

Type inference for PHP

A constraint based type inference written in Rascal

Ruud van der Weijde

January 8, 2017, 50 pages

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



WERKSPOT

HEERENGRACHT 496, AMSTERDAM
<http://www.werkspot.nl>



UNIVERSITEIT VAN AMSTERDAM

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

Contents

Abstract	3
Preface	4
1 Introduction	5
1.1 PHP	5
1.2 Position	5
1.3 Contribution	6
1.4 Plan	6
2 Background and Related Work	7
2.1 PHP Language Constructs	7
2.2 Rascal	9
2.3 Type systems	10
2.4 Related work	11
3 Research Context	12
3.1 Types	12
3.2 Type hierarchie	14
3.3 Research context	15
4 Design of PHP type inference	16
4.1 Type inference rules	16
4.1.1 Scalars	17
4.1.2 Assignments	18
4.1.3 Unary operators	20
4.1.4 Binary operators	23
4.1.5 Array	26
4.1.6 Casts	27
4.1.7 Clone	27
4.1.8 Class	28
4.1.9 Scope	29
4.2 Annotations	30
4.3 PHP built-ins	32
5 Implementation of PHP type inference	33
5.1 M^3 for PHP	33
5.1.1 Core elements	33
5.1.2 PHP specific elements	34
5.1.3 Constructing a model	35
5.2 Constraint extraction	36
5.2.1 Constraint definitions in Rascal	36
5.2.2 Constraint extraction in Rascal	37
5.3 Constraint solving	39
5.3.1 The algorithm	39
5.4 Additional constraints	39

6	Evaluation	41
6.1	Experiment setup	41
6.2	Steps to reproduce	42
6.3	Results	43
6.4	Analysis	44
7	Conclusion	46
7.1	Conclusion	46
7.2	Future work	46
7.3	Threats to validity	46
	Glossary	48

Abstract

Dynamic language like PHP are generally hard to statically analyse because of run-time dependencies. Despite the wide usage of PHP programs over the internet, there seem not many tools available to support all aspect completely. In statical analysis the program is analysed without running the program and there are many things undecided. In this thesis we present a constraint based type inference written in Rascal. Rascal is a programming language for meta-programming in the domain of software analysis and transformations. We created this type inference for PHP to be able to resolve the types of expressions used in programs. Follow up analysis can be performed when expressions are typable, for example to find vulnerabilities or provide programming context in IDEs. In a small experiment where we tested if adding type annotations of PHP docblock and PHP build-in information would help to infer more types. We saw that the number of inferred types increases when type hint annotations are taken into account.

Preface

This thesis could not have been completed without the help and support of various people. In this section I want to thank the people, who helped me during the process, by mentioning their name and shortly explain in which way they supported me. I start with a short story on how this research came to be.

The initial idea for this thesis was to dive into the topic of software analysis in order to find vulnerabilities in PHP programs. When expanding further on this topic, my interest in static analysis grew and I planned to reconstruct data flows in programs to allow taint analysis. In order to perform dataflow analysis for object oriented PHP programs, the types of expressions need to be known. This topic of resolving the run-time types at compile time was big enough to solely write a thesis on.

The First person I need to thank is Jurgen Vinju, who helped me throughout the whole process. I want to thank Jurgen for his endless enthusiasm, number of ideas and suggestions and personal help on coaching and mentoring me. Jurgen helped defining a research topic and gave me many directions where to go or look for whenever I got stuck on a subject.

Next I want to thank Mark Hills. The implementation of this research is build on top of the PHP Analysis in Rascal (PHP AiR) framework, created by Mark. This research uses the PHP AiR framework to parse PHP files to Rascal ASTs. We expanded the project with an implementation of an M^3 model for PHP programs and a constraint extractor and solver. I'm happy that I could contribute to the project in return.

I want to thank Bas Basten for the collaboration on improving and extending the M^3 model. After creating the initial version of the M^3 model for PHP, Bas helped to improved the model by adding more program information to the model. With his expertise on Rascal and and my knowledge of PHP we formed a solid team.

I would like to thank Winfred Peereboom for giving me the opportunity to write my thesis at Werkspot. During the period Winfred helped me on various aspects of coaching and mentoring me.

Finally I want to thank my girlfriend for her endless support and patience. And of course I want to thank everyone I forgot to mention here personally, but who did help me direct or indirect, conscious or unconscious.

Chapter 1

Introduction

1.1 PHP

PHP¹ is a server-side programming 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 Interpreter². The language is dynamically typed and allows objects to be changed during run-time. PHP uses duck-typing which means that there is no object type validation checking, but only validates if the attempted operation is permitted on the object.

Evolution The programming language PHP evolved since its creation in 1995. The first milestone was in the year 2000 when Object-Oriented (OO) language structures were added to the language with the release of PHP 4.0. The 5th version of PHP, released in 2004, provided an improved OO structure including the first type declarations for function parameters. Namespaces were added in PHP 5.3 in 2009, to resolve class naming conflicts between libraries and to create better readable class names. The OPcache extension is added was added in PHP 5.5 and speeds up the performance of including files on run-time by storing precompiled script byte-code in shared memory. The latest 5.x version is 5.6 and includes more internal performance optimisations and introduces a new debugger. The most recent stable version is 7, which is not taken into this research. In this latest version they achieved a mayor performance increase and memory decrease. Also type hints for scalar types are added, method/function return types, and strict typing can be enforced for the available type hints.

Popularity According to the Tiobe Index³ of July 2016, PHP is the 6th most popular language of all programming languages. 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 running php backend⁴. The majority of these websites use PHP version 5, rather than older or newer versions. We therefor focus our type inference on PHP version 5.

