

Type inference for PHP

Using annotations to provide more precise results

Ruud van der Weijde

August 17, 2014, 30 pages

Supervisor: Jurgén Vinju
Host organisation: Werkspot, <http://www.werkspot.nl>
Host supervisor: Winfred Peereboom



WERKSPOT
KLUS • KLIK • KLAAR
ruud.van.der.weijde@werkspot.nl
<http://www.werkspot.nl>



UNIVERSITEIT VAN AMSTERDAM
FACULTEIT DER NATUURWETENSCHAPPEN, WISKUNDE EN INFORMATICA
MASTER SOFTWARE ENGINEERING
<http://www.software-engineering-amsterdam.nl>

Contents

Abstract	2
Preface	3
1 Introduction	4
1.1 PHP	4
1.2 Position	4
1.3 Contribution	5
1.4 Plan	5
2 Background and related work	6
2.1 PHP Language Constructs	6
2.2 Rascal	8
2.3 M3	8
2.4 Related work	8
3 Research Method	10
3.1 Introduction	10
3.2 Research Question	10
3.3 Types	10
3.3.1 PHP types	11
3.3.2 Subtypes	11
3.4 Fact extraction	12
3.4.1 Type extraction	13
3.4.2 Constraint extraction	13
3.5 Annotations	25
3.6 Constraint solving	25
4 Results	26
4.1 Annotations	26
5 Case Study	27
6 Conclusion	28
6.1 Conclusion	28
6.2 Future work	28
Glossary	29

Abstract

PHP is widely used.

Some parts are hard to analyse.

Knowing (object) types can provide better results (no proof for this).

Annotations can (not) help to better resolve types.

Preface

In this section I will thank everyone who has helped me. Maybe also introduce some anecdote on how this research came to be.

Chapter 1

Introduction

1.1 PHP

PHP¹ is a server-side scripting language created by Rasmus Lerdorf in 1995. The original name ‘Personal Home Page’ changed to ‘PHP: Hypertext Preprocessor’ in 1998. PHP source files are executed using the PHP Interpreter. The language is dynamically typed, which means that the behaviour of the source code will be examined during run-time. Statically typed languages would apply these modification during compile type. PHP supports duck-typing, which means that the type of an expression can be transformed to another type at a certain point during execution.

Evolution The programming language PHP evolved after its creation in 1995. In the year 2000 Object-Oriented (OO) language structures were added to the langue with the release of PHP 4.0. The 5th version of PHP was release in 2004 and provided an improved OO structure. Namespace were added in PHP 5.3 in 2009, to be able to resolve class naming conflicts between library and create better readable class names. Namespaces are comparable to packages in JAVA. The most recent stable version is 5.5 in which the OPcache extension is added. OPcache speeds up the performance of including files on run-time by storing precompiled script bytecode in shared memory.

Popularity According to the Tiobe Index² of July 2014, PHP is the 7th most popular programming language. The language has been in the top 10 since its introduction in the Tiobe Index in 2001. More than 80 percent of the websites have a php backend³. The majority of these websites use PHP version 5, rather than version 4 or older versions. It is therefor wise to focus on PHP version from 5 and discard the older unsupported versions.

Analysability Although the popularity for more than a decade, there is still a lack of good PHP code analysis tools. Tools can help to reveal security vulnerabilities or find vulnerabilities or bugs in source code. The tools can also provide code completions or do automatic transformations which can be used to execute refactoring patterns. Source code analysis can be performed statically or dynamically or a combination of the two. More information on the analysability of php can be found in section 2.1.

1.2 Position

Todo: explain here why it is important to analyse programs: for security analysis and compiler optimisations. The next stuff needs to be rewritten.

As far as we know, there is no constraint based type inference research like this one performed for PHP. That makes this research unique. There have been similar analysis for other dynamic languages, like smalltalk, ruby and javascript.

¹<http://php.net>

²<http://www.tiobe.com/index.php/content/paperinfo/tpci/index.html>, July 2014

³<http://w3techs.com/technologies/details/pl-php/all/all>, July 2014

1.3 Contribution

TODO Review this part when the result of the analysis are performed.

- Created an M3 for php.
- Constraint system

Some idea's are that this analysis can help IDE tools to perform transformations on the source code. (But the performance may not be sufficient.)

The creation of the M3 model can help to compare researchers compare PHP programs with other programming languages. For now only Java is implemented, but more can follow (unchecked statement).

1.4 Plan

The rest of the thesis is as follows: chapter 2 contains background and context information about related work. Here we will explain the php language and explain similar research. In the next chapter 3 the research method is explained, which will explain the steps taken in this the research.

Chapter 2

Background and related work

This chapter explains a few language constructs which are important for this research. Further more Rascal and M3 will be explained. In the last section of this chapter the related work is described.

2.1 PHP Language Constructs

In this section for this research important language constructs are presented. Explanations of these constructs should help to understand the performed analysis.

Scoping PHP has a few scopes. "Define what scope is". Global, namespace, class, method, function. The global scope is contained in every file which is not inside a function or class. The global scope can contain namespaces. Namespaces are comparable to packages in Java. When namespaces are used, classes and functions will be scoped to the namespace. You can access them by providing the namespace name.

Todo: say something about the global statement and \$GLOBALS. (it is resolved in the M3 relation @uses). Refer to this link¹.

- In general:
- Functions are declared in

Includes Note to myself: how will I deal with this concept in my analysis (totally ignore it??? when maybe not add it to this background information. In our research we will assume that all files are included.). The problem of including files can be reduced using namespaces and autoloading. When a class which is not loaded in memory is instantiated, the autoloading will try to include a file and load the class. For this analysis we will include all files

* Refer to the analysis of mark hills, that most files can be resolved, but not all. We consider the use of including scripts for logic as bad practice. Every file should contain a class, and in this case, it is for our analysis not very interesting to resolve the includes.

