

# Static Type Analysis of Pattern Matching by Abstract Interpretation

Pietro Ferrara

► **To cite this version:**

Pietro Ferrara. Static Type Analysis of Pattern Matching by Abstract Interpretation. John Hatcliff; Elena Zucca. Joint 12th IFIP WG 6.1 International Conference on Formal Methods for Open Object-Based Distributed Systems (FMOODS) / 30th IFIP WG 6.1 International Conference on Formal Techniques for Networked and Distributed Systems (FORTE), Jun 2010, Amsterdam, Netherlands. Springer, Lecture Notes in Computer Science, LNCS-6117, pp.186-200, 2010, Formal Techniques for Distributed Systems. <10.1007/978-3-642-13464-7\_15>. <hal-01055154>

**HAL Id: hal-01055154**

**<https://hal.inria.fr/hal-01055154>**

Submitted on 11 Aug 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Static Type Analysis of Pattern Matching by Abstract Interpretation

Pietro Ferrara

ETH Zurich, Switzerland  
pietro.ferrara@inf.ethz.ch

**Abstract.** Pattern matching is one of the most attractive features of functional programming languages. Recently, pattern matching has been applied to programming languages supporting the main current object oriented features. In this paper, we present a static type analysis based on the abstract interpretation framework aimed at proving the exhaustiveness of pattern matchings and the safety of type casts. The analysis is composed by two distinct abstract domains. The first domain collects information about dynamic typing, while the second one tracks the types that an object cannot be instance of. The analysis has been implemented and applied to all the `Scala` library. The experimental results underline that our approach scales up since it analyzes a method in 90 msec in average. In addition, the analysis is precise in practice as well, since we prove the exhaustiveness of 42% of pattern matchings and 27% of the type casts without any manual annotation on the code.

**Key words:** Abstract interpretation, static analysis, pattern matching.

## 1 Introduction

Pattern matching is recognized to be one of the most expressive features of functional programming languages. Extending its full expressiveness to programming languages supporting the current object oriented features is not straightforward [12]. In fact, features like inheritance and information hiding are not supported by the functional pattern matching, since it dealt with algebraic data types. Nevertheless, recent work [14, 26] introduced pattern matching in programming languages supporting the main current object oriented features as `F#` [1] and `Scala` [23]. In addition, some extensions of `Java` supporting pattern matching appeared in the last few years [16, 24]. Pattern matching checks if an expression respects a given pattern looking at an ordered list of case expressions. One of the most common case expressions is the type case: we select the case to apply following the dynamic type of the expression. Consider the `Scala` program in Figure 1. The content of argument `x` is matched with respect to its type: if it is an instance of `Something`, the method returns the string representing the contained value. Otherwise, the type of `x` is `None`, and `extract` returns the empty string.

Languages like F# and Scala are compiled into a bytecode language (for instance, MSIL [11] or Java bytecode [18]) that supports the main imperative features. Thus type cases in pattern matching are translated into type tests and casts. For instance, the body of method `extract` in Figure 1 is translated into the code in Figure 2 by the Scala compiler [13]. Note that information about generics is erased. This code can be easily optimized adopting a specific static

```
abstract sealed class Option[T]
final case class Something[T](val y : T) extends Option[T]
final case object None extends Option[Nothing]

def extract [T](x : Option[T]) : String = x match {
  case Something(y) => return y.toString();
  case None => return "";
}
```

Fig. 1: Pattern matching with type cases

```
1 String extract (Option x) {
2     if (x instanceof Something)
3         return ((Something) x).y.toString ();
4     else if (x instanceof None)
5         return "";
6     else throw new MatchError ();
7 }
```

Fig. 2: Results of Scala compilation

analysis. When `x instanceof Something` is false, we know that `x instanceof None` is always true. In fact, the static type of `x` is `Option`, this class is `abstract`, and the only two classes that extend it are `Something` and `None`. Note that `Option` is declared as `sealed`. In Scala, a `sealed` class can be extended only by classes that are defined in the same source file. Therefore we know that `Option` cannot be extended by external code, and we can conclude that the `if` statement at line 4 is always evaluated to `true`. In addition, we know that the statement `throw new MatchError()` is unreachable, that is the pattern matching contained by the original Scala program is exhaustive. Finally, the type cast at line 3 is safe.

