

Abstract Domains for Type Juggling

Vincenzo Arceri¹

Department of Computer Science, University of Verona, Italy

Sergio Maffeis²

Department of Computing, Imperial College London, UK

Abstract

Web scripting languages, such as PHP and JavaScript, provide a wide range of dynamic features that make them both flexible and error-prone. In order to prevent bugs in web applications, there is a sore need for powerful static analysis tools. In this paper, we investigate how Abstract Interpretation may be leveraged to provide a precise value analysis providing rich typing information that can be a useful component for such tools.

In particular, we define the formal semantics for a core of PHP that illustrates *type juggling*, the implicit type conversions typical of PHP, and investigate the design of abstract domains and operations that, while still scalable, are expressive enough to cope with type juggling. We believe that our approach can also be applied to other languages with implicit type conversions.

Keywords: PHP, Static analysis, Abstract interpretation, Type conversions

1 Introduction

The success of web scripting languages such as PHP and JavaScript is also due to their wide range of dynamic features, which make them very flexible but unfortunately also error-prone. A key such feature is that language operations allow operands of any type, applying implicit type conversions when a specific type is needed. PHP, our example language, calls this feature *type juggling*.

In this paper, we investigate how the Abstract Interpretation approach to program analysis [3,4] may be leveraged to provide a precise value analysis in presence of type juggling. Since PHP is dynamically typed, meaning that the same variable can store values of different types at different points in the execution, our analysis does not aim to enforce type invariance, but instead aims to determine the most precise type for each variable in the final state.

Filaretti and Maffeis [6] define a formal operational semantics for most of the PHP language that is faithful to its mainstream Zend reference implementation [1]. In Section 2, we propose μPHP (*micro-PHP*), a much smaller core of the language that is still large enough to illustrate the main challenges related to type juggling. In fact, μPHP is valid PHP, and behaves exactly like the full language³, although the omission of certain language features from our formalisation (see Section 5) allows us to define a more straightforward semantics than the one in [6]. We present μPHP

¹ Email: vincenzo.arceri@studenti.univr.it

² Email: sergio.maffeis@imperial.ac.uk

³ All the examples in the paper are both derivable via our semantics and executable in PHP 5.4.

in *big-step* semantics style, as we are interested in properties of the final state.⁴ We show many examples that will reveal surprising behaviour of PHP to the non-expert.

In Section 3, we define an abstract semantics parametric on the domain, which defines a corresponding *flow-* and *path-sensitive* value analysis. We discuss assumptions on such domain under which we can argue that the analysis is sound with respect to the concrete semantics of μPHP . The design of our semantics makes it straightforward to implement an abstract interpreter to calculate the analysis result.

In Section 4, we define abstract domains and operations that capture the subtleties of type juggling. Rather than giving the definitions upfront, we expound the rationale behind our design, stressing expressivity, modularity and hopefully highlighting subtle points that can be useful to design domains for other languages with similar features. Some practical static analyses of realistic languages with dynamic type conversions, such as [9, 11], add to each type lattice extra points that represent information which can improve the precision of the analysis. Other analyses, such as [8], use powersets of values, limiting the set sizes by a parameter k in order to avoid infinite computations. That leads to very expressive domains when up-to- k values are analysed, that drastically loose precision for further values.

In contrast, we advocate an expressive and systematic approach that refines each type domain to include just the information necessary to obtain precise abstract operations and type juggling functions. Our analysis may not be highly efficient but is scalable, having polynomial complexity: we emphasise precision over performance. As argued in [4], in theory one should aim for the *best correct approximation* of a concrete operator f defined as $f^\sharp = \alpha \circ f \circ \gamma$, but f^\sharp is sometimes not computable, or practical. In defining the abstract operations of our type juggling domain we follow the spirit of this equation, striving to exploit at most the concrete information available, and delay as much as possible the loss of information caused by merging values with the \sqcup operator.

Related Work. Since the seminal work of [2], abstract interpretation has been used to define many value and type analyses, but we are not aware of any analysis designed to handle in particular the implicit type conversions for scripting languages. On the practical side, several static analysers for JavaScript and PHP are directly based, or at least inspired, by abstract interpretation [5, 8–12]. All aim to analyse real-world PHP programs, and focus most effort on prominent issues such as the analysis of associative arrays and functions, while paying less attention to implicit type conversions. As far as we can tell (sometimes essential details are missing from the cited references), none of the analyses in [5, 9–12] comes close to our level of precision, except for [8] which, as discussed above, uses expensive powerset domains. Nevertheless, we hope that our investigation may contribute to improve the precision of these analysers for programs that make intensive use of implicit type conversions. Moreover, none of the cited works above provides formal proofs of soundness, and some such as [10, 12] openly admit to be unsound.

Summarising, our main claim of novelty is to apply a systematic approach grounded in the theory of Abstract Interpretation to analyse, in a provably sound way, non-trivial features of (the core of) a practical programming language.

⁴ It would be easy, but notationally more cumbersome, to define an equivalent *small-step* semantics better able to represent trace properties.

2 Type Juggling in μPHP

We now define syntax and semantics of μPHP , a subset of PHP able to express most type juggling behaviour. Our examples can be verified in a PHP 5.x interpreter.

2.1 Syntax

To appreciate some subtle points of type juggling, we need to be somewhat precise about the representation of literals. Let CHAR be the finite set of characters used in PHP, and $\text{DIG} \subsetneq \text{CHAR}$ the set of digits $0, \dots, 9$. The literals of μPHP are partitioned in the sets

- **NULL**: the constant `NULL`, which is the default value of undefined variables.
- **BOOL**: the boolean constants `true` and `false`.
- **STR = CHAR***: strings such as `"hi!"`, `""`, `"bye!"`.
- **INT = $-^? \text{DIG}^+$** : signed integers such as `-5`, `0`, `1`, `00042`.
- **FLOAT = $-^? \text{DIG}^* . \text{DIG}^*$** : decimal notation numbers,⁵ such as `-1.3`, `0.`, `4.200`.