Conditional functions and classes Explain the code below.

```
1 if (!class_exists("Foo"))
2     class Foo { /* ... */ }
3
4 if (!function_exists("bar"))
5     function bar() { /* ... */ }
```

Listing 2.1: Conditional class and function definitions

Explain the code below.

¹<http://php.net/manual/en/language.variables.scope.php>, July 2014

```

1 function f() {
2     function g() {
3         class C {}
4     }
5 }
6
7 g(); f(); // will fail because 'g();' is not declared yet
8 f(); g(); // will work because 'g();' is declared when calling 'f();'
9 f(); new C(); // will fail because 'g();' needs to be called first
10 f(); g(); new C(); // will work because 'g();' is called and has declared 'f();'

```

Listing 2.2: Conditional function declaration

Dynamic features These include dynamic variables, dynamic class instantiations, dynamic function calls.

Late static binding Late static binding² is implemented in PHP since version 5.3 by adding the keyword 'static' to the language. It has the same function as 'parent' and 'self', because they both point to a class. The main difference is that 'parent' and 'self' can be resolved statically. 'static' can only be resolved on runtime and represents the exact class that is instantiated.

Magic methods In PHP it is allowed to call methods or use properties that do not exist. Normally this would result in a fatal error, but not with the use of magic methods. One of the magic methods is the constructor method '__construct'.

Dynamic class properties Although it is a good practice to define your class properties, it is not required. On runtime it is possible to add properties to classes, even without the implementation of magic methods.

```

1 class C {}
2 $p = (new C())->nonExistingProperty;
3 var_dump($p); // NULL
4 $c = new C();
5 $p = $c->nonExistingProperty = "property now exists";
6 var_dump($p); // string(19) "property now exists"

```

Listing 2.3: Dynamic class property

Annotations Explain here how annotations work in php. (in the next section I will explain something about parsing/reading annotations for the analysis)

Other concepts Todo: move this list to a different chapter, this is about my research. Not about php in general.

- We assume that the program is correct. This means that warnings can happen, but fatal errors not.
- We assume that all files are included.
- We assume that register globals is off! (maybe add some other runtime environments)
- Ignore warnings (because most production code has them off)

These items will not be covered by the analysis (maybe add this to threats/future work)

- Analysis is flow insensitive

²<http://php.net/manual/en/language.oop5.late-static-bindings.php>, July 2014

- Closure
- References
- Variable constructs (variable -variable, -method/function calls, -class instantiation, eval) :: todo: explain WHY not.
- Yields

2.2 Rascal

Explain what [Rascal](#) is and what we use it for... Some general stuff here. This research is performed in a Rascal environment.

2.3 M3

In this chapter, only describe what m3 is. In some next section I will explain the additional php m3 stuff.

But for now (need to create an image of the steps here):

- Parse all files (create [ASTs](#))
- For each file, create an m3 of the ast of the parsed file (aka extract information)
- Combine all the m3s
- Now run more analysis when all facts are collected in the m3.
- Done. M3 is finished!

*Maybe add this section to the next chapter, for the research methods. Many aspects will be useful for the type inference analysis

The M3-model is a generic model which can be used to analyse software programs. Our goal is to provide the results in an M3 model. Future research can use this to compare different programming languages.

The following core items are filled for M3:

- Containment
- Declarations (need to explain in more details)
- Modifiers
- Extends
- Uses (explain in details what is and what is NOT covered)

The following php specific items are added:

- Extends (class or interface and their extended class or interface)
- Implements (which class implements which interfaces)
- TraitUses (which class uses which trait)
- Parameters (methods and functions and their parameters)
- Constructors (which class uses which constructor, explain this more)
- Aliases (class aliases, for example the usage of class_alias)
- Annotations (constraints annotations on classes, methods, fields and variables)

Todo, explain in more details...

2.4 Related work

Describe these:

- ‘The HipHop Compiler for PHP’[\[Zha+12\]](#) (not much information available, only source code)
- ‘Phantom: PHP analyzer for type mismatch’[\[KSK10a; KSK10b\]](#) (investigate this in more details)
- PHPLint ³(uses a different kind of annotations, not the java like phpdocs)

- ‘Soft typing and analyses on PHP programs’[], code implementations: <https://github.com/henkerik/typing> and <https://github.com/marcelosousa/soft-typing-PHP5> (created for php4, code for php5, should check this out, might be able to compare results with this)
- ‘Design and Implementation of an Ahead-of-Time Compiler for PHP’[Big10] (to check in detail)

Also describe their differences with my research.

³<http://www.icosaedro.it/phplint>, july 2014

Chapter 3

Research Method

3.1 Introduction

<< Todo: properly introduction of this chapter >>

The type inference analysis will be performed in this order:

- Resolve types
- Resolve sub type relations
- Extract facts from the source code
- Extract constraints from source code
- Solve the constraints
- Add result to m3 (all decls)

In the next step, annotations will be added to see how the results be like.

3.2 Research Question

The research question will be something like:

Will the use of phpdoc annotations produce better (to be defined) results for constraint based type inference?

3.3 Types

Explain here what I mean with a type...

```
module lang::php::m3::TypeSymbol

data TypeSymbol
  = \any()      | array(TypeSymbol arrayType)
  | boolean()   | class(loc decl)
  | float()     | integer()
  | null()      | object()
  | resource()  | string()
;
```

3.3.1 PHP types

PHP has a similar class inheritance structure and interface implementation as Java. The main difference is that in PHP all class are **public** and that inner classes are not allowed in PHP.

The basis types in PHP are integers, floats (similar to doubles and reals), booleans, strings, arrays, resources and null. When variables are initialised without a values, they are null. The recourse type is a special one which is not important for this research.