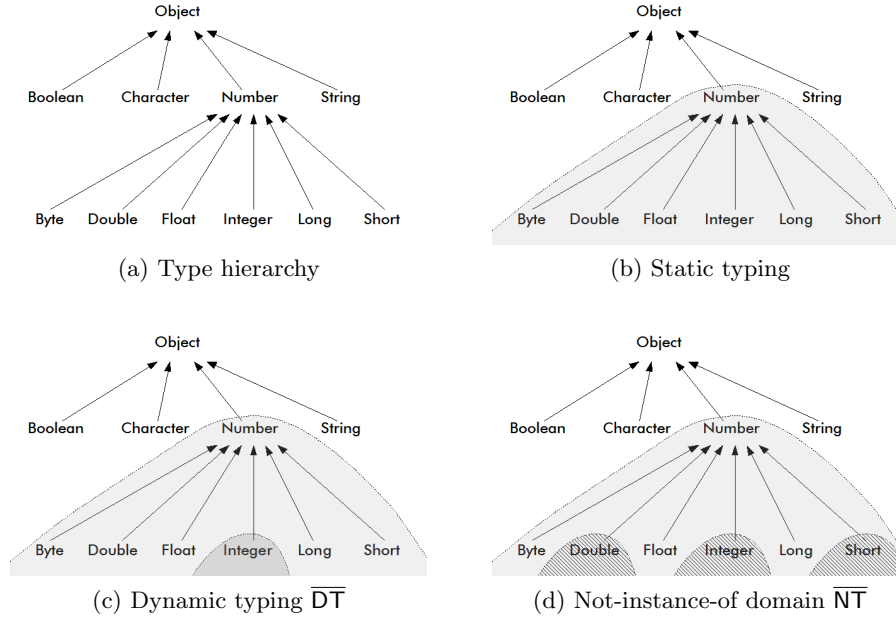


Fig. 3: Type system and abstract domains

## 1.1 Contribution

The main contribution of this paper is the introduction of two new abstract domains in order to capture precise information on the types of variables. The two main goals of our analysis are precision and efficiency. Precision is achieved through the development of abstract domains focused on the properties of interest, that is, exhaustiveness of pattern matching and safety of casts. Efficiency is achieved through a modular analysis. For this reason, our approach is based on the Design-by-Contract (DbC) methodology [21]. This means that when we analyze a method call we rely on class invariants, pre and post conditions. Thanks to this approach, we do not need to analyze the whole program, but we can analyze a method alone relying on its contracts. In a similar way, we rely on language constructs like `sealed` in `Scala` in order to be sure that a class cannot be extended by external code.

In this paper, we do not deal with issues related to the heap and its analysis. Hence the language on which we are going to define our analysis does not contain field accesses, aliasing, etc., but it deals only with variables. The object oriented programs we analyzed are preprocessed by a sound heap analysis that replaces heap accesses with variables. This means that method calls may have side effects on these variables. The heap analysis we perform is quite rough: intuitively, we approximate all the instances of the same class with the same variable, and we perform weak updates [4] as these abstract references may represent several

concrete references. Then the number of abstract references is bounded to the number of classes of the analyzed program. In this way we preserve the scalability of our approach. On the other hand, it may seem that this approach is too much coarse since we approximate all the fields having the same static type with the same variable. Indeed, this does not particularly affect the precision of our approach because, for the most part of the programs we analyzed, the content of fields is first stored in a local variable (and in this way we obtain a specific variable identifier), and then type checking and casts are performed on this local variable.

Consider now the type hierarchy depicted in Figure 3a. The class `Object` is extended by several subclasses that represent different types in a programming language, e.g., `Boolean`. The abstract class `Number` is extended by several subclasses, each representing one possible numerical type, e.g., `Integer`. For instance, Figure 3b represents a situation in which the static type of a variable is `Number`. The static type information guarantees that the objects assigned to a variable are instances of its static type or of one of its subclasses. Thus instances of `Long`, `Float`, etc. may be assigned to a variable whose static type is `Number`.

The first abstract domain we are going to define is  $\overline{DT}$ . It is aimed at approximating information about the dynamic type of variables. Figure 3c depicts a case in which we assign an object instance of the `Integer` class (darker area) to a variable of static type `Number` (light grey area). The information about dynamic typing comes mainly from assignments, type casts, `new` statement, and assumption of `instanceof` conditions. Since we rely on DbC, if a method has no postcondition about the dynamic type of variables (e.g., the returned value), we would abstract its type information using its static type after the method call.

The second domain we are going to define is  $\overline{NT}$ . It abstracts the types that a variable cannot be instance of. For instance, Figure 3d depicts a situation in which a variable of static type `Number` cannot be instance of `Double`, `Integer`, or `Short`. In this way we can discard some leaves of the type tree. This information is collected when we assume the negation of an `instanceof` boolean condition. We could have adopted a different approach capturing a set of types for each variable representing all the possible types a variable can have at a given program point. In this way, we may model that an object cannot be instance of a type simply removing the type from the set of types related to it. Unluckily this solution is not feasible for two main reasons. Imagine to be after a method call, where we know only the static type of a variable. In this context, assuming that an object cannot be instance of a given type (i) may affect the performance (e.g., if its static type is `Object` and we have to exclude `String` we should relate the variable to a set of types of huge dimensions), (ii) cannot be represented since our analysis is modular and external code usually defines and instantiates new classes. One may argue that, since we are interested in proving the exhaustiveness of pattern matching, we can restrict the analysis only on types that cannot be extended by external code. With such restriction the approach based on set of possible types would be feasible, and performance could be preserved supposing that we have few subclasses. First of all, we want to track information also on classes

that can be externally extended, in order to eventually provide information to developers like “if this class would be `sealed`, then the pattern matching would be exhaustive”. This solution would not collect the information necessary to provide this output. In addition, we do not expect in general that `sealed` classes (or more generally classes that cannot be extended by external code) have only few subclasses, since each subclass may be used to represent a particular case of a complex data structure.

Combining  $\overline{DT}$  and  $\overline{NT}$ , we are in position to obtain precise information about types and prove the safety of some cast operations, and the unreachability of some code, in particular when dealing with pattern matching.

The analysis works on a simple object oriented language. We designed it in order to translate the main object oriented programming languages into this language. In this way, we want to apply our analysis to several languages. Up to now, we have developed the translation of the `Scala` programming language. Thus our analysis is focused on the code obtained from the `Scala` pattern matching statements. In this context, if we prove that a `MatchError` exception is unreachable, this means that the pattern matching is exhaustive.

The analysis has been implemented and the experimental results underline that our approach scales up. We are in position to analyze `Scala` libraries (more than 30.000 methods) in less than one hour, and in average we require 90 milliseconds to analyze a single method. In addition, it turns out that the analysis is precise in practice. In fact, we are able to prove the exhaustiveness of 42% of pattern matchings (that may contain not only type cases but also other cases whose information is not captured by our analysis) and 27% of the type casts without any manual annotation of the code.