1.2 Position

Due to various dynamic features in PHP not all types and execution paths can be resolved without actually executing the program. Source code analysis tools need to know execution paths and types of expressions for optimal results in discovering security vulnerabilities or bugs. Such tools could also provide code completions or do automatic transformations when executing refactoring patterns. Source code analysis is performed statically or dynamically or a combination of the two. In static analysis the program is not executed.

¹<http://php.net>

²<https://github.com/php/php-src>

³http://www.tiobe.com/tiobe_index, July 2016

⁴<http://w3techs.com/technologies/details/pl-php/all/all>, July 2016

1.3 Contribution

This research contributes to the static analysis research field of dynamic programming languages by presenting a constraint based type inference analysis that over-approximates runtime values at compile time. Results of this analysis can be used to improve software analysis tools. In order to resolve the types of expressions we implemented a generic model, M^3 , which holds various program facts for PHP programs.

The main contributions of this thesis are:

- M^3 model for PHP programs
- constraint based type inference for PHP programs

An M^3 model for PHP programs contains various facts about the programs. The M^3 model, see section 5.1 for more information, was initially solely supporting Java programs. The addition of support for PHP programs was beneficial not only for this thesis, but also for other researchers. The model provides program context which is used during constrain extraction and solving. The model helped other researchers by providing program context when comparing the PHP programming usage with programs written in other languages.

The biggest contribution of this thesis is the *constrained based type inference for PHP programs*. In this type inference process we use an abstract syntax tree to generate type constraints on language constructs. We then solve these constraints to come to a set of types for each expression. We variated the inference process by adding context information of type annotations and php built-in to the constraint extraction to find out that this helps to resolve more types. The inferred types can help IDE tools and programmers by providing helpful tools, which could lead to better software development.

1.4 Plan

The rest of this thesis is organised as follows:

Chapter 2 contains background information and related work. The background information consists of important language constructs, information about annotations in PHP, brief introduction to Rascal, M^3 , and type systems. We end this chapter with related researches and their relation to this research.

Chapter 3, research context, describes the research approach and context. We explain under which assumptions this research is executed.

Chapter 4 describes the design of the type inference for PHP. We present the constraint rules on various language constructs. Next to that we give information about type annotations and php built-in information.

Chapter 5 contains implementation details. We show how we implemented the M^3 model for PHP. Next we explain how we implemented the constraint extracting. The constraint solving is explained by showing the used algorithm.

Chapter 6 shows type inference results on real world PHP programs. We present the results of multiple programs and analyse the results.

Chapter 7 ends this thesis with the conclusion, future work section and lists various threats to validity.

Chapter 2

Background and Related Work

Chapter 2 provides relevant background information. The first section, section 2.1, describes seven important language constructs which the reader needs to understand in order to understand the difficulties of analysing PHP. Section 2.2 explains [Rascal](#), the programming language used for the analysis. In this section we will explain M^3 , a programming language independent meta model which holds various facts about programs, in more details. Section 2.3 provides information about type systems and how type systems relate to this research. The last section of this chapter, section 2.4, presents related work and how these researches relate to this thesis.

2.1 PHP Language Constructs

PHP has various language constructs which complicate [static code analysis](#). This section presents language constructs and why these constructs are important for this research. Explanations of these constructs help you to understand the performed analysis. The discussed parts are scope, file includes, conditional classes and functions, dynamic features, late static binding, magic methods and dynamic class properties.

Scope In PHP, all classes and functions are globally accessible once they are declared. All classes and functions are implicitly public, inner classes are not allowed, and conditional functions (see upcoming paragraph about conditional classes and functions) will be available in the global scope. If a class or function is declared inside a namespace, their full qualified name includes the name of the namespace. Variables have three scope levels: global-, function-, and method-scope. Under normal circumstances when a variable is declared inside a function or method, their scope is limited to this function or method. Variables declared outside function or methods are available in the global scope, but not in the method or function scope. There is an exception for some predefined global variables which are available everywhere. Examples are `$GLOBALS`, `$_POST`, and `$_GET`. Variables inside a function or method can be aliased to a global variable by adding the keyword `GLOBAL` in front of the variable name. The variable are then linked to the global variable in the symbol table¹.

Script includes PHP allows files to include other PHP-files during program execution. The content of these files will be loaded at the place where the include statement is defined. This means that if you use an include in the middle of a file, the source code of this included file will be virtually inserted at that position.

File includes are mainly used for loading classes and for including templates to render output. PHP5 allows automated class loading based on the namespace, which is called autoloading classes. With autoloading classes there is no need for including files manually for each class.

Research by Mark Hills et al. [\[HKV14\]](#) has shown that most includes can be resolved with statical analysis. In this research we do not run such an analysis, we will assume that all files in the project are included during execution.

¹<http://php.net/manual/en/language.variables.scope.php>, July 2016

Conditional classes and functions Once a file is included in the execution, all the found classes and functions are declared in the top level scope. All class and function declarations within condition statements or within a method or function scope are only declared when the conditional statement is executed.