The capitalisation of `NULL`, `true`, and `false` above is irrelevant. An empty sequence of digits between the optional sign and the decimal point of a float is interpreted as 0, so for example `-.3` is an alternative representation for `-0.3`, and the degenerate case `."` is not a valid `FLOAT`. The syntax of μPHP is reported below:

<pre> Lit ::= NULL BOOL STR INT FLOAT Exp ::= Lit Var ① Exp Exp ② Exp </pre>	<pre> Var ::= \$ID Block ::= { } { Stmt } Stmt ::= Var = Exp ; if (Exp) Block else Block while (Exp) Block Stmt Stmt ; </pre>
---	---

where `ID` is a subset of `STR` suitable to define identifiers. We denote prefix unary operators by $\textcircled{1} \in \{!, -, +\}$ and infix binary operators by $\textcircled{2} \in \{+, -, *, /, \%, \&\&, ||, ==, !=, >, <, >=, <= \}$.

2.2 Semantics

Semantic values correspond to literals, but abstract away from representation details. In particular, leading zeros are dropped when parsing an `INT`, except for the literal `0`, and leading and trailing zeros are dropped when parsing a `FLOAT`, so `-004.20` is the semantic float `-4.2`. With a slight abuse of notation, we use the same font to denote literal and values, as the meaning should be clear from the context. `STR`, `INT`, and `FLOAT` are finite sets, and floating point numbers have limited pre-

⁵ PHP floats normally use the IEEE 754 double precision format. For simplicity, we use instead decimal numbers in μPHP .

cision. We denote by NUM the union $\text{INT} \cup \text{FLOAT}$, and by VAL the union of all the semantic values above. For any set S and $X \subseteq S$, we also define the notation \overline{X} for the complement of X with respect to S .

Program states $\text{STATE} : \text{ID} \rightarrow \text{VAL}$, ranged over by σ , are partial functions from identifiers to values. State updates and lookups are defined as follows:

$$\sigma[x \leftarrow v](y) = \begin{cases} v & \text{if } x = y \\ \sigma(y) & \text{otherwise} \end{cases}$$

Statements. The big-step semantics of blocks and statements is defined by the function $\llbracket \cdot \rrbracket : \text{Stmt} \times \text{STATE} \rightarrow \text{STATE}$ defined below

$$\begin{aligned} \llbracket \$x = e; \rrbracket \sigma &= \sigma[x \leftarrow \llbracket e \rrbracket \sigma] \\ \llbracket \text{if } (e) \text{ b11 else b12} \rrbracket \sigma &= \begin{cases} \llbracket \text{b11} \rrbracket \sigma & \text{if } \text{toBool}(\llbracket e \rrbracket \sigma) = \text{true} \\ \llbracket \text{b12} \rrbracket \sigma & \text{if } \text{toBool}(\llbracket e \rrbracket \sigma) = \text{false} \end{cases} \\ \llbracket \text{while } (e) \text{ b1} \rrbracket \sigma &= \llbracket \text{if } (e) \{ \text{b1 while } e \text{ b1 } \} \text{ else } \{ \} \rrbracket \sigma \\ \llbracket \{ S \} \rrbracket \sigma &= \llbracket S \rrbracket \sigma \\ \llbracket \{ \} \rrbracket \sigma &= \llbracket ; \rrbracket \sigma = \sigma \\ \llbracket S1 S2 \rrbracket \sigma &= \llbracket S2 \rrbracket (\llbracket S1 \rrbracket \sigma) \end{aligned}$$

All the rules are standard except for the if-else, which contains the first example of type juggling, where the value resulting from evaluating the guard expression e in state σ is then automatically converted to a boolean, using the function toBool defined below, where $\text{NUM}_0 = \{0, 0.0\}$, $\text{STR}_{\text{false}} = \{\text{"", "0"}\}$.

$$\text{toBool}(v) = \begin{cases} v & \text{if } v \in \text{BOOL} \\ \text{false} & \text{if } v \in \text{NULL} \cup \text{NUM}_0 \cup \text{STR}_{\text{false}} \\ \text{true} & \text{if } v \in \overline{\text{NUM}_0} \cup \overline{\text{STR}_{\text{false}}} \end{cases}$$

This leads us to our first example of odd behaviour in PHP:

```
php> if (0) {echo "yes";} else {echo "no";} // "no"
php> if ("0") {echo "yes";} else {echo "no";} // "no"
php> if (0.0) {echo "yes";} else {echo "no";} // "no"
php> if ("0.0") {echo "yes";} else {echo "no";} // "yes"
```

Expressions. The semantics of expressions is given by the function $\llbracket \cdot \rrbracket : \text{Exp} \times \text{STATE} \rightarrow \text{VAL}$ which we describe case-by-case below. The semantics of a literal is just the corresponding parsed value, as described at the beginning of this Section. The variable rule returns the value of the corresponding identifier, if it is defined in the current state, and NULL otherwise.

$$\llbracket \$x \rrbracket \sigma = \begin{cases} \sigma(x) & \text{if } x \in \text{dom}(\sigma) \\ \text{NULL} & \text{otherwise} \end{cases}$$

Arithmetic operations are defined on any type of operands:

$$\llbracket e1 \text{ @ } e2 \rrbracket \sigma = \text{toNum}(\llbracket e1 \rrbracket \sigma) \text{ @ } \text{toNum}(\llbracket e2 \rrbracket \sigma)$$

where the operands are converted to numbers (integers or floats) via another type juggling function `toNum`. Let $\text{parseNum} : \text{STR} \rightarrow (\text{NUM} + \{\perp\}) * \text{STR}$ be a function that returns the number that can be parsed as the largest prefix of a string (if any), and the remainder of the string that does not contribute to parsing the number. For example, $\text{parseNum}(".42000.37\text{hi}") = (0.42, ".37\text{hi}")$ and $\text{parseNum}(\text{"bye666"}) = (\perp, \text{"bye666"})$. The function `toNum` is defined by

$$\text{toNum}(v) = \begin{cases} v & \text{if } v \in \text{INT} \cup \text{FLOAT} \\ 1 & \text{if } v = \text{True} \\ 0 & \text{if } v \in \text{NULL} \cup \{\text{False}\} \\ 0 & \text{if } \text{parseNum}(v) = (\perp, v) \\ n & \text{if } \text{parseNum}(v) = (n, s) \text{ for some } s \end{cases}$$

When $\textcircled{2} \in \{+, -, *\}$, $\boxed{2}$ corresponds to the most precise corresponding primitive operation between integers and floats (denoted by $\{+, -, *\}$). So, for example:

```
php> var_dump(3.2*"hi" + 45 - "3bye"*true); // float(42)
```