The following of this section introduces some preliminaries on abstract interpretation. Section 2 presents the language that we analyze. Section 3 introduces the concrete domain and semantics, while Section 4 formalizes the two domains adopted by our analysis, and explains how they work on an example. Section 5 presents the experimental results, and Section 6 introduces the related work. Finally, Section 7 concludes and discusses the future work.

## 1.2 Abstract Interpretation

Our approach is based on abstract interpretation [7, 8], a theory for defining and soundly approximating the semantics of a program. A concrete semantics, aimed at specifying the runtime properties of interest, is defined. Then it is approximated to obtain an abstract semantics computable but still precise enough to capture the property of interest, and specified by an abstract domain and an abstract transition function.

Formally, the concrete domain  $D$  is a lattice  $\langle D, \sqsubseteq_C, \perp_C, \top_C, \sqcup_C, \sqcap_C \rangle$ . The concrete elements are related to the abstract domain  $\langle \overline{D}, \sqsubseteq_A, \perp_A, \top_A, \sqcup_A, \sqcap_A \rangle$  by a concretization function  $\gamma$  and an abstraction function  $\alpha$ . In order to obtain a sound analysis, we require that they form a Galois connection; we denote it by  $\langle D, \sqsubseteq_C, \perp_C, \top_C, \sqcup_C, \sqcap_C \rangle \xrightarrow[\alpha]{\gamma} \langle \overline{D}, \sqsubseteq_A, \perp_A, \top_A, \sqcup_A, \sqcap_A \rangle$ .

$C ::= x = E$	$(C1)$	$E ::= \text{new } T() \text{ (E1)}$	$(E1)$	$B ::= x \text{ instanceof } T \text{ (B1)}$	$(B1)$
$\text{declare } x : T \text{ (C2)}$	$(C2)$	$x$	$(E2)$	$! x \text{ instanceof } T \text{ (B2)}$	$(B2)$
$\text{assume}(B) \text{ (C3)}$	$(C3)$	$(T) x$	$(E3)$		
$x.M(y1, \dots, yn) \text{ (C4)}$	$(C4)$				

Table 1: OO core language

A semantics is defined on the abstract and concrete domains. Given a state of the domain and a statement, it produces the result of the execution of the statement starting from the given state. The abstract semantics  $\bar{\mathbb{S}}$  has to soundly approximate the concrete one, that is  $\forall \bar{d} \in \bar{D} : \mathbb{S}[\llbracket \gamma(\bar{d}) \rrbracket] \sqsubseteq_A \gamma \circ \bar{\mathbb{S}}[\bar{d}]$ . Usually, the semantics is defined as the computation of a fixpoint [8].

Different domains can be combined in Cartesian products. Let  $\bar{A} \times \bar{B}$  be the Cartesian product of abstract domains  $\bar{A}$  and  $\bar{B}$ . The lattice operators are defined as the pointwise application on the components of the pairs of the lattice operator on each abstract domain. The reduced product  $\bar{A} \otimes \bar{B}$  is a Cartesian product on which a *reduce* function is provided: given a state of the Cartesian product, it mutually refines the information contained in each domain using the information provided by the other domain.

**Generic analyzers** Abstract interpretation can be applied in order to develop generic analyzers [5]. In particular, this theory allows one to define a compositional analysis, e.g., an analysis that can be instantiated with different numerical domains, and in order to analyze different properties. Many different generic analyzers have been proposed recently [15, 19, 25]. The type analysis we present in this paper has been implemented in **Sample** (Static Analyzer of Multiple Programming Languages). This generic analyzer can be plugged with different numerical domains and heap analysis. In addition, it can be extended with new analyses like our type analysis.

## 2 Language

We introduce a core language in order to make explicit the main characteristics of our abstract domains. It consists of assignments (C1), declarations of variables (C2), assumptions of boolean conditions (C3), and method calls (C4). An expression that can be assigned to a variable is the instantiation of a class (E1), another variable (E2), or the cast of a variable (E3). A condition can be the check that the dynamic type is instance of (B1) or it is not instance of (B2) a type.

We focus our attention only on the main aspects of the type analysis we are interested in, that is, casts and runtime checks of dynamic types. Nevertheless, our implementation covers all the features of current programming languages, as numerical operations, field accesses, etc. Intuitively, we represent the body of a

$\mathbb{E}[\text{new } T(), d] = T$ $\mathbb{E}[x, d] = d(x)$ $\mathbb{E}[(T) x, d] =$ $= \begin{cases} d(x) & \text{if } d(x) \sqsubseteq_{\overline{T}} T \\ \perp & \text{otherwise} \end{cases}$	$\mathbb{B}[x \text{ instanceof } T, d] =$ $= d(x) \sqsubseteq_{\overline{T}} T$ $\mathbb{B}[\neg B, d] = \neg \mathbb{B}[B, d]$	$\mathbb{S}[x = E, d] = d[x \mapsto \mathbb{E}[E, d]]$ $\mathbb{S}[\text{declare } x : T, d] = d[x \mapsto T]$ $\mathbb{S}[\text{assume}(B), d] =$ $= \begin{cases} \perp & \text{if } \mathbb{B}[B, d] = \text{false} \\ d & \text{if } \mathbb{B}[B, d] = \text{true} \end{cases}$ $\mathbb{S}[x.M(y_1, \dots, y_n), d] =$ $= \text{execute}(x.M(y_1, \dots, y_n), d)$
---	--	--

Table 2: Collecting semantics

method through a Control Flow Graph (CFG). Each block in the CFG contains a list of statements. Different blocks are connected through edges that could eventually contain a boolean condition to represent conditional jumps. In this way, we support control structures like if statements and while loops.