3.3.2 Subtypes

Explain something about subtypes here. For now, only this figure 3.1

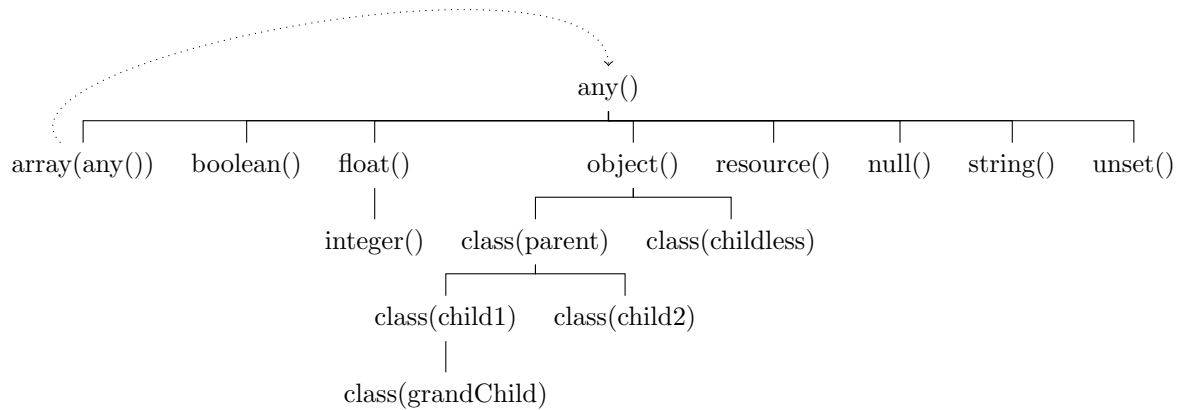


Figure 3.1: Subtype hierarchy

The subtype relation of class inheritance is a **reflexive transitive closure** relation. A class extension of class A on class C will define class A as a subtype of class C in our analysis, as you can see in figure 3.2. If a class does not extend another class, it will implicitly extend the **stdClass** class. You can see that this happens with class D in the example. The **stdClass** is represented as the type **object()** in our analysis.

The basic PHP types also contain a subtype relation. Integers are subtypes of floats.

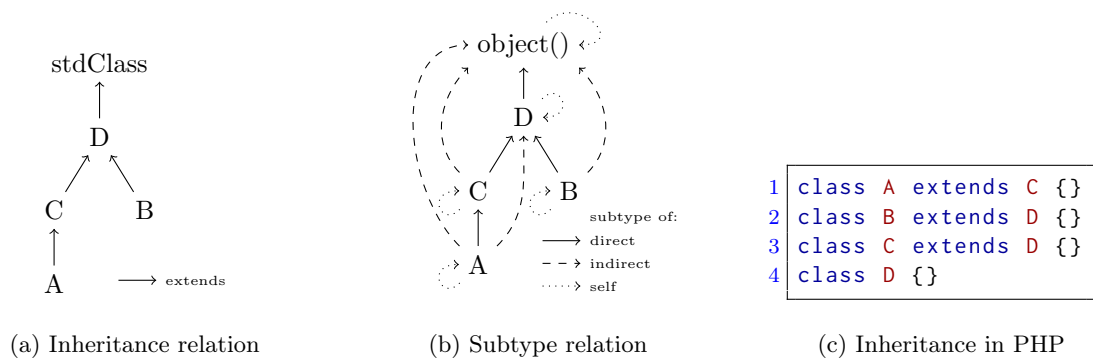


Figure 3.2: Relation of subtypes among classes

3.4 Fact extraction

We can extract fact about classes, class-constants/fields/methods, functions, parameters. For these facts, we can use a relation, so we have a many-to-many relation. On the left side we will have the class, function or method. On the right side we have their attribute.

A list of properties: (todo: rewrite this list into a 'normal' section.)

- `<loc classDecl, className(str name)>`
- `<loc classDecl, classMethod(str name, set[Modifier] modifiers)>`
- `<loc classDecl, classProperty(str name, set[Modifier] modifiers)>`
- `<loc classDecl, classConstant(str name, set[Modifier] modifiers)>`
- `<loc classDecl, classConstructorParameters(list[PhpParam] params)>`
- `<loc methodDecl, methodName(str name)>`
- `<loc methodDecl, methodParameters(list[PhpParam] params)>`
- `<loc functionDecl, functionName(str name)>`
- `<loc functionDecl, functionParameters(list[PhpParam] params)>`

In Rascal:

```
alias TypeFacts = rel[loc decl, Fact fact];
```

```
data Fact
  = className(str name) // = FQN = fully qualified name
  | classMethod(str name)
  | classProperty(str name)
  | classConstant(str name)
  | classConstructorParameters(list[PhpParam] params)
  | methodName(str name)
  | methodParameters(list[PhpParam] params)
  | functionName(str name)
  | functionParameters(list[PhpParam] params)
  ;
```

Other facts that will be used:

```
alias PhpParams = lrel[loc decl, set[loc] typeHints, bool isRequired, bool byRef];
data Annotation = returnType(set[TypeSymbol]) | parameterType(loc var, set[TypeSymbol])
  | varType(loc var, set[TypeSymbol]);

anno rel[loc from, loc to] M3@containment; // 'from' directly contains 'to'
anno rel[loc from, loc to] M3@extends; // 'from' extends 'to'
anno rel[loc from, loc to] M3@implements; // 'from' implements 'to'
anno rel[loc decl, PhpParams params] M3@parameters; // formal parameters of functions/methods
anno rel[loc decl, loc to] M3@constructors; // 'decl' and its constructor 'to'
anno rel[loc decl, Annotation annotation] M3@annotations; // result of parsed php docs
```

3.4.1 Type extraction

In order to define the subtype relations in class extensions, we will need to declare all existing class types. We can do this in rascal like is done in the example below:

```
visit (system) {  
  case c:class(_, _, _, _, _): types += class(c@decl);  
}
```