When $\textcircled{2} \in \{/, \%\}$ instead, $\boxed{2}$ implements a μPHIP -specific function that returns `false` when division by zero occurs.

$$n_1 \boxed{/} n_2 = \begin{cases} n_1/n_2 & \text{if } n_2 \in \overline{\text{NUM}_0} \\ \text{false} & \text{if } n_2 \in \text{NUM}_0 \end{cases}$$

The semantics of comparison operators is tricky, as it depends on the type of the operands. For example, to compare a string with a boolean, first it is converted to a boolean, and then both booleans are compared after being converted to numbers, leading to the perhaps surprising example below.

```
php> var_dump("0" < true); // bool(true)
php> var_dump("0.0" < true); // bool(false)
```

More formally, we define the semantics for the less-than operator as follows (the other comparison operators follow a similar pattern):

$$\llbracket e1 < e2 \rrbracket \sigma = \llbracket e1 \rrbracket \sigma \boxed{<} \llbracket e2 \rrbracket \sigma$$

When $e1$ and $e2$ reach final values v_1 and v_2 , the semantics rules reported in Figure 1 are applied, where $<$ is the primitive operator of less-than for numbers, and $<_{\text{STR}}$ is a non-standard comparison between strings. If two strings can be parsed exactly as numbers, they are compared using $<$ on the parsed numbers; otherwise, they are compared in the lexicographic order $<_L$.

$$s_1 <_S s_2 = \begin{cases} n_1 < n_2 & \text{if } \text{parseNum}(s_1) = (n_1, "") \text{ and } \text{parseNum}(s_2) = (n_2, "") \\ s_1 <_L s_2 & \text{otherwise} \end{cases}$$

This leads to more surprising behaviour. For example,

```
php> var_dump("10" < "9"); // bool(false)
php> var_dump("10LOW" < "9HIGH"); // bool(true)
```

\sqsubseteq	INT	FLOAT	BOOL
INT	$v_1 < v_2$	$v_1 < v_2$	$\text{toNum}(\text{toBool}(v_1)) < \text{toNum}(v_2)$
FLOAT	$v_1 < v_2$	$v_1 < v_2$	$\text{toNum}(\text{toBool}(v_1)) < \text{toNum}(v_2)$
BOOL	$\text{toNum}(v_1) < \text{toNum}(\text{toBool}(v_2))$	$\text{toNum}(v_1) < \text{toNum}(\text{toBool}(v_2))$	$\text{toNum}(v_1) < \text{toNum}(v_2)$
STR	$\text{toNum}(v_1) < v_2$	$\text{toNum}(v_1) < v_2$	$\text{toNum}(\text{toBool}(v_1)) < \text{toBool}(v_2)$
NULL	$\text{toNum}(v_1) < v_2$	$\text{toNum}(v_1) < v_2$	$\text{toNum}(v_1) < \text{toNum}(v_2)$

\sqsubseteq	STR	NULL
INT	$v_1 < \text{toNum}(v_2)$	$v_1 < \text{toNum}(v_2)$
FLOAT	$v_1 < \text{toNum}(v_2)$	$v_1 < \text{toNum}(v_2)$
BOOL	$\text{toNum}(v_1) < \text{toNum}(\text{toBool}(v_2))$	$\text{toNum}(v_1) < \text{toNum}(v_2)$
STR	$\text{toNum}(v_1) < v_2$	$\text{toNum}(\text{toBool}(v_1)) < \text{toBool}(v_2)$
NULL	$\text{toStr}(v_1) <_{\text{STR}} v_2$	false

Figure 1. Tables with semantics rules for the less-than operator applied to basic values

```
php> var_dump(0+"10LOW"<"9HIGH"); // bool(false)
```

where the use of `+` in the third example forces the use of `toNum` on the first string (hence on the second one too), and the use of `<` instead of `<STR` in the comparison.

The semantics of string concatenation is defined as follows

$$\llbracket e1.e2 \rrbracket \sigma = \text{toStr}(\llbracket e1 \rrbracket \sigma) \sqsubseteq \text{toStr}(\llbracket e2 \rrbracket \sigma)$$

where \sqsubseteq is the primitive operation of string concatenation. The type juggling function `toStr` is defined below, where $\text{FLOAT}_{\text{INT}} = \text{INT}?.0^*$ (excluding the degenerate case “.”) represents the floats that can be interpreted as integers without approximation, such as `.00`, `42.`, `0.0`. When an element of $\text{FLOAT}_{\text{INT}}$ is concatenated with a string, only its integer part is concatenated.

$$\text{toStr}(v) = \begin{cases} "1" & \text{if } v = \text{true} \\ "" & \text{if } v = \text{false} \\ "v" & \text{if } v \in \text{INT} \cup \overline{\text{FLOAT}_{\text{INT}}} \\ "u" & \text{if } v \in \text{FLOAT}_{\text{INT}} \text{ and } u = \text{floor}(v) \\ v & \text{if } v \in \text{STR} \end{cases}$$

Above, $\text{floor} : \text{NUM} \rightarrow \text{INT}$ rounds down its argument to the nearest integer.

3 Abstract Interpretation of μPHIP

Our goal is to design an efficient value analysis that retains precise information on the type of variables. Hence, our *concrete domain* representing the properties of interest is the standard complete lattice $\langle 2^{\text{VAL}}, \subseteq \rangle$. With the above goal in mind, we now define an abstract semantics for μPHIP that is parametric in the choice of an abstract domain of values $\langle \text{VAL}^\sharp, \sqsubseteq \rangle$.

3.1 Abstract Semantics

Our analysis is non-relational, hence we can somewhat simplify the design of the abstract semantics and the definition of its soundness properties. In particular, abstract program states $\text{STATE}^\sharp : \text{ID} \rightarrow \text{VAL}^\sharp$, ranged over by ξ , can partition the available information *per identifier*, and be defined as partial functions from identifiers to abstract values. State updates and lookups are defined as for the concrete semantics.