An example of an conditional statement can be found in listing 2.1. If the class `Foo` or function `bar` do not exist before the statements are executed, then the class and function will not yet be declared. When you try to use the class or function before the code is executed, the script will exit with an fatal error.

```
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

Another example of dynamic function and class loading is displayed in listing 2.2. If the first call is `g()` as you can see in line 8, the script will result in a fatal error because function `g()` will only be declared after function `f()` is executed. `class C` will be declared once function `g()` is executed. As soon as the functions and classes are declared, they are available in the top scope, possibly prefixed with the name of the namespace they are declared within.

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

Listing 2.2: Conditional function declaration

Dynamic features PHP comes with dynamic built-in features like: include dynamic variables, dynamic class instantiations, dynamic function calls, dynamic function creation, reflection, and eval. New functions and classes can be declared on the fly during run-time. Method calls, or even whole pieces of code, can be executed based on variable strings.

A previous study by Mark Hills[HKV13] has shown that most real world applications make use of dynamic features. Dynamic features are powerful, but can complicate the statical analysis. Analysis like constant propagation is needed to help resolving most of these dynamic features. This is not in scope for this research, but could be added in future work.

Late static binding Late static binding² is implemented in PHP since version 5.3 by adding the keyword `static` to the language. Its usage is similar to the keyword `self`, which refers to the current class. The main difference is that `self` refers to the class where the code is located, while `static` refers to the actual instantiated class and can be a descendant of the class. The keyword `self` can be statically resolved without running the program while `static` can only be resolved on runtime.

Magic methods PHP allows calls and property access on methods and fields that don't exist on a class. Normally a call to a non-existing method or property would result in a fatal error, but with the use of magic methods you can specify the wanted behaviour. Listing 2.3 shows an example of the `__call` method. This method is triggered when a inaccessible properties are read. In this example the code will try to return the value of a private property based on the provided name. The full list of magic

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

methods is `__construct`, `__destruct`, `__call`, `__callStatic`, `__get`, `__set`, `__isset`, `__unset`, `__wakeup`, `__toString`, `__invoke`, `__set_state`, `__clone`, and `__debugInfo`.