Once all types are defined, we can add the subtype relation. We will need to have the subtype of `int()` and `float()` and the class extensions. You can see that in the code below:

```
public rel[TypeSymbol, TypeSymbol] getSubTypes(M3 m3, System system)  
{  
  rel[TypeSymbol, TypeSymbol] subtypes  
    // add int() as subtype of float()  
    = { <\int(), float(> }  
    // use the extends relation from M3  
    + { <class(c), class(e)> | <c,e> <- m3@extends }  
    // add subtype of object for all classes which do not extends a class  
    + { <class(c@decl), object(> | 1 <- system, /c:class(n,_,noName(),_,_) <- system[1]  
};  
  
  // compute reflexive transitive closure and return the result  
  return subtypes*;  
}
```

3.4.2 Constraint extraction

Introduction is needed here... for now I will just list the types that I have found. Maybe this needs to be moved to a different chapter.

This is a list of items which are not supported (yet):

- References (in PHP they are symbol table aliases)
 - on expression assignments :: `$a = &$b`
 - on functions :: `function &f() {...}`
 - on parameters :: `function f(&$a) {...}`
- Variable structures:
 - ~~Variable variables~~ :: `$$a;`
 - ~~Variable class instantiation~~ :: `new $a;`
 - ~~Variable method or function calls~~ :: `$a();`
- List assign :: `list($a,$b) = array("one","two");` (we can assume that the rhs is of type array, when the program is correct)
- ~~Method or function parameters (including type hints)~~
- ~~Class structures, method calls~~
- ~~Class Constants~~
- ~~The global statement~~ (should be resolved by the usage relation from M3)
- ~~Casts of expressions~~

- Parameters
- ~~Predefined variables~~ (\$this, self, parent, static)
- ~~Eval~~ (will not be supported)
- ~~Closures~~ (not used much in production code)
- ~~Traits~~ (not used much in production code)
- Callable (introduced in 5.4 as typehint, not used much in production code)
- Foreach(\$a as ... (=> ...)) => \$a is an array or an object;
- ~~return;=> return type is null~~ (is added to the situation when there are no return statements)
- add predefined globals (and their type: \$[GLOBALS, _SERVER, _GET, _POST, _REQUEST, _COOKIE, _ENV, _SESSION, php_errormsg] (all in global scope))
- add magic constants: __[DIR, FILE, LINE,_NAMESPACE, FUNCTION, CLASS, METHOD]__
- ~~predefined constants: TRUE(b), FALSE(b), NAN(f), INF(f), NULL(n), STDIN(r), STDOUT(r), STDERR(r)~~
- define("name", value) mixed with constants (?out of scope?)
- ~~keywords: self, parent, static in a class~~ (is included in method and property calls)
- ADD CONSTANTS! RECORD THE TYPE OF THE DEFINED CONSTANTS AND TRY TO READ THEIR TYPE.

Legend

=	=	Equal to (type)	C	=	A class
<:	=	Is subTypeOf	$\rightarrow c$	=	A class constant
E_k	=	An expression	$\rightarrow p$	=	A class property
$[E_k]$	=	Type of some expression	$\rightarrow m$	=	A class method
f	=	A function	$[m]$	=	(Return) type of a method call
$[f]$	=	(Return) type of a function	(A_n)	=	The n'th actual argument
$:: c$	=	Static property fetch	(P_n)	=	The n'th formal parameter
$:: m$	=	Static method call	th	=	Type hint
$:: p$	=	Static property fetch	v	=	Default value
Mfs	=	Modifiers	Γ	=	Whole program

Table 3.1: Constraint legend

A list of predefined items can be found here:

- [Constants](#)
- [Variables](#)
- todo: functions
- todo: classes

Expressions

Normal assignment

$$\frac{E \equiv (E_1 = E_2)}{\begin{array}{l} [E_2] <: [E_1], \\ [E_1] <: [E] \end{array}}$$

```
1 $a = $b; // [$b] <: [$a], [$a] <: [$a=$b]
2 $c = $d = $e; // [$e] <: [$d], [$d] <: [$c],
3           // [$d] <: [$d=$e], [$c] <: [$c=$d=$e]
```

Listing 3.1: Normal assignment

Ternary

$$\frac{E \equiv (E_1 ? E_2 : E_3)}{[E] <: [E_2] \vee [E] <: [E_3]}$$
$$\frac{E \equiv (E_1 ? : E_3)}{[E] <: [E_1] \vee [E] <: [E_3]}$$

```
1 $expr ? $b : $c; // typeOf(E) is subtypeOf($b) or subtypeOf($c)
2 $expr ? : $c; // typeOf(E) is subtypeOf($expr) or subtypeOf($c)
```

Listing 3.2: Ternary

Assignments with operators (1) always resulting in ints

$$\begin{array}{l} E_1 \&= E_2 \\ E_1 |= E_2 \\ E_1 \wedge= E_2 \\ E_1 <<= E_2 \\ E_1 >>= E_2 \\ E_1 \% = E_2 \end{array}$$
$$\frac{}{[E_1] = integer()}$$

```
1 $a &= $b; /* $a = integer() */
2 $a |= $b; /* $a = integer() */
3 $a ^= $b; /* $a = integer() */
4 $a <<= $b; /* $a = integer() */
5 $a >>= $b; /* $a = integer() */
6 $a %= $b; /* $a = integer() */
```

Listing 3.3: Assignments with operators (1)

Assignments with operators (2) string concat (.=)

$$\frac{E_1 .= E_2}{[E_1] = string(),}$$
$$if([E_2] <: object()) => hasMethod([E_2], "__toString")$$

```
1 $a .= $b; /* $a = string() */
2 // An error occurs when $b is of type object() and
3 // __toString is not defined or does not return a string
```