### 3 Concrete Domain and Collecting Semantics

The concrete domain is composed by environments that collect the dynamic type for each variable. Formally,  $\text{Env} : [\text{Var} \rightarrow \overline{T}]$ . The concrete domain is composed by sets of environments. Each environment represents one possible execution of the program in a given point. The lattice is obtained using set operators. Formally,  $(\wp(\text{Env}), \subseteq, \cup, \cap, \text{Env}, \emptyset)$ .

We define the collecting semantics as a function that, given an environment and a statement, returns the environment resulting from the execution of the given statement on the given environment. Table 2 reports the formal definition of this collecting semantics<sup>1</sup>.  $\sqsubseteq_{\overline{T}}$  corresponds to the subtype relation, while  $\perp$  is used to represent a situation in which the execution is stopped because of an unsafe dynamic cast. The collecting semantics is aimed at formalizing the runtime behaviors focusing on the information we are interested. The extension of this semantics to our concrete domain (that is made by set of environments) is the application of this semantics on each environment in the initial state.

$\mathbb{S}$  is the basic step adopted in order to compute the semantics of a method on its CFG representation. Intuitively, we build up incrementally the traces representing the executions of a method relying on  $\mathbb{S}$  to perform single computational steps [6]. Since non-deterministic behaviors are possible (e.g., because of inputs from the user) and these may also affect the type information, we consider sets of traces.

### 4 Abstract Domain and Semantics

In this section, we present and formalize the two abstract domains of our analysis, how we combine them, and we explain how they work in practice on the example

<sup>1</sup>  $\text{execute}(x.M(y_1, \dots, y_n), d)$  resolves and executes method  $M$  on  $d$



presented in Figure 2. Note that we will define the abstract semantics only on single statements. Its extension to blocks and CFGs is obtained as usual in approaches based on abstract interpretation and trace semantics relying on the upper bound operators [7].

#### 4.1 Static Typing

The programming language introduced in Section 2 provides some information on the static types of variables. We suppose that a subtype relation  $\sqsubseteq_{\bar{T}}$  is provided, and that this is a partial ordering. Given two types  $\bar{t}_1, \bar{t}_2 \in \bar{T}$ ,  $\bar{t}_1 \sqsubseteq_{\bar{T}} \bar{t}_2$  if and only if  $\bar{t}_1$  is subtype of  $\bar{t}_2$ . We denote by  $\sqcup_{\bar{T}}$  and  $\sqcap_{\bar{T}}$  respectively the upper and lower bound operators univocally identified by  $\sqsubseteq_{\bar{T}}$ . These operators form a lattice. In this way, we support the most part of common object oriented type systems [3], like the ones of `Java`, `C#`, `Scala`, and `F#`.

Each type represents itself and all the types that are its subtype. Note that this is exactly the semantics of a static type: at runtime a variable of static type  $\bar{t}$  can have a type that is  $\bar{t}$  or one of its subtypes. Formally, the concretization of a type  $\bar{t}$  is defined as  $\gamma_{\bar{T}}(\bar{t}) = \{\bar{t}' : \bar{t}' \sqsubseteq_{\bar{T}} \bar{t}\}$ . Since information about static typing can be considered as syntactic sugar, we suppose that a function *statictype* :  $[\text{Var} \rightarrow \bar{T}]$  is provided. Given a variable, it returns its static type.

#### 4.2 $\overline{\text{DT}}$ : Dynamic Typing

$\overline{\text{DT}}$  abstracts the dynamic type of variables. Relying on  $\bar{T}$  we build up the domain  $\overline{\text{DT}} : [\text{Var} \rightarrow \bar{T}]$  that relates each variable to its dynamic type. The operators on the lattice  $\langle \overline{\text{DT}}, \sqsubseteq_{\overline{\text{DT}}}, \sqcup_{\overline{\text{DT}}}, \sqcap_{\overline{\text{DT}}}, \top_{\overline{\text{DT}}}, \perp_{\overline{\text{DT}}} \rangle$  are obtained as the functional extension of the ones of  $\bar{T}$ . The information inferred by  $\overline{\text{DT}}$  is refined with the static type information provided by the *statictype* function: if the type of a variable  $v$  is not yet defined, the dynamic type domain returns its static type, that is, *statictype*( $v$ ). The concretization of  $\overline{\text{DT}}$  is defined as the functional extension of  $\gamma_{\bar{T}}$  as well. Formally,  $\gamma_{\overline{\text{DT}}}(\bar{f}) = \{[v \mapsto \bar{t}] : v \in \text{dom}(\bar{f}) \wedge \bar{t} \in \gamma_{\bar{T}}(\bar{f}(v))\}$ .