```
1 class Car {
2     private $maxSpeed = 210;
3     function __get($name) { return @$this->$name; }
4 }
5 var_dump((new Car)->maxSpeed); // 210
6 var_dump((new Car)->numberOfWheels); // NULL
```

Listing 2.3: Magic methods in PHP

Dynamic class properties Although it is a good practice to define your class properties, it is not required to do so in PHP. After instantiating a class it is possible to add properties to objects, even without magic method usage. In listing 2.4 you can see a code sample of adding a property to an object. The access of the non-existing property `nonExistingProperty` will result in a warning, but code execution will continue and will just return `NULL`. The code on line 4 creates the property by writing to a nonexisting property. The object of variable `$c` will have the `nonExistingProperty` publicly available now. But in a new class instantiation, as you can see on line 6, will not have the property.

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

Listing 2.4: Dynamic class property

2.2 Rascal

Rascal[KSV09] is a meta programming language developed by Centrum Wiskunde & Informatica (CWI). Rascal is designed to analyse, transform and visualise source code. The language is build on top of Java and implements various concepts of existing programming languages. In this research, Rascal is the main programming language. Rascal is used for gathering facts about the program and to solve constraints. The facts are gathered by visiting AST tree representing the program and hold semantic information about the program. Constraints are generated based on the collected facts and these constraints are solved with an in Rascal created constraint solver. The only part that does not use Rascal is the PHP parser. Although this could be implemented in Rascal, there was an existing library written in PHP available.

M^3 [Izm+13; Bas+15] is a model which holds various information of source code and is implemented in Rascal. This model is created to gain insights in the quality of open-source projects. For our research we use the M^3 -model to store facts about the program in a structured way, so we can easily use it at a later stage of the analysis.

The core elements of the M^3 -model are **containment**, **declarations**, **documentation**, **modifiers**, **names**, **types**, **uses**, **messages**. The **declarations** relation contains class, method, variable- information with their logical name and their real location. The type of the relation are **locations** and represent the logical name of the declaration and will be used in the rest of the M^3 . The **containment** relation has information on what declarations are contained in each other. For example a package can contain a class; a class can contain fields and methods or an inner class; a method can contain variables. The **documentation** relation contains all comments from the source code and its source location. The **modifiers** relation has information about the modifiers of declarations. Modifiers can be abstract, final, public, protected, or private. The **names** relation contains a simplified name of the full declarations. The **types** relation holds information about the type of the source code elements. The **uses** relation describes what references use an object. For instance when a field of a class is used in some expression, the **uses** relation links the

field in the expression to the declaration of the field in the class. And lastly, `messages` contains global error, warning, and info statements.

2.3 Type systems

A **type** is a set of possible values and a set of operations that can be performed on them. PHP is a dynamically typed language, which means that the types of the expressions are not examined at compile time. PHP implements duck typing, so the type of objects are not examined during run-time, but only checks if the operations are allowed on the data. In this thesis we are interested in the run-time types of values so we are able to perform static source code analysis.

Type systems Type systems define how a set of rules are applied to types in their context. A type system validates the type usage with **type checking**. In order to perform type checking the types of the expressions needs to be known. The process of resolving the types of the expression is called **type inference**. Both type checking and type inference are explained in more details below.

Type checking Type checking is a mechanism which validates and/or enforces the constraints of a type in their specific context. There is a difference between static type checking and dynamic type checking. **Static type checking** is a process of checking the types based on the source code. The static type checker will ensure that a program is type safe before executing the program, which means that there will occur no type errors during runtime. **Dynamic type checking** performs the type checking during runtime. This means that the program needs to run to gain feedback on the usage of types. PHP is a dynamically typed language, which means that there are no types checked before actually running the program. Although PHP also implements some static type features, like parameter type hints, it cannot be perfectly determine all types at compile time.

Type inference Type inference is the process of resolving types of variables and expressions. The inference process is a prerequisite to perform type checking. Being able to infer the type before running the program enables you to optimise code execution by applying compiler optimisations. These optimisation can be performance improvements or memory optimisations. In dynamic languages like PHP, it can be hard to resolve the type of a variable or expression without running the program. In statically typed languages, type inference happens at compile time. In the next paragraph we will briefly explain some type inference systems.

The **Hindley-Milner**[\[Hin69\]](#) (HM) type system was found in 1969 by Roger J. Hindley and almost 10 years later rediscovered[\[Mil78\]](#) by Robin Milner. The first implementation was created four years later by PhD student Luis Damas. Damas proved the soundness and completeness of the HM type system with **Algorithm W**[\[DM82\]](#) in the context of the programming language ML. The HM type system deduces the types of the variables to their most abstract type, based on their usage. Type declarations and hints are not necessarily to perform type inference. The type system is used for various functional languages. Haskell for example uses the Hinley-Milner type system as a foundation for the Haskell type system.

Control Flow Analysis[\[NNH99\]](#) (CFA) is concerned with resolving sound approximate run-time values at compile time. CFA is build on top of data flow analysis[\[ASU86\]](#) and tries to resolve the control-flow problem for high order programming languages. The control-flow problem deals with resolving which caller can call which callee in a program. One of the earlier CFA algorithms was Shivers' 0CFA algorithm[\[Shi88\]](#), a flow-sensitive constraint based algorithm. Shiver then defined *k*-CFA[\[Shi91\]](#), where the precision of the analysis is increased by taking the context of the expressions into account. The *k*-CFA algorithms compute an conservative over-approximations of run-time values during compile type.

The **Cartesian Product Algorithm**[\[Age95\]](#) (CPA) is a type inference algorithm created by Ole Agesen in 1995. Agesens work was based on Palsberg and Schwartzbach' **basic type inference algorithm**[\[PS91\]](#). This basic type inference algorithm derives a set of constraints based on *trace graphs* and solves the constraints using a fix-point algorithm. Agesen extended the basic algorithm with *templates*. These templates are based on control flow and have start and end nodes with their possible in- and outputs. The CPA calculates the possibles output types for each template by taking the cartesian product, the set of all possible ordered pairs, of the input types.

2.4 Related work

In this section we briefly describe related research work in order to get a better understanding of similar performed researches.

Similar work has been presented by Patrick Camphuijsen[Cam07; CHH09]. Patrick created a constraint-based type inference analysis for his master’s thesis. The inference algorithm combines possible results of the constraints and takes the union to define the types. To guarantee termination the algorithm uses widening, by replacing the current result with the result of the union, to make sure that there will be a fixed-point. Further work improved the implementation by adding support for objects[VH15].

Paul Biggar created an Ahead-Of-Time (AOT) compiler for PHP[Big10]. The main goal of this compiler is to improve the performance of PHP programs. The AOT compiler starts by parsing a PHP program into an AST. This AST is transformed into an High-level Intermediate Representation (HIR) to remove all redundant constructs and then transformed into a Medium-level Intermediate Representation (MIR). Using dataflow analysis, alias analysis, static single assignment (SSA), and type analysis the compiler performs optimisations on the MIR. After the optimisations, the compiler generates C code, which then can be executed to run the program.

PHANTM[KSK10a; KSK10b] (PHp ANalyzer for Type Mismatch) is an open source PHP analyser written in Scala. Because of PHP’s dynamic nature, without compiler or interpreter type checking, it is easy to make typing errors that result in unexpected behaviour or in fatal errors. PHANTM performs a hybrid flow-sensitive analysis to find type errors in PHP5. The hybrid analysis combines static and dynamic analysis. A program can be annotated to start a static analysis at a specific point. The analyser collects run-time type information while running the program and then starts the static analysis. PHANTM uses data-flow analysis to infer types. Although PHANTM has proven to be able to find a decent number of type errors on scalar usages in three different programs, there is a lack of finding errors in object oriented structures.

Facebook improved the performance of PHP programs with a static compiler, called HipHop Virtual Machine[Zha+12] (HHVM). This static compiler extracts the program into an AST, traverses this AST to collect information, performs pre-optimisations, performs type inference, performs post-optimisations, and lastly generates C++ code. During the pre-optimisations the compiler removes unneeded actions, for example constant inlining, logical-expression simplifications, and dead-code elimination. The type inference process is based on the Hindley-Milner constraint based algorithm[DM82], to infer types of constants, variables, functions parameters, and return types. These new inferred types are then used in the post-optimisation. In the last step the AST is traversed to generate C++ code. Although the compiler does not cover all functions of PHP, it does covers most of the features. The performance benefits on the other hand are significantly better, showing on average 5.5x more efficiency for PHP5.

PHPLint³ extends the PHP syntax with type hints where PHP lacks support for it, using custom inline comment blocks to add extra typing information. These doc blocks with type information can be used in the analysis, allowing more strict type checking. The used syntax for the type hints are `/*.*`, for example: `/*.* string */ $s = $var;` which means that the variable `$s` is of type `string`. PHPLint solves the lack of type hint support on scalar and array types in PHP5. PHPLint can generate type hints based on information retrieved from simple type inference.

³<http://www.icosaedro.it/phplint>, July 2016

Chapter 3

Research Context

This chapter describes the research context. In section 3.1 we explain more about our defined types. The relation between the types are described in 3.2. Section 3.3 explains in which context the research is executed.

3.1 Types

The basis types in PHP are integers, floats, booleans, strings, arrays, resources and objects. 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.

Because PHP has no explicit type system, we define our own type system for PHP. In the Rascal code below in Rascal 1 you can see our defined types. Here we define the `TypeSymbol` data type in Rascal with our defined types as possible values. A brief description of all the defined types are listed below.

Rascal 1 `TypeSymbol` definitions in Rascal

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

data TypeSymbol
  = \any()1 // unknown, can be any of the types below
  | arrayType(TypeSymbol arrayType) // array of a type
  | booleanType() // boolean values
  | classType(loc decl) // a specific class
  | floatType() // float, double or real
  | integerType() // integer numbers
  | interfaceType(loc decl) // a specific interface
  | numberType() // a float or integer
  | nullType() // empty or undefined value
  | objectType() // any class type
  | resourceType() // a build-in type
  | scalarType() // any number, string, resource or boolean
  | stringType() // text values
;
```

any As you can see in the comments in the code above, Rascal 1, `any()` represents the combination of all possible types. This type will be used for mixed and unknown types, for example when variables are used, but are never defined.