Listing 3.4: Assignments with operators (2)

Assignments with operators (3) resulting in int where rhs is no array

$$\frac{E_1 / = E_2 \quad E_1 - = E_2}{[E_1] = integer(), \quad [E_2] \neq array(any())}$$

```

1 $a /= $b; /* $a = integer() */
2 $a -= $b; /* $a = integer() */
3 // An error occurs when $b is of type array() for /= and -=
4 // Fatal error: Unsupported operand types

```

Listing 3.5: Assignments with operators (3)

Assignments with operators (4) resulting in int or float

$$\frac{E_1 * = E_2 \quad E_1 + = E_2}{[E_1] <: float()}$$

```

1 $a *= $b; /* when $b == (boolean()|integer()|null()) */ /* $a = integer() */
2 $a *= $b; /* when $b != (boolean()|integer()|null()) */ /* $a = float() */
3 $a += $b; /* when $b == (boolean()|integer()|null()) */ /* $a = integer() */
4 $a += $b; /* when $b != (boolean()|integer()|null()) */ /* $a = float() */

```

Listing 3.6: Assignments with operators (4)

Unary operators

$$\frac{E \equiv (+E_1) \vee (-E_1)}{[E] <: float(), \quad [E_1] \neq array(\backslash any())}$$

$$\frac{E \equiv (!E_1)}{[E] = boolean()}$$

$$\frac{E \equiv (\sim E_1)}{[E_1] = float() \vee [E_1] = integer() \vee [E_1] = string(), \quad [E] = integer() \vee [E] = string()}$$

$$\frac{E \equiv (E_1 + +) \vee (E_1 - -)}{if([E_1] <: array(\backslash any())) => [E] <: array(\backslash any()), \quad if([E_1] = boolean()) => [E] = boolean(), \quad if([E_1] = float()) => [E] = float(), \quad if([E_1] = integer()) => [E] = integer(), \quad if([E_1] = null()) => [E] = integer() \vee [E] = null(), \quad if([E_1] <: object()) => [E] <: object(), \quad if([E_1] = resource()) => [E] = resource(), \quad if([E_1] = string()) => [E] = integer() \vee [E] = float() \vee [E] = string()}$$

$$\begin{array}{c}
E \equiv (+ + E_1) \vee (- - E_1) \\
\hline
if([E_1] <: array(\text{any}()) => [E] <: array(\text{any}()), \\
if([E_1] = boolean()) => [E] = boolean(), \\
if([E_1] = float()) => [E] = float(), \\
if([E_1] = integer()) => [E] = integer(), \\
if([E_1] = null()) => [E] = null, \\
if([E_1] <: object()) => [E] <: object(), \\
if([E_1] = resource()) => [E] = resource(), \\
if([E_1] = string()) => [E] = integer() \vee [E] = float() \vee [E] = string()
\end{array}$$

```

1 +$a // positive
2 -$a // negation
3 !$a // not
4 ~$a // bitwise not
5 $a++ // post increase
6 $a-- // post decrease
7 ++$a // pre increase
8 --$a // pre decrease

```

Listing 3.7: Unary operators

Binary operators

$$\begin{array}{c}
E \equiv (E_1 + E_2) \\
\hline
[E] <: array(_) \vee [E] <: float(), \\
if([E_1] <: array(_) \wedge [E_2] <: array(_) => [E] <: array(_), \\
if([E_1]! <: array(_) \vee [E_2]! <: array(_) => [E] <: float() \\
\\
E \equiv (E_1 - E_2) \vee (E_1 * E_2) \vee (E_1 / E_2) \\
\hline
[E] <: float(), \\
[E_1] \neq array(_), \\
[E_2] \neq array(_) \\
\\
E \equiv (E_1 \% E_2) \vee (E_1 << E_2) \vee (E_1 >> E_2) \\
\hline
[E] = integer() \\
\\
E \equiv (E_1 \& E_2) \vee (E_1 | E_2) \vee (E_1 \wedge E_2) \\
\hline
[E] = string() \vee [E] = integer(), \\
if([E_1] = string() \wedge [E_2] = string()) => [E] = string(), \\
if([E_1] \neq string() \vee [E_2] \neq string()) => [E] = integer()
\end{array}$$

```

1 $a + $b // addition
2 $a - $b // subtraction
3 $a * $b // mulitiplication
4 $a / $b // division
5 $a \% $b // modulus
6 $a \& $b // bitwise And
7 $a | $b // bitwise Or
8 $a ^ $b // bitwise Xor
9 $a << $b // bitwise shift left
10 $a >> $b // bitwise shift right

```

Listing 3.8: Binary operators

Comparison operators

$$\begin{aligned} E &\equiv (E_1 == E_2) \\ E &\equiv (E_1 === E_2) \\ E &\equiv (E_1 != E_2) \\ E &\equiv (E_1 <> E_2) \\ E &\equiv (E_1 !== E_2) \\ E &\equiv (E_1 < E_2) \\ E &\equiv (E_1 > E_2) \\ E &\equiv (E_1 <= E_2) \\ E &\equiv (E_1 >= E_2) \\ \hline [E] &= \text{boolean}() \end{aligned}$$

```
1 $a == $b /* boolean() */
2 $a === $b /* boolean() */
3 $a != $b /* boolean() */
4 $a <> $b /* boolean() */
5 $a !== $b /* boolean() */
6 $a < $b /* boolean() */
7 $a > $b /* boolean() */
8 $a <= $b /* boolean() */
9 $a >= $b /* boolean() */
```