#### Semantics

$$\begin{aligned}
\overline{\text{DT}}[x = y, \bar{dt}] &= \bar{dt}[x \mapsto \bar{dt}(y)] \\
\overline{\text{DT}}[x = \text{new } T(), \bar{dt}] &= \bar{dt}[x \mapsto T] \\
\overline{\text{DT}}[x = (T)y, \bar{dt}] &= \bar{dt}[x \mapsto T \sqcap_{\bar{T}} \bar{dt}(y), y \mapsto T \sqcap_{\bar{T}} \bar{dt}(y)] \\
\overline{\text{DT}}[\text{assume}(x \text{ instanceof } T), \bar{dt}] &= \bar{dt}[x \mapsto \bar{dt}(x) \sqcap_{\bar{T}} T] \\
\overline{\text{DT}}[x.M(y_1, \dots, y_n), \bar{dt}] &= \overline{\text{DT}}[\text{assume}(\text{rename}(\text{postcondition}_M)), \top_{\overline{\text{DT}}}]
\end{aligned}$$

When a variable is assigned to another one ( $x = y$ ), we capture that the dynamic type of  $x$  is the type of  $y$ . Similarly, when we assign to a variable the instance of a fresh object, we relate this variable to the type of the new object. The dynamic type of the assignment of a cast to type  $T$  of a variable  $y$  is the lower bound between the dynamic type of  $y$  and type  $T$ . Thus, it may be the case in

which it contains a type that is more approximated than  $\top$  even if the cast is safe. After the cast, we know that its dynamic type will be surely  $\top$  or one of its subtypes. This is why we take the lower bound between the dynamic type of the variable and  $\top$ . In a similar way, when we test to `true` a boolean condition like `x instanceof T`, we track that the dynamic type of `x` is the lower bound of its type and  $\top$ . When we analyze a method call, we assume the postconditions of this method call, that may contain some type information, on the top state of the  $\overline{DT}$  domain. Since we do not consider framing information, we suppose that potentially a method may access and assign all the variables of the current state of computation. Then we assume the postcondition on a top state in order to preserve the soundness of our analysis. In addition, we suppose that a `rename` function is provided, and it renames all the variables contained in the postcondition matching the calling environment.

**Running Example** At the beginning of the analysis of method `extract` we have no information about dynamic typing, but we know that the static type of `x` is `Option`. When we analyze the condition of the `if` branch at line 2 we calculate that the state leading to the `then` statement of our dynamic type domain is  $[x \mapsto \text{Something}]$ , since the condition is evaluated to `true`. In the same way, at line 4 we obtain that  $[x \mapsto \text{None}]$ . Thus we prove that the cast at line 3 is safe, while we are not able yet to detect that line 6 is unreachable.

### 4.3 $\overline{NT}$ : Not-instance-of Domain

The second domain we introduce is  $\overline{NT}$ . This domain is aimed at collecting the types an object cannot be instance of at a given point of the program.

First of all, we define a lattice that collects sets of types. The intuition behind them is that they contain all the types of which a variable cannot be instance of. Thus let  $\overline{NT}$  be this domain, i.e.,  $\overline{NT} = \wp(\overline{T})$ . Its concretization contains all the types that are not in the given set or that are subtypes of one of the types in the set. Formally,  $\gamma_{\overline{NT}}(\overline{f}) = \{[v \mapsto \overline{t} : v \in \text{dom}(\overline{f}) \wedge \nexists \overline{t}' \in \overline{f}(v) : \overline{t} \sqsubseteq_{\overline{T}} \overline{t}']\}$ . The ordering operator on  $\overline{NT}$  is not the subset operator, as more types we have, more precise we are. On the other hand, it is neither the superset operator, as several types may be subtypes of the same type, and they are different elements of the set. Thus the ordering operator is defined as

$$\overline{nt}_1 \sqsubseteq_{\overline{NT}} \overline{nt}_2 \Leftrightarrow \forall v \in \text{dom}(\overline{nt}_2) : v \in \text{dom}(\overline{nt}_1) \wedge \forall \overline{t}_2 \in \overline{nt}_2(v) : \exists \overline{t}_1 \in \overline{nt}_1(v) : \overline{t}_2 \sqsubseteq_{\overline{T}} \overline{t}_1$$

The other lattice operators are the ones induced by  $\sqsubseteq_{\overline{NT}}$ . The bottom element for a single variable is the set containing only the top type (e.g., `Object` in Java or `Any` in Scala). The concretization of this element would be all the types that are not subtype of the top type, i.e., the empty set. We denote by  $\perp_{\overline{NT}}$  a function that relates all the variables to  $\{\top_{\overline{T}}\}$ . Formally,  $\perp_{\overline{NT}} = \lambda x. \{\top_{\overline{T}}\}$ . The top element is the empty function. Formally,  $\langle \overline{NT}, \sqsubseteq_{\overline{NT}}, \perp_{\overline{NT}}, \sqcap_{\overline{NT}}, \sqcup_{\overline{NT}}, \emptyset \rangle$ .

## Semantics

$$\begin{aligned}
\overline{\text{NT}}[x = y, \overline{\text{nt}}] &= \overline{\text{nt}}[x \mapsto \overline{\text{nt}}(y)] \\
\overline{\text{NT}}[x = \text{new } T, \overline{\text{nt}}] &= \overline{\text{nt}}[x \mapsto \emptyset] \\
\overline{\text{NT}}[x = (T)y, \overline{\text{nt}}] &= \overline{\text{nt}}[x \mapsto \overline{\text{nt}}(y)] \\
\overline{\text{NT}}[\text{assume}(! x \text{ instanceof } T, \overline{\text{nt}})] &= \overline{\text{nt}}[x \mapsto \overline{\text{nt}}(x) \cup \{T\}] \\
\overline{\text{NT}}[x.M(y_1, \dots, y_n), \overline{\text{dt}}] &= \overline{\text{DT}}[\text{assume}(\text{rename}(\text{postcondition}_M)), \top_{\overline{\text{NT}}}]
\end{aligned}$$

We are interested in collecting which types a variable is not instance of. When a variable  $y$  is assigned to  $x$  (with or without cast), we track that the types  $x$  cannot be instance of are the same of  $y$ . When we assign a new instance of a class to a variable, we forget everything we knew about that variable. In this way, we lose some information: we could consider that the variable cannot be instance of all the types that are not subtype of the instantiated class. On the other hand, this would be computationally expensive, the same information is already contained in  $\overline{\text{DT}}$ , and we are going to combine this domain with  $\overline{\text{NT}}$ .  $\overline{\text{NT}}$  is also interested to approximate information when a boolean condition like  $x \text{ instanceof List}$  is tested to false. For instance, when we analyze the if statement  $\text{if}(x \text{ instanceof List}) C1; \text{ else } C2$ , we know that in the else branch  $x$  cannot be instance of `List`, and we add `List` to the state of the  $\overline{\text{NT}}$  domain before analyzing  $C2$ . The method call is managed as in  $\overline{\text{DT}}$ : method calls may have side effects on variables, arbitrarily changing their dynamic type. In order to preserve the soundness of our approach, the only information we have after the method call is that the postcondition of the called method holds.