¹any is escaped with a backslash because it is a reserved keyword.

arrayType The type `arrayType(TypeSymbol arrayType)` is a recursive declaration. The argument of the type is the type of the array. For example, an array of strings is declared as `arrayType(stringType())` and for an unknown array the type is `arrayType(\any())`.

booleanType The type `booleanType()` is the type for boolean values. Boolean values are `true` and `false`.

classType The type `classType(loc decl)` represents a specific class. The argument is the declaration, which represents the logical name of the class. An example of the `Exception` class is `classType(|php+class:///exception|)`.

floatType Floating point numbers, represented as floats, reals, and doubles are defined by the `floatType()`. Example are 1.234, 1.2e3, and 7E-10.

integerType Integers are whole numbers in decimal, hexadecimal, octal or binary notation. The `integerType` values can be positive or negative.

interfaceType The `interfaceType(loc decl)` represents a specific interface. Interfaces can be provided as type hints.

numberType The type `numberType()` covers the `integerType()` and `floatType()`. Because of type coercion in PHP, these types can be easily mixed.

nullType The type `nullType()` is used for the value `null`.

objectType The type `objectType()` is the parent type for all class types. This type represent the object type and could also been written as `classType(|php+class:///object|)`.

resourceType The type `resourceType()` represents the build-in PHP resource type. Various function return the `resourceType` from build-in PHP functions.

scalarType The type `scalarType()` is the generic type for `resourceType()`, `booleanType()`, `numberType()`, and `stringType()`.

stringType The type `stringType()` represents strings, a sequence of characters.

3.2 Type hierarchie

As described in some of the previous descriptions, types can relate to other types. The full relation schema of the types is shown in figure 3.1. In this diagram the `-Type` postfix is omitted to save space. We speak of **subtypes** when the types are descendant of the given type. The subtypes of the root node any are `scalarType`, `arrayType`, and `objectType`.

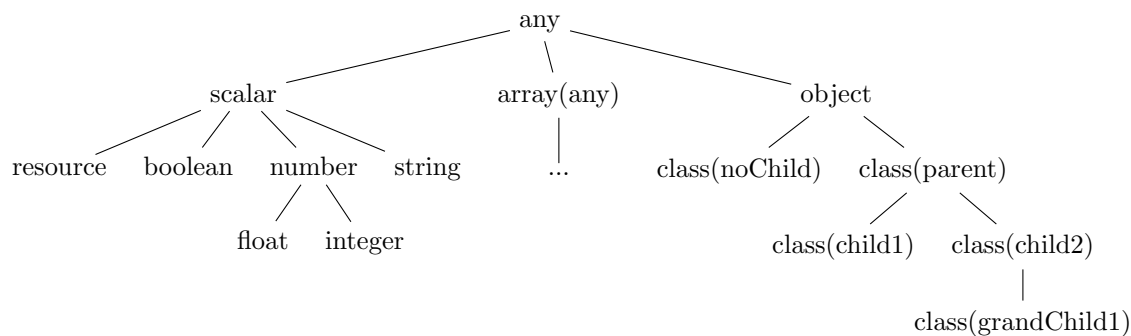


Figure 3.1: Type hierarchy

The *scalar type* is the super type for the non-complex types `resourceType`, `booleanType`, `numberTypes`, and `stringType`. These types can in practise be combined because of coercion. If they are used mixed up, they will be classified as scalar types.

The *array type* in the subtype diagram is the most generic type of array, the array of any type. We have omitted the other array types to reduce complexity of the hierarchy. The array type is a recursive type and can go to infinite depth.

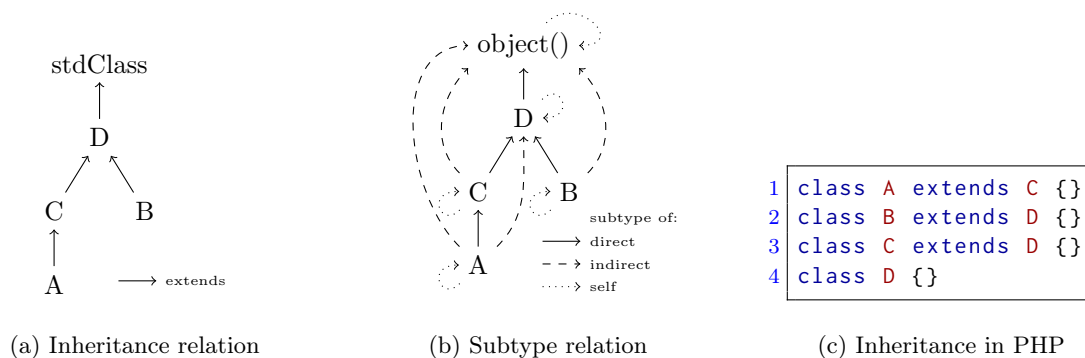


Figure 3.2: Relation of subtypes among classes

The *object type* is the most generic object type, which represents the `stdClass` in PHP. The class inheritance relation in PHP 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.

3.3 Research context

In order to let our research take place, we need to make sure that some environment variables are constant.

Program correctness In order to be able to execute this research we assume that the programs are correct and works as intended. This is needed to be able to reason about the programs we analyse, without having to question whether the program works as intended.