Statements. The abstract semantics of blocks and statements $\llbracket \cdot \rrbracket^\sharp : \text{Stmt} \times \text{STATE}^\sharp \rightarrow \text{STATE}^\sharp$ is similar to the concrete one.

$$\begin{aligned} \llbracket \$x = e; \rrbracket^\sharp \xi &= \xi[x \leftarrow \llbracket e \rrbracket^\sharp \xi] \\ \llbracket \{ s \} \rrbracket^\sharp \xi &= \llbracket s \rrbracket^\sharp \xi \\ \llbracket \{ \} \rrbracket^\sharp \xi &= \llbracket ; \rrbracket^\sharp \xi = \xi \\ \llbracket s1 \ s2 \rrbracket^\sharp \xi &= \llbracket s2 \rrbracket^\sharp (\llbracket s1 \rrbracket^\sharp \xi) \end{aligned}$$

The rules for assignment, blocks and sequences are analogous to the ones for the concrete semantics. Note that in particular we are considering *strong* updates to the state: our analysis is *flow-sensitive*.

$$\llbracket \text{if } (e) \text{ b11 else b12} \rrbracket^\sharp \xi = \begin{cases} \llbracket \text{b11} \rrbracket^\sharp \xi & \text{if } \gamma(\text{toBool}^\sharp(\llbracket e \rrbracket^\sharp \xi)) = \{\text{true}\} \\ \llbracket \text{b12} \rrbracket^\sharp \xi & \text{if } \gamma(\text{toBool}^\sharp(\llbracket e \rrbracket^\sharp \xi)) = \{\text{false}\} \\ \llbracket \text{b11} \rrbracket^\sharp \xi \sqcup \llbracket \text{b12} \rrbracket^\sharp \xi & \text{if } \gamma(\text{toBool}^\sharp(\llbracket e \rrbracket^\sharp \xi)) \supseteq \text{BOOL} \end{cases}$$

The rule for if-else is *path-sensitive*, mimicking the concrete one, yet includes a conservative extra case when the evaluation of the guard does not result in a precise boolean value. It relies on an abstract type juggling function toBool^\sharp which is to be defined together with the abstract domain $\langle \text{VAL}^\sharp, \sqsubseteq \rangle$, as discussed in Section 4.2.

$$\llbracket \text{while } (e) \text{ b1} \rrbracket^\sharp \xi = \text{lf}_\rho(\lambda \rho. (\rho \sqcup \llbracket \text{if } (e) \text{ b1 else } \{ \} \rrbracket^\sharp \rho))$$

The rule for while loops in the concrete semantics can be equivalently formulated as $\llbracket \text{while } (e) \text{ b1} \rrbracket \sigma = \text{lf}_\rho(\llbracket \text{if } (e) \text{ then b1 else } \{ \} \rrbracket)$: the abstract rule is simply a conservative approximation, whose computability depends on the definition of abstract domain.⁶

Expressions. The abstract evaluation of expressions is denoted by $\llbracket \cdot \rrbracket^\sharp : \text{Exp} \times \text{STATE}^\sharp \rightarrow \text{VAL}^\sharp$. Literal values are simply abstracted by the rule $\llbracket \text{Lit} \rrbracket^\sharp \xi = \alpha(\text{Lit})$. The abstract variable look-up rule is analogous to the concrete one, except that looking up an undefined identifier returns $\alpha(\text{NULL})$ instead of NULL . Depending on the choice of VAL^\sharp and the definition of α , that could be a specific element NULL^\sharp ,

⁶ Our abstract semantics does not use boolean filter functions, because it is not practical to define realistic ones for a programming language as complicated as PHP. This choice has the downside of sacrificing some precision in the semantics of while loops, because we do not refine the information in the abstract state at the end of the loop to reflect that the guard has to be false.

or \top , or a different abstract element.

$$\llbracket \$x \rrbracket^\# \xi = \begin{cases} \xi(x) & \text{if } x \in \text{dom}(\xi) \\ \alpha(\text{NULL}) & \text{otherwise} \end{cases}$$

The abstract evaluation of arithmetic expressions is analogous to the concrete case

$$\llbracket e1 \text{ @ } e2 \rrbracket^\# \xi = \text{toNum}^\#(\llbracket e1 \rrbracket^\# \xi) \text{ @ }^\# \text{toNum}^\#(\llbracket e2 \rrbracket^\# \xi)$$

where $\text{@}^\#$ is the abstract operation corresponding to @ , and $\text{toNum}^\#$ is the abstract type juggling function corresponding to toNum . Both $\text{@}^\#$ and $\text{toNum}^\#$ are to be defined along with the abstract domain on which they depend. The abstract semantics of the other expressions follows a similar pattern.

3.2 Soundness of the analysis

We argue that the class of analyses defined by our abstract semantics is sound, assuming that the abstract domain has the right structure, and that the abstract operations provided with such domain satisfy some local soundness conditions.

Assumption 3.1 (Abstract Domain) *The abstract domain $\langle \text{VAL}^\#, \sqsubseteq \rangle$ is a complete lattice, and it forms a Galois connection $2^{\text{VAL}} \xleftrightarrow[\alpha]{\gamma} \text{VAL}^\#$ with the concrete domain $\langle 2^{\text{VAL}}, \subseteq \rangle$.*

Assumption 3.2 (Abstract Operations) *The abstract operations provided with the domain $\langle \text{VAL}^\#, \sqsubseteq \rangle$ are monotonic and locally sound approximations of the concrete ones: $\forall f^\#. \forall u, v \in \text{VAL}^\# : u \sqsubseteq v \Rightarrow f^\#(u) \sqsubseteq f^\#(v)$ and $\forall f, f^\#. \forall v \in \text{VAL} : \alpha(f(v)) \sqsubseteq f^\#(\alpha(v))$.*

We can take advantage of the big-step style of our semantics, and of our interest in properties of the final state, to bypass the standard definition of a collecting semantics and state our soundness theorem directly in terms of the concrete and abstract semantics. We only need to lift the definition of α from values to states: $\alpha(\sigma) = \alpha \circ \sigma$, and similarly for γ, \sqsubseteq .