**Running Example** The  $\overline{\text{NT}}$  domain models information on the example presented in Figure 2 when we test to `false` conditions containing `instanceof`. Thus at line 4 (i.e., the else branch of the if statement at line 2) its state is  $[x \mapsto \{\text{Something}\}]$ . At line 6 (i.e., the else branch of the if statement at line 4) its state is  $[x \mapsto \{\text{Something}, \text{None}\}]$ . Unluckily,  $\overline{\text{NT}}$  alone is not enough in order to prove that line 6 is unreachable, as it knows nothing about the static or dynamic types of  $x$ .

### 4.4 $\overline{\text{FT}}$ : Reduced Product of $\overline{\text{DT}}$ and $\overline{\text{NT}}$

The  $\overline{\text{DT}}$  domain infers the dynamic type of a variable.  $\overline{\text{NT}}$  collects the types a variable cannot be instance of. Thus these domains model different types of information, but we can mutually refine them. In particular, we are interested in checking if the information contained by  $\overline{\text{NT}}$  about a variable is not compatible with the one contained in  $\overline{\text{DT}}$ . In this case, we know that this point of the program is unreachable, and thus we can refine our domain to the bottom state. For instance, consider an abstract sealed class `I` (thus it cannot be instantiated nor extended by external code) that is implemented by two classes `O1` and `O2`. If we know from  $\overline{\text{DT}}$  that the dynamic type of a variable  $x$  is `I`, but from  $\overline{\text{NT}}$  we know that  $x$  cannot be instance of `O1` nor `O2`, then there is no possible type for  $x$ .

Line	$\overline{DT}$	$\overline{NT}$	Reduced $\overline{FT}$
2	$[x \mapsto \text{Option}]$	$\emptyset$	
3	$[x \mapsto \text{Something}]$	$\emptyset$	
4	$[x \mapsto \text{Option}]$	$[x \mapsto \{\text{Something}\}]$	
5	$[x \mapsto \text{None}]$	$[x \mapsto \{\text{Something}\}]$	
6	$[x \mapsto \text{Option}]$	$[x \mapsto \{\text{Something}, \text{None}\}]$	$\perp_{\overline{FT}}$

Table 3: Results of the analysis of the example

Formally, let be  $\overline{FT}$  the reduced product of  $\overline{DT}$  and  $\overline{NT}$ , i.e.,  $\overline{FT} = \overline{DT} \otimes \overline{NT}$ . The lattice operators of  $(\overline{FT}, \sqsubseteq_{\overline{FT}}, \sqcup_{\overline{FT}}, \sqcap_{\overline{FT}}, \top_{\overline{FT}}, \perp_{\overline{FT}})$  are defined as usual when dealing with the product of domains. The reduction function  $\overline{red}_{\overline{FT}}$  is defined as

$$\overline{red}_{\overline{FT}} : [\overline{FT} \rightarrow \overline{FT}]$$

$$\overline{red}_{\overline{FT}}((\overline{d}, \overline{n})) = \begin{cases} \perp_{\overline{FT}} & \text{if } \exists x \in \text{dom}(\overline{d}) : \overline{getType}(x, \overline{d}) = \overline{t} \wedge \\ & \forall \overline{t}_1 \in \gamma_{\overline{NT}}(\overline{t}) : \overline{abstract}(\overline{t}_1) = \text{false} \wedge \exists \overline{t}_2 \in \overline{n} : \overline{t}_1 \sqsubseteq_{\overline{NT}} \overline{t}_2 \wedge \\ & \quad \overline{t} \text{ cannot be extended by external code} \\ (\overline{d}, \overline{n}) & \text{otherwise} \end{cases}$$

where  $\overline{abstract} : [\overline{T} \rightarrow \{\text{true}, \text{false}\}]$  is a function that, given a type, checks if the given type cannot be instantiated, e.g., it is an interface or an abstract class in Java, or a trait in Scala.

The reduction function is aimed at discovering if the information contained in  $\overline{NT}$  is not compatible with the information contained in  $\overline{DT}$ . It returns  $\perp_{\overline{FT}}$  if and only if (i) all the types a variable can be instance of following the information contained in  $\overline{DT}$  are excluded at least by one type contained in  $\overline{NT}$ , (ii) the type related to the variable in  $\overline{DT}$  cannot be extended by external code (e.g., it is declared as `sealed` in Scala). In this way we achieve the modularity in our analysis. Note that this reduction is partial, as there may be other ways to refine the information contained in the different domains. Anyway, we did not find useful other reductions of the information contained in this domain in order to improve the precision of our analysis.