File includes In this research we assume that all PHP file are included during runtime within the project folder. When a PHP system is constructed of classes with namespaces, the files will be logically loaded using PHP's autoloader. Because most recent systems use namespaces, we will assume that all files are included. For our experiment we are sure that all chosen projects comply to this.

Register globals Register globals allows variables to be magically be created from GET and POST values. Since it is discouraged to use this setting, we will assume that all software products have this setting disabled. This feature is disabled by default in version 4.2 and has been completely removed in version 5.4.

PHP warnings For this research we will ignore all warnings. Warnings do not alter the behaviour of the program. In a most production environment these warnings are suppressed and will not change the behaviour of the program. We do take fatal errors into account, which lead to runtime errors.

Flow insensitive Our analysis is intra-procedural flow-insensitive, which means that we don't take the order of statements into account within a function or method scope. We do assume type consistency within a scope.

Chapter 4

Design of PHP type inference

In this chapter we present the type constraint rules, in section 4.1, for PHP language constructs with supporting code examples. In the last two sections we provide more information on annotations in section 4.2, and on PHP built-ins in sections 4.3.

4.1 Type inference rules

The constraint definitions we use in this section are based on the definition of Palsberg and Schwartzbach[PS94]. We have extended the definition to conform to the PHP language. A legend with all symbols is displayed in table 4.1, followed by the constraint definitions for PHP.

symbol	description	symbol	description
\equiv	= equivalent expression	$=$	= equivalent type
$:=$	= assignment	$!$	= negation
$<:$	= (lhs) is subtype (of rhs)	$:>$	= (lhs) is supertype (of rhs)
C	= a class	$\rightarrow c$	= class constant
E_k	= an expression	$\rightarrow p$	= class property
$\llbracket E_k \rrbracket$	= typeset of an expression	$\rightarrow m$	= class method
f	= a function	$\llbracket m \rrbracket$	= (return) typeset of a method call
$\llbracket f \rrbracket$	= (return) typeset of a function	(A_n)	= the n'th actual argument
$:: c$	= static class property	(P_n)	= the n'th formal parameter
$:: m$	= static class method	th	= type hint
$:: p$	= static class property	v	= default value
Mfs	= modifiers	\in	= is defined in
$\{ \}$	= set of types	is_a	= (lhs) is of (rhs) type
\wedge	= conjunction	\vee	= disjunction

Table 4.1: Constraint definition legend

We write the definitions in the following form:

$$\frac{\text{premiss 1} \quad \text{premiss 2}}{\text{constraint 1,} \\ \text{constraint 2}}$$

Above the horizontal line we write the premisses. In our case premisses are PHP expressions which are true or false depending on the context. If the premiss is true for a PHP statement or expression we can define the constraints below the horizontal line.

4.1.1 Scalars

Extracting constraints from the scalar types is straight forward. We show the constraint rules for strings, integers, floats, booleans, and null values.

Strings Strings in PHP can be written with single or double quotes. If an expression is a literal string, which we know from the AST, we can add the constraint that the typeset of that expression is equal to a string type.

$$\frac{E \text{ is_a string}}{\llbracket E \rrbracket = \{ \text{stringType}() \}}$$

Code sample:

```
1 "Str"; // stringType()
2 'abc'; // stringType()
```

Listing 4.1: Strings

Integers In the listing below you can see different types of integers. If we encounter a expression that represents one of these integer formats, we can extract the constraint that the typeset of the expression should be equal to an integer type.

$$\frac{E \text{ is_a integer}}{\llbracket E \rrbracket = \{ \text{integerType}() \}}$$

Code sample:

```
1 1234; // integerType() (decimal number)
2 -123; // integerType() (negative number)
3 0123; // integerType() (octal number)
4 0x1A; // integerType() (hexadecimal number)
5 0b11111111; // integerType() (binary number)
```

Listing 4.2: Integers

Floats If we see php syntax that represents a floating number, we can extract the constraint that this expression is of floating type.

$$\frac{E \text{ is_a float}}{\llbracket E \rrbracket = \{ \text{floatType}() \}}$$

Code sample:

```
1 1.4; // floatType()
2 1.2e3; // floatType()
3 7E-10; // floatType()
```

Listing 4.3: Floats

Boolean values Boolean values in PHP are case sensitive, as you can see in the examples. If we encounter a boolean value, we can extract the constraint the this expression is of boolean type.

$$\frac{E \equiv \text{true} \vee E \equiv \text{false}}{\llbracket E \rrbracket = \{ \text{booleanType}() \}}$$

Code sample:

```
1 true; // booleanType()
2 false; // booleanType()
3 TRUE; // booleanType()
4 FALSE; // booleanType()
```

Listing 4.4: Boolean values

Null values null is a reserved keyword in PHP. When we encounter null in the source code we can add the nullType type constraint.

$$\frac{E \equiv \text{null}}{\llbracket E \rrbracket = \{ \text{nullType}() \}}$$

Code sample:

```
1 null; // nullType()
2 NULL; // nullType()
```

Listing 4.5: Null values

4.1.2 Assignments

Assign statements transfer values from one expression or variable into another. PHP uses the = symbol as assignment syntax. In the premiss we use := for assigns.