Listing 3.9: Comparison operators

Logical operators

$$\begin{aligned} E &\equiv (E_1 \text{ and } E_2) \\ E &\equiv (E_1 \text{ or } E_2) \\ E &\equiv (E_1 \text{ xor } E_2) \\ E &\equiv (E_1 \&\& E_2) \\ E &\equiv (E_1 || E_2) \\ \hline [E] &= \text{boolean}() \end{aligned}$$

```
1 $a and $b /* boolean() */
2 $a or $b /* boolean() */
3 $a xor $b /* boolean() */
4 $a && $b /* boolean() */
5 $a || $b /* boolean() */
```

Listing 3.10: Logical operators

Array

Array value fetch

$$\begin{aligned} E &\equiv E_1[E_2] \\ \hline [E_1] &\neq \text{object}(), \\ \text{if}([E_1] = \text{string}()) &=> [E] = \text{string}(), \\ \text{if}([E_1] = \text{array}(\{\text{types}\})) &=> [E] <: \{\text{types}\}, \\ \text{if}([E_1] \neq \text{string}() \wedge [E_1] \neq \text{array}(_)) &=> [E] = \text{null}() \end{aligned}$$

```

1 $a[0];
2 // typeof($a) != object()
3 // when typeof($a) == string() => typeof($a[/.*...*/]) is string()
4 // when typeof($a) == array() => typeof($a[/.*...*/]) is mixed()
5 // when typeof($a) != string|array => typeof($a[/.*...*/]) is null()

```

Listing 3.11: Array value fetch

Array declaration

$$\frac{E', \text{ where } E' \text{ is an array declaration}}{[E'] = array(any())}$$

```

1 array(/.*...*/); // typeof() = array();
2 // Rascal: array(_) => array()

```

Listing 3.12: Array declaration

Scalars

Scalars

$$\frac{E, E \text{ is a string}}{[E] = string()}$$

$$\frac{E, E \text{ is a float}}{[E] = float()}$$

$$\frac{E, E \text{ is an integer}}{[E] = integer()}$$

```

1 "Str" // string()
2 'abc' // string()
3 100 // integer()
4 1.4 // float()

```

Listing 3.13: Scalars

Encapsulated strings

$$\frac{E, E \text{ is an encapsd string}^*}{[E] = string()}$$

* When a string contains expression(/variables), it is processed as encapsd.

```

1 "$var"

```

Listing 3.14: Encapsulated strings

Casts

Casts

Note: PHP Warnings are ignored

$$\frac{(array)E_1}{[(cast)E_1] <: array(any())}$$
$$\frac{(bool)E_1 \vee (boolean)E_1}{[(cast)E_1] = boolean()}$$
$$\frac{(float)E_1 \vee (double)E_1 \vee (real)E_1}{[(cast)E_1] = float()}$$
$$\frac{(int)E_1 \vee (integer)E_1}{[(cast)E_1] = integer()}$$
$$\frac{(object)E_1}{[(cast)E_1] <: object()}$$
$$\frac{(string)E_1}{[(cast)E_1] = string(),}$$
$$if([E_1] <: object()) => hasMethod([E_1], "__toString")$$
$$\frac{(unset)E_1}{[(cast)E_1] = null()}$$

```
1 (array)$a // <: array(\any())
2 (bool)$a // boolean()
3 (float)$a // float()
4 (int)$a // integer()
5 (object)$a // object()
6 (string)$a // string(), when $a == object() the object needs to have __toString()
7 (unset)$a // null()
```

Listing 3.15: Casts

Clone

Clone

$$\frac{E \equiv clone(E_1)}{[E] <: object(), [E_1] <: object()}$$

```
1 clone($a) // typeof($a) = object, typeof(clone($a)) = object
```

Listing 3.16: Clone

Class

Class instantiation (1) matching the class name

$$\frac{new\ C_1() \in \Gamma, C.hasName(C_1.name)}{[new\ C_1] = class(C.decl)}$$

```
1 new C;
```

Listing 3.17: Class instantiation (1)

Class instantiation (2) of an expression

$$\frac{new\ E_1 \in \Gamma}{[new\ E_1] <: object()}$$

```
1 $c = "C";
2 new $c;
```

Listing 3.18: Class instantiation (2)

Special keywords `self` `parent` `static`

$$\frac{\$this \in C}{[\$this] = class(C)}$$

$$\frac{self \in C}{[self] = class(C)}$$

$$\frac{parent \in C}{[parent] <: class(C.p^{+1})}$$

$$\frac{static \in C}{[static] <: class(C^{+2})}$$

```
1 // $this can only be used within a class
2 $this // in class C -> class(C)
3 self  // in class C -> class(C)
4 parent // in class C -> parentOf(class(C))
5 static // in class C -> class(C) or parentOf(class(C))
```

Listing 3.19: Special keywords

Class property fetch

* Possible add fact that the field `E` is declared in class `C`, when it is on the left side of an assignment.

$$\frac{\$this \rightarrow E_1 \subseteq C_1}{[E_1] = C_1.hasProperty(E_1.name, static \notin Mfs) \vee [E_1] = C_1.parent.hasProperty(E_1.name, public|protected \in Mfs \wedge static \notin Mfs) \vee [E_1] = C_1[parent].hasMethod("__get")}$$

$$\frac{self :: E_1 \subseteq C_1}{[E_1] = C_1.hasProperty(E_1.name, static \in Mfs)}$$

$$\frac{parent :: E_1 \subseteq m}{[E_1] = C.parent.hasProperty(E_1.name, static \in Mfs)}$$

¹all direct or indirect parents

²all direct or indirect parents, including the class itself

$$\frac{E_1 \rightarrow E_2 \subseteq C_1^*}{[E_1] = C_1.\text{hasProperty}(E_2.\text{name}, \text{static} \notin \text{Mfs}) \vee [E_1] = C_1.\text{parent}.\text{hasProperty}(E_2.\text{name}, \text{public}|\text{protected} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs}) \vee [E_1] = C.\text{hasProperty}(E_2.\text{name}, \text{public} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs})}$$

