reps. Shape analysis as a generalized path problem. In *Proceedings of the 1995 ACM SIGPLAN symposium on Partial evaluation and semantics-based program manipulation, PEPM '95*, pages 1–11. ACM, 1995.
- [16] J. C. Reynolds. Automatic computation of data set definitions. In *IFIP Congress (1)*, pages 456–461, 1968.
- [17] J. Saraiva. HaLeX: A Haskell Library to Model, Manipulate and Animate Regular Languages. In *Proc. ACM Workshop on Functional and Declarative Programming in Education*, University of Kiel. Tech. Report 0210, pages 133–140, 2002.
- [18] D. N. Turner, P. L. Wadler, and C. Mossin. Once upon a type. In *7th International Conference on Functional Programming and Computer Architecture*, pages 1–11, La Jolla, California, June 1995. ACM Press.
- [19] K. Wansbrough and S. L. P. Jones. Once upon a polymorphic type. In *The Twenty-sixth ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, San Antonio, Texas, January 1999.