The concretization of  $\overline{FT}$  is defined as the intersection of the pointwise application of  $\gamma_{\overline{DT}}$  and  $\gamma_{\overline{NT}}$ . Formally,  $\gamma_{\overline{FT}}((\overline{d}, \overline{n})) = \gamma_{\overline{DT}}(\overline{d}) \cap \gamma_{\overline{NT}}(\overline{n})$ . The semantics  $\overline{\mathbb{F}T}$  is defined as the pointwise application of  $\overline{\mathbb{D}T}$  and  $\overline{\mathbb{N}T}$ .

**Theorem 1.** *The fixpoint trace semantics based on  $\overline{\mathbb{F}T}$  is sound with respect to the one based on  $\mathbb{S}$ .*

**Running Example** Table 3 reports the results of the  $\overline{FT}$  domain when applied to method `extract` of the example presented in Figure 2. We already pointed out how the information is captured by the abstract domains. The reduced product refines the information contained in the two domains at line 6 discovering that

it is not possible for variable `x` to be instance of `Option`, and not to be instance of `Something` and `None` at the same time. In fact, these two classes represent all the possible instances of `Option`. Note that `Option` is declared as `sealed`, hence it cannot be extended by external code. Therefore the reduced product is able to discover that line 6 is unreachable, and `MatchError` cannot be thrown. Previously, we discovered that  $\overline{DT}$  proved that the cast at line 3 is safe, thus, in a Java environment, a `ClassCastException` cannot be thrown. Using this information, we may improve the runtime performance (for instance removing the last `if` statement since the `else` branch is unreachable) and provide useful information to developers in order to debug programs, since we guarantee that pattern matching is exhaustive in the original Scala code.

## 5 Experimental Results

	#m	t	tm	Casts			MatchErrors		
				v	nv	p	v	nv	p
actors	1626	3'56"	145.35	47	104	31.13%	20	17	54.05%
collection	14578	10'21"	42.63	183	509	26.45%	42	74	36.21%
util	4926	8'21"	101.75	126	508	10.87%	36	65	35.64%
xml	2786	6'21"	136.59	108	286	27.41%	5	21	19.23%
mainlib	8218	12'02"	87.82	97	196	33.11%	37	12	75.51%
scala lib	35732	54'02"	90.72	725	1951	27.09%	156	208	42.86%

We applied our analysis to all the Scala library v. 2.7.7. We executed the analysis on an Intel Code 2 Quad CPU 2.83 GHz with 4 GB of RAM, running Windows 7, and the Java SE Runtime Environment 1.6.0\_16-b01. The preceding table reports the experimental results. Column `#m` reports the number of analyzed methods. Column `t` reports the time required to analyze these methods<sup>2</sup>, while `tm` reports the time required in average to analyze a single method in milliseconds. We report the number of validated cases (column `v`), not validated (`nv`) cases, and the overall precision (`p`) of our analysis when analyzing both type casts and reachability of `MatchError` exceptions.

The Scala library contains more than 35.000 methods. In average we analyze a method in 90 milliseconds. Thus our analysis is quite efficient, and it scales up. When we analyze the type casts, we are not able to distinguish between casts due to the compilation of pattern matchings, and other casts, since our analysis works on the code obtained after the transformations and simplifications performed by the Scala compiler. Hence our analysis takes into account all the casts contained in the program. We are able to automatically prove safe more than 27% of all the casts. We think that this result goes beyond our initial goal (that was, to precisely analyze pattern matching), since we proved the safety of more than 700 casts while there were only 364 `MatchError` exceptions (that means that there were 364 pattern matchings). On the other hand, in general a precision less than 30% is not particularly satisfying. We analyzed these warnings, and it turns out that

<sup>2</sup> This time takes into account also the heap analysis

the most part is due to erasure of type generics of `Scala` compiler. For instance, when we access the next element of a list of `Integer` objects, this is compiled into a program that takes the next element of a list of `Object` instances and casts it to `Integer`. Obviously, since we do not have information about generics, we cannot prove the safety of this code. We plan to apply our analysis to languages that preserve generic type information in order to study how this feature can improve the precision of our analysis. A minor part of the warnings is induced by the lack of contracts. Since the `Scala` library is not annotated, each time we have a method call we forget everything is contained in  $\overline{DT}$  and  $\overline{NT}$ . We expect the we may improve the precision of our analysis adding some contracts concerning framing information in particular.

About the reachability of `MatchError`, we proved that about 43% of these exceptions is unreachable. This result is particular encouraging: even if we took into account a part of the information that can be used in pattern matching and we do not annotate the code with contracts, we proved that almost half of the existing pattern matchings is safe, and we can remove the `MatchError` exception. Usually we were not able to prove the exhaustiveness of pattern matchings when they contain some checks on numerical information. In a minor part of the cases, we failed to prove the exhaustiveness because the pattern matching was not exhaustive since it expected to receive a particular type of expression. Contracts would be the solution in order to express that a method expects that an argument can have only some specific types (e.g., not all the ones that are subtype of its static type), thus we expect that we could improve the precision adding this annotation.

## 6 Related Work

The related work is addressed mainly into two directions: the static analysis of type information in order to optimize virtual calls, and static analysis of pattern matching.

The first topic has been studied during the last 15 years. Object oriented languages introduced the idea of virtual calls: we do not know at compile time exactly which method we are calling, but we need some runtime information (for instance, the dynamic type of the object on which we are calling the method). Thus the binding of the called method is dynamic, and this step may require an overhead. Applying static analysis to reduce and bound the number of virtual calls improves the runtime performance of programs. In this context, many approaches were proposed [2]. Some of them performs an inter-procedural analysis [17], while our approach relies on contracts. Other approaches analyze the class hierarchy statically [9], and check if the virtual call can refer only to one implementation of the method in this context. These approaches usually consider information only about dynamic types.