Theorem 3.3 (Soundness) *The abstract semantics is a sound approximation of the concrete semantics: $\forall s \in \text{Stmt} : \alpha \circ \llbracket s \rrbracket \sqsubseteq \llbracket s \rrbracket^\# \circ \alpha$.*

Proof By induction on the derivation of $\llbracket \cdot \rrbracket$ (joining the definition for statements and expressions), using Assumption 3.1, Assumption 3.2 and standard properties of lattices. We show the case for if-else which is representative of the other cases. Assume that $\text{toBool}(\llbracket e \rrbracket) \sigma = \text{true}$ (the case when $\text{toBool}(\llbracket e \rrbracket) = \text{false}$ is analogous). Let $\xi = \alpha(\sigma)$. By inductive hypothesis, $\alpha(\llbracket b11 \rrbracket \sigma) \sqsubseteq \llbracket b11 \rrbracket^\# \xi$. By definition, $\llbracket b11 \rrbracket^\# \xi \sqsubseteq \llbracket b11 \rrbracket^\# \xi \sqcup \llbracket b12 \rrbracket^\# \xi$. Hence, we only need to exclude the case where $\gamma(\text{toBool}^\#(\llbracket e \rrbracket^\# \xi)) = \{\text{false}\}$. By Assumption 3.2, $\alpha(\text{toBool}(\llbracket e \rrbracket) \sigma) \sqsubseteq \text{toBool}^\#(\alpha(\llbracket e \rrbracket) \sigma)$. By inductive hypothesis, $\alpha(\llbracket e \rrbracket) \sigma \sqsubseteq \llbracket e \rrbracket^\# \xi$. By monotonicity of $\text{toBool}^\#$, $\text{toBool}^\#(\alpha(\llbracket e \rrbracket) \sigma) \sqsubseteq \text{toBool}^\#(\llbracket e \rrbracket^\# \xi)$. By transitivity of \sqsubseteq , $\alpha(\text{toBool}(\llbracket e \rrbracket) \sigma) \sqsubseteq \text{toBool}^\#(\llbracket e \rrbracket^\# \xi)$. By assumption, $\text{toBool}(\llbracket e \rrbracket) \sigma = \text{true}$. If $\gamma(\text{toBool}^\#(\llbracket e \rrbracket^\# \xi)) = \{\text{false}\}$, substituting in the equations above, we obtain

$\gamma(\alpha(\mathbf{true})) \subseteq \{\mathbf{false}\}$. By Assumption 3.1, $\{\mathbf{true}\} \subseteq \gamma(\alpha(\mathbf{true}))$, which leads to the contradiction $\{\mathbf{true}\} \subseteq \{\mathbf{false}\}$. \square

Proposition 3.4 (Incompleteness) *The abstract semantics is not complete:*

$$\exists s. \alpha \circ \llbracket s \rrbracket \not\subseteq \llbracket s \rrbracket^\# \circ \alpha.$$

Proof We show that there is a counterexample even for the most precise abstract domain possible: $\langle 2^{\text{VAL}}, \subseteq \rangle$ itself, where α and γ are the identity function. Let P be the μPHIP program $\$x=1; \text{ while } (\$x>0)\{ \$x=\$x-1; \}$. For any $\sigma \in \text{STATE}$, we have

$$(\alpha \circ \llbracket P \rrbracket)(\sigma) = \sigma[x \leftarrow \{0\}] \not\subseteq \sigma[x \leftarrow \{0, 1\}] = (\llbracket P \rrbracket^\# \circ \alpha)(\sigma). \quad \square$$

The informal meaning of our formal results is that if our analysis finds that a certain property holds, then that property (or possibly a stronger one) also holds across all the concrete executions compatible with the initial abstract state.

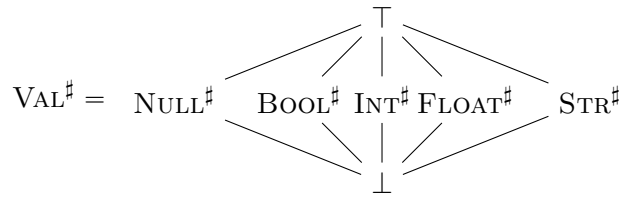
4 Abstract Domains for Type Juggling

Equipped with the abstract semantics of Section 3, we can design abstract domains and operations that capture the subtlety of type juggling in μPHIP . Rather than giving the definitions upfront, we expound the rationale behind our design, stressing expressivity, modularity and hopefully highlighting subtle points that can be useful to design domains for other languages with similar features.

4.1 Abstract Domains

We face three main design choices: how to combine the abstraction of the various types of μPHIP ; how to abstract each type; how to ensure that we can represent as much of the information relevant to type juggling as possible.

Type combination. Let us assume that for each set of basic values T we have defined an abstract type lattice $T^\#$. A typical analysis for statically-typed languages may combine abstract types using the *coalesced sum* lattice, which in our case yields



This choice is not appropriate for a dynamically-typed language such as μPHIP , as the resulting lattice cannot represent union types. For example, in the lattice above it must be the case that $\alpha(5) \sqcup \alpha(3.2) = \top$, leading to an unnecessary loss of precision when we convert such value to a string, because it has to be the case that $\text{toStr}^\#(\top) = \top$. In contrast, in a domain with the union type $\text{INT}^\# + \text{FLOAT}^\#$, $\text{toStr}^\#$ could have retained the information that numbers are never converted to empty strings, allowing to derive $\text{toStr}^\#(\text{INT}^\# + \text{FLOAT}^\#) = \text{STR}^\#_{\neq ""}$, assuming that the type abstraction of strings was able to account for such elements. A common

solution to this problem consists in switching to the *cartesian product* lattice of the abstract types

$$\text{VAL}^\sharp = \text{NULL}^\sharp \times \text{BOOL}^\sharp \times \text{INT}^\sharp \times \text{FLOAT}^\sharp \times \text{STR}^\sharp$$

where, for example, the union type $\text{INT}^\sharp + \text{FLOAT}^\sharp$ is implicitly represented by the vector $(\perp, \perp, \text{INT}^\sharp, \text{FLOAT}^\sharp, \perp)$.⁷

Type abstraction. Another key design choice is how to abstract the types themselves. For example, consider the μPHP semantics of division. It normally returns a `NUM` except for the case of division-by-zero, where it returns `false`. Abstracting a value directly to its type, as in $\alpha(5) = \text{INT}^\sharp$, is too imprecise because it prevents an analysis from detecting the division-by-zero case, and it forces the return type to be at best $\text{NUM}^\sharp + \text{BOOL}^\sharp$, instead of the more precise $\text{NUM}^\sharp + \text{false}^\sharp$. Hence, we include also the constants of each type to the product lattice. We define NULL^\sharp as the *lift*, and BOOL^\sharp and STR^\sharp as the *flat* lattices built from the corresponding sets:

$$\text{NULL}^\sharp = \text{lift}(\text{NULL}) \quad \text{BOOL}^\sharp = \text{flat}(\text{BOOL}) \quad \text{STR}^\sharp = \text{flat}(\text{STR})$$