*The same goes for static property fetches, except for the ‘static \notin Mfs’ part: ‘static \in Mfs’.

$$\frac{E_1 \rightarrow E_2 \not\subseteq C \subseteq \Gamma^*}{[E_1] = C.\text{hasProperty}(E_2.\text{name}, \text{public} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs})}$$

*Property fetch outside a class scope, also for static properties.

```
1 $this->prop // name = prop, vis = public|protected, !static || mm
2 self::$prop // static property in class
3 parent::$prop // static property in the parent(s)
4 $a->prop // non-static property fetch
5 $a::$pro // static property fetch
```

Listing 3.20: Class property fetch

Class property fetch variable

$$\frac{E_1 \rightarrow E_2, E_2 \text{ is an expression}}{[E_1] = \text{object}()}$$

```
1 $b = "b";
2 $a->$b
```

Listing 3.21: Class property fetch variable

Class method call

$$\frac{\$this \rightarrow E_1 \subseteq C_1}{[E_1] = C_1.\text{hasMethod}(E_1.\text{name}, \text{static} \notin \text{Mfs}) \vee [E_1] = C_1.\text{parent}.\text{hasMethod}(E_1.\text{name}, \text{public}|\text{protected} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs}) \vee [E_1] = C_1[\text{parent}].\text{hasMethod}(\text{"__call"})}$$

$$\frac{self :: E_1 \subseteq C_1}{[E_1] = C_1.\text{hasMethod}(E_1.\text{name}, \text{static} \in \text{Mfs}) \vee [E_1] = C_1.\text{parent}.\text{hasMethod}(E_1.\text{name}, \text{public}|\text{protected} \in \text{Mfs} \wedge \text{static} \in \text{Mfs}) \vee [E_1] = C_1.\text{hasMethod}(\text{"__callStatic"})}$$

$$\frac{parent :: E_1 \subseteq C_1}{[E_1] = C_1.\text{parent}.\text{hasMethod}(E_1.\text{name}, \text{public}|\text{protected} \in \text{Mfs}) \vee [E_1] = C_1.\text{parent}.\text{hasMethod}(\text{"__callStatic"})}$$

$$\frac{E_1 \rightarrow E_2 \subseteq C_1^*}{[E_1] = C_1.\text{hasMethod}(E_2.\text{name}, \text{static} \notin \text{Mfs}) \vee [E_1] = C_1.\text{parent}.\text{hasMethod}(E_2.\text{name}, \text{public}|\text{protected} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs}) \vee [E_1] = C.\text{hasMethod}(E_2.\text{name}, \text{public} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs})}$$

*The same goes for static method calls, except for the ‘static \notin Mfs’ part: ‘static \in Mfs’.

$$\frac{E_1 \rightarrow E_2 \not\subseteq C \subseteq \Gamma^*}{[E_1] = C.\text{hasMethod}(E_2.\text{name}, \text{public} \in \text{Mfs} \wedge \text{static} \notin \text{Mfs})}$$

*method call outside a class scope, also for static methods.

```
1 $this->methodCall();
2 self::methodCall();
3 parent::methodCall();
4 $a->methodCall();
5 $a::methodCall();
```

Listing 3.22: Class method call

Class method call variable

$$\frac{E_1 \rightarrow E_2(), E_2 \text{ is an expression}}{[E_1] <: \text{object}()}$$

```
1 $a->$methodCall()
```

Listing 3.23: Class method call variable

Class constants (needs to be reviewed)

$$\frac{\text{self}::c_1 \subseteq \Gamma}{[\text{self}::c_1] = C_1.\text{hasConstant}(E_2.\text{name}) \vee [\text{self}::c_1] = C_1.\text{parent}.\text{hasConstant}(E_2.\text{name}, \text{public|protected} \in \text{Mfs})}$$

$$\frac{\text{parent}::c_1 \subseteq \Gamma}{[\text{self}::c_1] = C_1.\text{parent}.\text{hasConstant}(E_2.\text{name}, \text{public|protected} \in \text{Mfs})}$$

$$\frac{E_1::c_1 \subseteq \Gamma}{[E_1] = \text{object}()}$$

```
1 self::CONST
2 parent::CONST
3 SOMECLASS::CONST
```

Listing 3.24: Class constants (needs to be reviewed)

Parameters

Parameters in class instantiation

*These parameters are just examples for what happens if they have typeHints (*th*), default values (*v*) or none *The constructor can be found in the m3 model (@constructors(loc classDecl, loc constructorMethodDecl))

$$\frac{\begin{array}{l} \text{new } C_1 (A_1, A_2, \dots, A_k) \subseteq \Gamma \\ \$a \rightarrow m() (A_1, A_2, \dots, A_k) \subseteq \Gamma \\ \text{function}_1 (A_1, A_2, \dots, A_k) \subseteq \Gamma \end{array}}{[P_1] <: [A_1], [A_1] <: [th_1], [P_1] <: [th_1], \text{hasRequiredParam}(P_1), \text{hasRequiredParam}(P_k)}$$

$$\frac{\begin{array}{l} \text{class } C (th_1 P_1, P_2 = v, \dots, P_k) \subseteq \Gamma \\ \text{public function } m() (th_1 P_1, P_2 = v, \dots, P_k) \subseteq \Gamma \\ \text{function } (th_1 P_1, P_2 = v, \dots, P_k) \subseteq \Gamma \end{array}}{[P_1] <: [A_1], [A_1] <: [th_1], [P_1] <: [th_1], \text{hasRequiredParam}(P_1), \text{hasRequiredParam}(P_k)}$$

```
1 new C($foo);
```