A recent work [20] proposed a type analysis based on the abstract interpretation theory in order to optimize JavaScript code. This analysis is focused on the numerical information that can be assigned to a variable in order to infer

if we can adopt an `Int32` type instead of a `Float64` type. Using this information the code is optimized. This approach is focused on the numerical information and adopts this information to infer the types of variables, while our approach approximates information on the types of variables to check which casts are safe and if some statements are unreachable.

Recent work applied static analysis techniques to prove properties of pattern matching. Mitchell and Runciman [22] proposed an analysis to check if non-exhaustive pattern matching in Haskell could lead to failures. In particular, it infers preconditions that are strong enough to prove that the pattern matching never fails. Our approach is aimed at discovering which pattern matchings are non-exhaustive, instead of inferring which constraints are necessary in order to make exhaustive a pattern matching that is not.

Dotta *et al.* [10] introduced a verification system that checks the disjointness, reachability, and exhaustiveness of `Scala` pattern matching. Global information about the program is inferred through different kinds of formulas. The authors do not formalize how these formulas are inferred, and there is no proof about the soundness of the approach. Instead, we fully formalized our approach and we proved its soundness. On the other hand, their analysis tracks more information than ours. Their implementation requires up to some seconds to prove the exhaustiveness of some case studies. Our analysis tracks information only of type-based pattern matchings, but it scales up.

## 7 Conclusion and Future Work

In this paper we presented a static type analysis focused on pattern matching. We introduced and combined two distinct abstract domains, each abstracting different information (dynamic typing, and not-instance-of information). We proved the soundness of our approach relying on the abstract interpretation framework. The analysis has been implemented in `Sample`, and the experimental results proved the scalability and the precision of our approach.

As future work, we plan to combine our analysis with some numerical domains in order to study if considering this information improves the precision of our analysis. In addition, we plan to exploit the information captured by our analysis to optimize the programs resulting from the `Scala` compiler, and to study how much this optimization may improve the runtime performance of some `Scala` benchmarks.

## References

1. F#. <http://research.microsoft.com/fsharp>.
2. D. F. Bacon and P. F. Sweeney. Fast static analysis of C++ virtual function calls. In *OOPSLA '96*, pages 324–341. ACM, 1996.
3. L. Cardelli. Type systems. In A. B. Tucker, editor, *The Computer Science and Engineering Handbook*, chapter 97. CRC Press, 2004.

4. D. R. Chase, M. Wegman, and F. K. Zadeck. Analysis of pointers and structures. In *PLDI '90*. ACM, 1990.
5. P. Cousot. The calculational design of a generic abstract interpreter. In *Calculational System Design*. NATO ASI Series F. IOS Press, Amsterdam, 1999.
6. P. Cousot. Constructive design of a hierarchy of semantics of a transition system by abstract interpretation. In *Theoretical Computer Science*, volume 277, pages 47–103. Elsevier Science publishers, 2002.
7. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL '77*. ACM, 1977.
8. P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *POPL '79*. ACM, 1979.
9. J. Dean, D. Grove, and C. Chambers. Optimization of object-oriented programs using static class hierarchy analysis. In *ECOOP '95*. Springer-Verlag, 1995.
10. M. Dotta, P. Suter, and V. Kuncak. On static analysis for expressive pattern matching. Technical report, EPFL, 2008.
11. S. ECMA-335. *Common Language Infrastructure (CLI)*. ECMA, 4th edition, June 2006.
12. B. Emir. *Object-oriented pattern matching*. PhD thesis, EPFL, 2007.
13. B. Emir, Q. Ma, and M. Odersky. Translation correctness for first-order object-oriented pattern matching. In *APLAS '07*, LNCS, pages 54–70. Springer, 2007.
14. B. Emir, M. Odersky, and J. Williams. Matching objects with patterns. In *ECOOP '07*, pages 273–298. Springer, 2007.
15. P. Ferrara. Checkmate: a generic static analyzer of java multithreaded programs. In *Proceedings of SEFM '09*. IEEE Computer Society, 2009.
16. M. Hirzel, N. Nystrom, B. Bloom, and J. Vitek. Matchete: Paths through the pattern matching jungle. In *PADL '08*, LNCS. Springer, 2008.
17. R. Kumar and S. S. Chakraborty. Precise static type analysis for object oriented programs. *SIGPLAN Not.*, 42:17–26, 2007.
18. T. Lindholm and F. Yellin. *Java Virtual Machine Specification*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1999.
19. F. Logozzo and M. Fähndrich. On the relative completeness of bytecode analysis versus source code analysis. In *Proceedings of CC '08*, LNCS. Springer-Verlag, 2008.
20. F. Logozzo and H. Venter. Rata: Rapid atomic type analysis by abstract interpretation. Application to JavaScript optimization. In *CC '10*, LNCS. Springer, 2010.
21. B. Meyer. *Object-Oriented Software Construction (2nd Edition)*. Prentice Hall, 1997.
22. N. Mitchell and C. Runciman. Not all patterns, but enough: an automatic verifier for partial but sufficient pattern matching. *SIGPLAN Not.*, 44:49–60, 2009.
23. M. Odersky. *The Scala Language Specification*, 2008.
24. A. Richard and O. Lhotak. Oomatch: pattern matching as dispatch in Java. In *OOPSLA '07*, pages 771–772. ACM, 2007.
25. F. Spoto. JULIA: A Generic Static Analyser for the Java Bytecode. In *Proceedings of FTfJP'2005*, 2005.
26. D. Syme, G. Neverov, and J. Margetson. Extensible pattern matching via a lightweight language extension. *SIGPLAN Not.*, 42(9):29–40, 2007.