By a judicious definition of INT^\sharp as the product lattice of signs, the constant 0, and natural numbers, we obtain the discriminating power of the traditional sign domain, plus the precision of numeric constants.

$$\text{INT}^\sharp = \text{flat}(\{+, -\}) \times \text{lift}(\{0\}) \times \text{flat}(\mathbb{N})$$

For example, (\top, \perp, \top) denotes non-zero integers, and $(+, 0, \top)$ represents non-negative integers.⁸ A similar argument applies to FLOAT^\sharp , which we define as

$$\text{FLOAT}^\sharp = \begin{array}{c} \top \\ / \quad \backslash \\ + \quad - \\ \backslash \quad / \\ \perp \end{array} \times \begin{array}{c} 0 \\ | \\ \perp \end{array} \times \begin{array}{c} \top \\ / \quad \backslash \\ 1 \quad 2 \quad 3 \quad \dots \\ \backslash \quad / \\ \perp \end{array} \times \begin{array}{c} 0 \\ | \\ \perp \end{array} \times \begin{array}{c} \top \\ / \quad \backslash \\ 001 \quad 12 \quad 266 \quad \dots \\ \backslash \quad / \\ \perp \end{array} \times \begin{array}{c} 0.0 \\ | \\ \perp \end{array}$$

and is isomorphic to $\text{INT}^\sharp \times \text{lift}(\{0\}) \times \text{flat}(\text{FRAC}) \times \text{lift}(\{0.0\})$, where FRAC is the set of non-zero “fractional parts” denoted by the regular expression $[0..9]^*[1..9]$. The last component of the product is necessary to distinguish $\alpha(0.0) \sqcup \alpha(1.2)$, which can be zero, from $\alpha(0.1) \sqcup \alpha(1.0)$, which cannot. For notational convenience, we denote the abstraction n^\sharp within the type domain T^\sharp by $\alpha_{T^\sharp}(n)$. We also abbreviate the bottom element of a product type, such as $(\perp, \perp, \perp) : \text{INT}^\sharp$ simply by \perp (and similar for \top). Finally, we use the shorthand $\top_{\text{BOOL}^\sharp}$ for the element $(\perp, \top, \perp, \perp, \perp) : \text{VAL}^\sharp$, with the obvious generalisation to other elements or domains.

Type juggling. Thanks to the definitions above, most of our domains already include enough information to handle type juggling. For example, the definition

⁷ The definition of α and γ will be left implicit as it can be understood from the context, as the obvious best approximation.

⁸ Some points in our lattice, such as $(+, 0, \perp)$ are redundant (zero has no sign). It is possible to optimise the domains to remove such points, slightly increasing the efficiency of the analysis (although the precision remains the same). We leave investigating that direction to future work.

of `toBool` depends on the set $\text{NUM}_0 = \{0, 0.0\}$. In order to define a precise abstract `toBool`[#], we should avoid loss of precision when deciding if an abstract value, once concretised, belongs to NUM_0 . Our domain achieves that, because for example $\gamma(\alpha(0) \sqcup \alpha(0.0)) = \text{NUM}_0$, and similarly $\gamma(\alpha(5) \sqcup \alpha(-3.2)) = \overline{\text{NUM}_0}$. The only domain which we need to refine explicitly is that of strings. In fact, `toBool` also relies on the set $\text{STR}_{\text{false}} = \{\text{"", "0"}\}$, but if $\text{STR}^{\#}$ is just the flat string domain, then $\gamma(\alpha(\text{STR}_{\text{false}})) = \text{STR} \neq \text{STR}_{\text{false}}$. A solution to this specific problem is to add to $\text{STR}^{\#}$ elements representing exactly $\alpha(\text{STR}_{\text{false}})$ and $\alpha(\overline{\text{STR}_{\text{false}}})$. The downside is that repeating this process for the other operations leads to a proliferation of special cases. For example, the division operation needs to decide if the result of `toNum` is in NUM_0 . Hence, for a precise `toNum`[#] we need two new points in $\text{STR}^{\#}$ representing precisely $\alpha(\{\text{"0"}, \text{"0.0"}\})$ and its complement. Moreover, we would need to introduce additional structure in the lattice to compare these points and the ones representing $\alpha(\text{STR}_{\text{false}})$, $\alpha(\overline{\text{STR}_{\text{false}}})$, and so on. Our proposal is instead to simply add all the information that is missing from the $\text{STR}^{\#}$ domain by adding to strings additional properties reflecting their value *after* an hypothetical type juggling. We re-define $\text{STR}^{\#}$ as a product involving also booleans, integers and floats, interpreted as properties of the corresponding abstract string:

$$\text{STR}^{\#} = \text{flat}(\text{STR}) \times \text{BOOL}^{\#} \times \text{INT}^{\#} \times \text{FLOAT}^{\#}$$

All the points representing properties of interest hypothesised above now are included in the lattice, with the correct ordering relation. For example, the string type of $\alpha(\text{STR}_{\text{false}})$ is $(\top, \text{false}, 0^{\#}, \perp)$, whereas the one of $\alpha(\text{"0"}, \text{"0.0"})$ is $(\top, \top, 0^{\#}, 0.0^{\#})$. As a final example of the expressivity of our type juggling domain $\text{VAL}^{\#}$, let x be the abstract value $\alpha(\text{"0.0doh"}) \sqcup \alpha(42)$, which in our domain is $(\perp, \perp, 42^{\#}, \perp, (\text{"0.0doh"}, \text{true}, \perp, 0.0^{\#}))$. Our domain contains enough information to be able to infer that x is not `NULL`, that it is `true` if converted to a boolean, and that the abstract evaluation of `84/x` yields $(\perp, \text{false}^{\#}, 2^{\#}, \perp, \perp)$, assuming a suitable definition of $\boxed{/}^{\#}$ (see Section 4.2).