Listing 3.25: Parameters in class instantiation

Scope

Type of a certain variable within some scope

this applies to global- class- function- and method- scope

$$\frac{E, E', E'', E''' \dots \text{etc} \subseteq f \quad E \text{ is a variable}}{[E] = [E] \vee [E'] \vee [E''] \vee [E'''] \dots \text{etc}}$$

```
1 function f() {
2   $a = 1;
3   $a = "true";
4 }
5 // typeof($a) is typeof($a1, $a2, ..., $an);
```

Listing 3.26: Type of a certain variable within some scope

Return type of function or method (1) having no return statements or return;

$$\frac{\text{return} \not\subseteq f \vee \text{return}; \subseteq f}{[f] = \text{null}()}$$

```
1 function f() {} // no return = null()
2 function f() { return; } // return; = null()
```

Listing 3.27: Return type of function or method (1)

Return type of function or method (2) every exit path ends with a return statement

$$\frac{(\text{return } E_1) \vee (\text{return } E_2) \vee \dots \vee (\text{return } E_k) \subseteq f}{[f] <: [E_1] \vee [E_2] \vee \dots \vee [E_k]}$$

```
1 function f() {
2   if (rand(0,1))
3     return $a;
4   else
5     return $b;
6 }
7 // returns typeOf($a) or typeOf($b)
```

Listing 3.28: Return type of function or method (2)

Return type of function or method (3) possible no return value

$$\frac{(\text{return } E_1) \vee (\text{return } E_2) \vee \dots \vee (\text{return } E_k) \vee (\neg \text{return}) \subseteq f}{[f] <: [E_1] \vee [E_2] \vee \dots \vee [E_k] \vee \text{null}()}$$

```
1 function f() {
2   if (rand(0,1))
3     return $a;
4   else if (rand(0,1))
5     return $b;
6 }
7 // returns typeOf($a) or typeOf($b) or null()
```

Listing 3.29: Return type of function or method (3)

Function calls

Function call

$$\frac{f() \subseteq \Gamma}{[f()] <: \text{return of } [f]}$$

```
1 function f() {}
2 f();
```

Listing 3.30: Function call

Function call variable

$$\frac{E \equiv E_1() \subseteq \Gamma}{[E] = \text{any}(),}$$

$$[E_1] <: \text{object}() \vee [E_1] = \text{string}(),$$

$$\text{if}([E_1] <: \text{object}()) \Rightarrow \text{hasMethod}(\text{"__invoke"})$$

```
1 function f() {}
2 $f = "f";
3 $f(); // unknown what function will be called
4 // [$f] <: object(with __invoke method) | [$f] = string()
```

Listing 3.31: Function call variable

How to resolve expressions:

- Find all expressions which are defined above and annotate them with @type.
- Annotate the rest of the expressions with @type = any(); (should only be for relevant expressions)

3.5 Annotations

Explain how the annotations are added to the constraints.

3.6 Constraint solving

Explain how we will solve constraints.

Chapter 4

Results

For the results we picked X software products to see how it performs. For each product we performed the type inference with and without reading annotations from the doc blocks.

4.1 Annotations

The results of the analysis when adding the annotations to the analysis. Compare the results with the results of the analysis without the annotation information.

Chapter 5

Case Study

Explain how the case study is performed.

Chapter 6

Conclusion

Summary of the whole work, with conclusions. T.B.A.

6.1 Conclusion

6.2 Future work

Explain something about combining this analysis to other analysis (like dead code elimination, constant folding/propagation resolve, alias analysis, array analysis) to gain more precise results.

Something about performance optimisations... Explain what is already done to boost the performance and what still can be done.

Use a bigger corpus to gains better results of the analysis by doing analysis on more programs.

Glossary

AST

Abstract Syntax Tree, explain... .

Rascal

Rascal is a meta-programming language developed by SWAT (Software analyse and transformation) team at CWI in the Netherlands. See <http://www.rascal-mpl.org/> for more information.

reflexive transitive closure

A relation is transitive if $\langle a, b \rangle \in R$ then $\langle b, a \rangle \in R$.

A relation is reflexive if $\langle a, b \rangle \in R$ and $\langle b, c \rangle \in R$ then $\langle a, c \rangle \in R$.

A reflexive transitive closure can be established by creating direct paths for all indirect paths and adding self references, until a fixed point is reached.

stdClass

A predefined class in the PHP library. The class is the root of the class hierarchy. It is comparable to the Object class in Java.

Bibliography

- [Big10] Paul Biggar. “Design and Implementation of an Ahead-of-Time Compiler for PHP”. In: (2010).
- [KSK10a] Etienne Kneuss, Philippe Suter, and Viktor Kuncak. “Phantm: PHP Analyzer for Type Mismatch”. In: *Proceedings of the Eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering*. FSE ’10. Santa Fe, New Mexico, USA: ACM, 2010, pp. 373–374. ISBN: 978-1-60558-791-2. DOI: [10.1145/1882291.1882355](https://doi.org/10.1145/1882291.1882355). URL: <http://doi.acm.org/10.1145/1882291.1882355>.
- [KSK10b] Etienne Kneuss, Philippe Suter, and Viktor Kuncak. “Runtime Instrumentation for Precise Flow-Sensitive Type Analysis”. English. In: *Runtime Verification*. Ed. by Howard Barringer et al. Vol. 6418. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 2010, pp. 300–314. ISBN: 978-3-642-16611-2. DOI: [10.1007/978-3-642-16612-9_23](https://doi.org/10.1007/978-3-642-16612-9_23). URL: http://dx.doi.org/10.1007/978-3-642-16612-9_23.
- [Zha+12] Haiping Zhao et al. “The HipHop Compiler for PHP”. In: *SIGPLAN Not.* 47.10 (Oct. 2012), pp. 575–586. ISSN: 0362-1340. DOI: [10.1145/2398857.2384658](https://doi.org/10.1145/2398857.2384658). URL: <http://doi.acm.org/10.1145/2398857.2384658>.