4.2 Abstract Operations

We now discuss how to implement abstract operations that take advantage of the information represented by $\text{VAL}^{\#}$.

Type juggling functions. We focus on the example of `toNum`[#] as it illustrates all the main issues at hand. Since an abstract value is actually a 5-tuple of individual abstract types, in order to retain precision, we convert each component independently, using specialised functions such as $\text{StrToNum}^{\#} : \text{STR}^{\#} \rightarrow \text{VAL}^{\#}$, where the result is either an abstract number or \perp . Hence, the type of `toNum`[#] is $\text{VAL}^{\#} \rightarrow (\text{VAL}^{\#})^5$. Note that we do not collapse the resulting 5-tuple into a single $\text{VAL}^{\#}$ so that the operation that invoked the type juggling operation can leverage the information at best. For example,

$$\text{toNum}^{\#}((\perp, \perp, 4^{\#}, \perp, \text{"6doh"}^{\#})) = (\perp_{\text{VAL}^{\#}}, \perp_{\text{VAL}^{\#}}, 4_{\text{INT}^{\#}}, \perp_{\text{VAL}^{\#}}, 6_{\text{INT}^{\#}})$$

and a division by $2_{\text{INT}^{\#}}$ can return “positive integer” instead of “positive number”. The specialised conversions $\text{StrToNum}^{\#}$, $\text{BoolToNum}^{\#}$, etc. are straightforward to

define, following their concrete counterparts. For example, the latter returns respectively $\perp_{\text{VAL}^\#}, 0_{\text{INT}^\#}, 1_{\text{INT}^\#}, (+, 0, 1)$ on the inputs $\perp, \text{true}^\#, \text{false}^\#, \top$. Without loss of precision, we define $\text{StrToNum}^\#$ as the function $\lambda x. \pi_3(x) \sqcup \pi_4(x)$ that joins the pre-computed conversions to integer and float associated to the $\text{STR}^\#$ value.

Semantic operations. We now discuss how abstract operations can leverage the expressiveness of our domain. We give the example of division, which is representative of the other cases. Since $\text{toNum}^\#$ has already been applied by the abstract semantics of expressions, we now have to divide two 5-tuples of $\text{VAL}^\#$, hence $\square^\# : \text{VAL}^\#{}^5 \times \text{VAL}^\#{}^5 \longrightarrow \text{VAL}^\#$.

The first step of $\square^\#$ is to *normalise* each tuple by removing any $\perp_{\text{VAL}^\#}$ value, and retaining only its numeric components greater than \perp , obtaining two vectors of at most 6 elements each. For example, let $v = \text{toNum}^\#(\alpha(\text{true}) \sqcup \alpha("\text{-5foo}") \sqcup \alpha("\text{4.2doh}")) = (\perp_{\text{VAL}^\#}, 1_{\text{INT}^\#}, \perp_{\text{VAL}^\#}, \perp_{\text{VAL}^\#}, \alpha(-5) \sqcup \alpha(4.2))$. By normalising, we obtain $\mathfrak{n}(v) = [1^\#, -5^\#, 4.2^\#]$.

Once we have two normalised (row) vectors z and w , we can compute the analogous of the matrix product $z^t \times 1/w$, effectively obtaining a matrix r of dimension $|z| \times |w|$ where $r_{i,j} = z[i] /^\# w[j]$, and $/^\# : (\text{INT}^\# + \text{FLOAT}^\#)^2 \longrightarrow \text{VAL}^\#$ is the abstract division operator defined below

$$n_1 /^\# n_2 = \begin{cases} \alpha(m_1 \square m_2) & \text{if } \gamma(n_1) = \{m_1\} \text{ and } \gamma(n_2) = \{m_2\} \\ n_1 /^\#_{\text{INT}^\#} n_2 & \text{else, if } n_1, n_2 \text{ are both INT}^\# \\ \text{toFloat}^\#(n_1) /^\#_{\text{FLOAT}^\#} \text{toFloat}^\#(n_2) & \text{otherwise} \end{cases}$$

where $\text{toFloat}^\# : \text{INT}^\# \rightarrow \text{FLOAT}^\#$ maps $\alpha_{\text{INT}^\#}(k)$ to $\alpha_{\text{FLOAT}^\#}(k.0)$. The final step of $\square^\#$ is to join all the elements of r into a single $\text{VAL}^\#$. We define

$$u \square^\# v = \bigsqcup_{\substack{i \in 1..|x| \\ j \in 1..|y|}} x[i] /^\# y[j] \quad \text{where } x = \mathfrak{n}(u) \text{ and } y = \mathfrak{n}(v).$$

The abstract division operator $/^\#$ relies on specialised abstract divisions for integers and floats (respectively $/^\#_{\text{INT}^\#}$ and $/^\#_{\text{FLOAT}^\#}$). When both operands are abstract integers, we perform a further normalisation, separating the information about 0 from the information about \mathbb{N} encoded in each operand. For example, $\mathfrak{n}((- , 0, 5)) = [(-, \perp, 5), (\perp, 0, \perp)]$. Then, we compute $u /^\#_{\text{INT}^\#} v = \bigsqcup_{\substack{i \in 1..|x| \\ j \in 1..|y|}} x[i] /^\#_{\text{INT}^\#} y[j]$ where $x = \mathfrak{n}(u)$ and $y = \mathfrak{n}(v)$, and the inner $/^\#_{\text{INT}^\#}$ is computed using the rules in Figure 2. The case for $/^\#_{\text{FLOAT}^\#}$ is analogous.

For example, let us revisit the example of $84/x$ from Section 4.1, where this time

$$x = \alpha("\text{0.0doh}") \sqcup \alpha("\text{1argh}") \sqcup \alpha(42) \sqcup \alpha(0) = (\perp, \perp, (+, 0, 42), \perp, (\top, \text{true}^\#, 1^\#, 0.0^\#))$$

We have that $\mathfrak{n}(\text{toNum}^\#(x)) = [1^\#, 0.0^\#, (+, 0, 42)]$. The first two divisions are computed directly as $\alpha(84 \square 1) = \alpha(84)$ and $\alpha(84 \square 0.0) = \alpha(\text{false})$. The third division is computed as $84 /^\#_{\text{INT}^\#} (+, 0, 42)$. The denominator is normalised to $[0^\#, 42^\#]$, leading to two further divisions $\alpha(84 \square 0) = \alpha(\text{false})$ and $\alpha(84 \square 42) = \alpha(2)$. Hence, the final result is $\alpha(84) \sqcup \alpha(\text{false}) \sqcup \alpha(2) = (\perp, \text{false}^\#, (+, \perp, \top), \top, \top)$, where we

$/^{\sharp}_{\text{INT}^{\sharp}}$	0^{\sharp}	1^{\sharp}	n_2^{\sharp}	(\top, \perp, n_2)
0^{\sharp}	false[#]	0^{\sharp}	0^{\sharp}	0^{\sharp}
1^{\sharp}	false[#]	1^{\sharp}	$\alpha(1 \sqcap n_2)$	$\alpha(1 \sqcap n_2) \sqcup \alpha(1 \sqcap -n_2)$
n_1^{\sharp}	false[#]	n_1^{\sharp}	$\alpha(n_1 \sqcap n_2)$	$\alpha(n_1 \sqcap n_2) \sqcup \alpha(n_1 \sqcap -n_2)$
(\top, \perp, n_1)	false[#]	(\top, \perp, n_1)	$\alpha(-n_1 \sqcap n_2) \sqcup \alpha(n_1 \sqcap n_2)$	$\alpha(n_1 \sqcap n_2) \sqcup \alpha(-n_1 \sqcap n_2) \sqcup \alpha(n_1 \sqcap -n_2) \sqcup \alpha(-n_1 \sqcap -n_2)$
$(+, \perp, \top)$	false[#]	$(+, \perp, \top)_{\text{INT}^{\sharp}}$	$(\pi_1(n_2^{\sharp}), \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\pi_1(n_2^{\sharp}), \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
$(-, \perp, \top)$	false[#]	$(-, \perp, \top)_{\text{INT}^{\sharp}}$	$(\pi_1(n_2^{\sharp}), \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\pi_1(n_2^{\sharp}), \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
(\top, \perp, \top)	false[#]	$(\top, \perp, \top)_{\text{INT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
$(\top, 0, \top)$	false[#]	$(\top, 0, \top)_{\text{INT}^{\sharp}}$	$(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$

$/^{\sharp}_{\text{INT}^{\sharp}}$	$(+, \perp, \top)$	$(-, \perp, \top)$	(\top, \perp, \top)	$(\top, 0, \top)$
0^{\sharp}	0^{\sharp}	0^{\sharp}	0^{\sharp}	false[#] \sqcup 0^{\sharp}
1^{\sharp}	$(+, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (+, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(-, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (-, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(+, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (+, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
n_1^{\sharp}	$(\pi_1(n_1^{\sharp}), \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\pi_1(n_1^{\sharp}), \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\pi_1(n_1^{\sharp}), \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\pi_1(n_1^{\sharp}), \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
(\top, \perp, n_1)	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
$(+, \perp, \top)$	$(+, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (+, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(-, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (-, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
$(-, \perp, \top)$	$(-, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (-, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(+, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (+, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
(\top, \perp, \top)	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, \perp, \top)_{\text{INT}^{\sharp}} \sqcup (\top, \perp, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$
$(\top, 0, \top)$	$(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	$(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$	false[#] \sqcup $(\top, 0, \top)_{\text{INT}^{\sharp}} \sqcup (\top, 0, \top, 0, \top, \perp)_{\text{FLOAT}^{\sharp}}$

 Figure 2. Tables for the abstract operation $/^{\sharp}_{\text{INT}^{\sharp}}$.

know that, unless there was a division by zero, we obtain a positive integer. Note that both normalisation steps introduced by \sqcap^{\sharp} and $/^{\sharp}_{\text{INT}^{\sharp}}$ were essential to retain this level of precision.

5 Conclusions

We have defined the formal semantics of $\mu\text{P}\text{HP}$, a subset of PHP that precisely represents type juggling behaviour, as a basis to explore new and expressive abstract domains for type/value analysis. We have also defined an abstract interpreter that implements, parametrically on the domain, a non-relational, path-sensitive analysis to leverage our abstract domains. We have shown with various examples that our value analysis is more expressive than comparable ones present in the literature. To the best of our knowledge, a novelty of our approach is the definition of the string domain as the product of the string type with other abstract types (integers and floats). This construction helps retaining more precise information about strings

after type juggling.

The main limitations of our current work also suggest natural directions for future work. μPHP covers only a small subset of PHP, and it will be interesting to see how our type juggling domain interacts with the analyses of other challenging language features such as aliasing, functions, objects and exceptions. Our construction of the type juggling domain strives to be systematic but we do not investigate how an analysis of completeness of the abstract operations, along the lines of [7], may lend further justification to our current design choices, or lead to the completely automated construction of a more precise domain.

Acknowledgments. This work is partially supported by EPSRC grant EP/K032089/1.

References

- [1] The PHP Group. PHP Zend Engine. <http://php.net>. Accessed: 2016-06-09.
- [2] P. Cousot. Types as abstract interpretations. In *POPL'97*, 1997.
- [3] 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*, 1977.
- [4] P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *POPL'79*, 1979.
- [5] J. Dahse and T. Holz. Simulation of built-in PHP features for precise static code analysis. In *NDSS'14*, 2014.
- [6] D. Filaretti and S. Maffeis. An executable formal semantics of PHP. In *ECOOP'14*, 2014.
- [7] R. Giacobazzi, F. Ranzato, and F. Scozzari. Making abstract interpretations complete. *J. ACM*, 2000.
- [8] D. Hauzar and J. Kofron. Framework for static analysis of PHP applications. In *ECOOP'15*, 2015.
- [9] S. H. Jensen, A. Møller, and P. Thiemann. Type analysis for javascript. In *SAS'09*, 2009.
- [10] N. Jovanovic, C. Krügel, and E. Kirda. Pixy: A static analysis tool for detecting web application vulnerabilities (short paper). In *(SE'06)*, 2006.
- [11] V. Kashyap, K. Dewey, E. A. Kuefner, J. Wagner, K. Gibbons, J. Sarracino, B. Wiedermann, and B. Hardekopf. JSAI: a static analysis platform for javascript. In *FSE'14*, 2014.
- [12] E. Kneuss, P. Suter, and V. Kuncak. Phantm: PHP analyzer for type mismatch. In *FSE'10*, 2010.