Assignment When an assignment is used, we can extract the following constraint: the right hand side (E_2) of the assignment is a subtype of the left hand side (E_1). This relation is a subtype relation, not an is equal relation, because of the subclass relations of inheritance. The whole expression (E) is equal to the newly assigned value.

$$\frac{E \equiv (E_1 := E_2)}{\llbracket E_2 \rrbracket <: \llbracket E_1 \rrbracket, \llbracket E_1 \rrbracket = \llbracket E \rrbracket}$$

Code sample:

```
1 $a = $b; // [$b] <: [$a]
2 // [$a] = [$a = $b]
3
4 $c = $d = $e; // [$e] <: [$d]
5 // [$d] <: [$c],
6 // [$d] = [$d = $e]
7 // [$c] = [$c = $d = $e]
```

Listing 4.6: Assignment

Ternary operator The *ternary* operator is a conditional assignment. If the expression E_1 is evaluated as **true**, the left hand side (E_2) is the value of the whole ternary expression. If E_1 is evaluated as **false**, the right hand side (E_3) is the value. The constraint we can extract from the ternary expression is that the type of the whole expression should be the type of E_2 or E_3 (i).

The ternary operator without a left hand side value, also known as the elvis operator, returns the value of E_1 when E_1 is evaluated as **true**. Here the type of the expression should be either the type of E_1 or E_3 (ii).

$$\frac{E \equiv (E_1 ? E_2 : E_3)}{\llbracket E \rrbracket = \llbracket E_2 \rrbracket \vee \llbracket E_3 \rrbracket} \text{ (i)} \quad \frac{E \equiv (E_1 ? : E_3)}{\llbracket E \rrbracket = \llbracket E_1 \rrbracket \vee \llbracket E_3 \rrbracket} \text{ (ii)}$$

Code sample:

```
1 $a ? $b : $c; // [$a ? $b : $c] = ([ $b ] || [ $c ])
2 $a ?: $c;      // [$a ?: $c] = ([ $a ] || [ $c ])
```

Listing 4.7: Ternary operator

Assignments resulting in integers PHP provides several assignment statements combined with operators. The type of the left hand side (E_1) is in the cases of *bitwise and* (i), *bitwise inclusive or* (ii), *bitwise exclusive or* (iii), *bitwise shift left* (iv), *bitwise shift right* (v), and *modulus* (vi) always of *integer* type.

$$\begin{array}{lll} \frac{E_1 \&= E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}} \text{ (i)} & \frac{E_1 |= E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}} \text{ (ii)} & \frac{E_1 \^{}= E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}} \text{ (iii)} \\ \frac{E_1 <<= E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}} \text{ (iv)} & \frac{E_1 >>= E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}} \text{ (v)} & \frac{E_1 \% = E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}} \text{ (vi)} \end{array}$$

Code sample:

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

Listing 4.8: Assignments resulting in integers

Assignment with string concat When the *string concat* operator is used, in combination with the assignment operator (i), the type of the left hand side (E_1) is always a string.

About the right hand side (E_2) we can say that **if** the type of E_2 is a subtype of object, then this object should have the method `__toString()` (ii).

$$\frac{E_1 .= E_2}{\llbracket E_1 \rrbracket = \{ \text{stringType}() \}} \text{ (i)} \quad \frac{(E_1 .= E_2) \quad (\llbracket E_2 \rrbracket <: \{ \text{objectType}() \})}{\llbracket E_2 \rrbracket \text{ hasMethod } \text{"__toString"}} \text{ (ii)}$$

Code sample:

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

Listing 4.9: Assignment with string concat

Assignments with division or subtraction operator *Division* (i) and *subtraction* (ii) assignment in PHP will always result in an integer type. This is the case for all values, except for array's. A fatal error will occur when the right hand side value is of type array.

$$\frac{E_1 / = E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}, \llbracket E_2 \rrbracket \neq \{ \text{arrayType}(_) \}} \text{ (i)} \quad \frac{E_1 - = E_2}{\llbracket E_1 \rrbracket = \{ \text{integerType}() \}, \llbracket E_2 \rrbracket \neq \{ \text{arrayType}(_) \}} \text{ (ii)}$$

Code sample:

```
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 4.10: Assignments with division or subtraction operator

Assignments resulting in numbers The result of an *multiplication* (i) and *addition* (ii) assignment is either a float or an integer. When the type of the right hand side (E_2) is either `booleanType`, `integerType`, or `nullType`, the result of the assignment (E_1) will be of `integerType`. If E_2 is of any other type, E_1 will be of type `floatType`. Float and integer are both subtypes of integers, so we can use the subtype relation for `numberType` for this.

$$\frac{E_1 * = E_2}{\llbracket E_1 \rrbracket <: \{ \text{numberType}() \}} \text{ (i)} \quad \frac{E_1 += E_2}{\llbracket E_1 \rrbracket <: \{ \text{numberType}() \}} \text{ (ii)}$$

Code sample:

```
1 $a *= $b; // [$a] <: numberType()
2 $a += $b; // [$a] <: numberType()
```

Listing 4.11: Assignments resulting in numbers

4.1.3 Unary operators

Unary operators in PHP consist of positive and negative numbers, negation operators, and increase and decrease operators.

Positive and negative numbers When a *plus* (i) or *minus* (ii) sign is used in PHP in front of a variable, the type of the whole expression must be of `numberType`. The type of the variable cannot be of any `arrayType`.

$$\frac{E \equiv (+E_1)}{\llbracket E \rrbracket <: \{ \text{numberType}() \}, \llbracket E_1 \rrbracket \neq \{ \text{arrayType}(_) \}} \text{ (i)} \quad \frac{E \equiv (-E_1)}{\llbracket E \rrbracket <: \{ \text{numberType}() \}, \llbracket E_1 \rrbracket \neq \{ \text{arrayType}(_) \}} \text{ (ii)}$$

Code sample:

```

1 +$a; // [$a] <: numberType();
2     // [$a] != arrayType();
3 -$a; // [$a] <: numberType();
4     // [$a] != arrayType();

```

Listing 4.12: Positive and negative numbers

Negation operators The PHP language holds two types of negation operators. The type of the whole expression for *normal negation* operator (i) is boolean. For the *bitwise negation* operator (ii) the type of attached variable is either a number or a string. The type of the whole expression is an integer or string.

$$\begin{array}{c}
\frac{E \equiv (!E_1)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \quad \text{(i)} \quad \frac{E \equiv (\sim E_1)}{\begin{array}{l} \llbracket E_1 \rrbracket = \{ \text{numberType}(), \text{stringType}() \}, \\ \llbracket E \rrbracket = \{ \text{integerType}(), \text{stringType}() \} \end{array}} \quad \text{(ii)}
\end{array}$$

Code sample:

```

1 !$a // [!$a] = booleanType()
2 ~$a // [$a] = numberType() or stringType()
3     // [~$a] = integerType() or stringType()

```

Listing 4.13: Negation operators

Post increment operators From post increment and decrement operators we can only extract conditional constraints.

If the type of E_1 is of any **array** type, the result of the expression is also of any **array** type (i).

If the type of E_1 is of **boolean** type, the result of the expression is also of **boolean** type (ii).

If the type of E_1 is of **float** type, the result of the expression is also of **float** type (iii).

If the type of E_1 is of **integer** type, the result of the expression is also of **integer** type (iv).

If the type of E_1 is of **null** type, the result of the expression is either of **integer** or **boolean** type (v).

If the type of E_1 is of any **object** type, the result of the expression is also of any **object** type (vi).

If the type of E_1 is of **resource** type, the result of the expression is also of any **resource** type (vii).

If the type of E_1 is of **string** type, the result of the expression is either of **number** or **string** type (vii).

The rules below are only written for the post increment, but also apply on the post decrement.

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket <: \{ \text{arrayType}(_) \})}{\llbracket E \rrbracket <: \{ \text{arrayType}(_) \}} \quad \text{(i)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket = \{ \text{booleanType}() \})}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \quad \text{(ii)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket = \{ \text{floatType}() \})}{\llbracket E \rrbracket = \{ \text{floatType}() \}} \quad \text{(iii)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket = \{ \text{integerType}() \})}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \quad \text{(iv)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket = \{ \text{nullType}() \})}{\llbracket E \rrbracket = \{ \text{integerType}(), \text{nullType}() \}} \quad \text{(v)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket <: \{ \text{objectType}() \})}{\llbracket E \rrbracket <: \{ \text{objectType}() \}} \text{ (vi)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket = \{ \text{resourceType}() \})}{\llbracket E \rrbracket = \{ \text{resourceType}() \}} \text{ (vii)}$$

$$\frac{(E \equiv (E_1 + +)) \quad (\llbracket E_1 \rrbracket = \{ \text{stringType}() \})}{\llbracket E \rrbracket <: \{ \text{numberType}(), \text{stringType}() \}} \text{ (viii)}$$

Code sample:

```

1 $a++ // (post increase)
2 // if ([ $a ] <: arrayType()) => [ $a++ ] <: arrayType()
3 // if ([ $a ] = booleanType()) => [ $a++ ] = booleanType()
4 // if ([ $a ] = floatType()) => [ $a++ ] = floatType()
5 // if ([ $a ] = integerType()) => [ $a++ ] = integerType()
6 // if ([ $a ] = nullType()) => [ $a++ ] = integerType() or nullType()
7 // if ([ $a ] <: objectType()) => [ $a++ ] <: objectType()
8 // if ([ $a ] = resourceType()) => [ $a++ ] = resourceType()
9 // if ([ $a ] = stringType()) => [ $a++ ] <: numberType() or stringType()
10 $a-- // (post decrease)
11 // same rules as above apply for $a--

```

Listing 4.14: Post increment operators

Pre increment operators From pre increment and decrement operators we can also only extract conditional constraints. The rules are similar to the rules for the post increment, except for the `nullType()`. If the type of E_1 is of null type, the result of the expression is either of null type (v).

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket <: \{ \text{arrayType}(_) \})}{\llbracket E \rrbracket <: \{ \text{arrayType}(_) \}} \text{ (i)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket = \{ \text{booleanType}() \})}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (ii)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket = \{ \text{floatType}() \})}{\llbracket E \rrbracket = \{ \text{floatType}() \}} \text{ (iii)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket = \{ \text{integerType}() \})}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (iv)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket = \{ \text{nullType}() \})}{\llbracket E \rrbracket = \{ \text{nullType}() \}} \text{ (v)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket <: \{ \text{objectType}() \})}{\llbracket E \rrbracket <: \{ \text{objectType}() \}} \text{ (vi)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket = \{ \text{resourceType}() \})}{\llbracket E \rrbracket = \{ \text{resourceType}() \}} \text{ (vii)}$$

$$\frac{(E \equiv (+ + E_1)) \quad (\llbracket E_1 \rrbracket = \{ \text{stringType}() \})}{\llbracket E \rrbracket <: \{ \text{numberType}(), \text{stringType}() \}} \text{ (viii)}$$

Code sample:

```

1 ++$a // (pre increase)
2 // if ([ $a ] <: arrayType()) => [ ++$a ] <: arrayType()
3 // if ([ $a ] = booleanType()) => [ ++$a ] = booleanType()
4 // if ([ $a ] = floatType()) => [ ++$a ] = floatType()
5 // if ([ $a ] = integerType()) => [ ++$a ] = integerType()
6 // if ([ $a ] = nullType()) => [ ++$a ] = nullType()
7 // if ([ $a ] <: objectType()) => [ ++$a ] <: objectType()
8 // if ([ $a ] = resourceType()) => [ ++$a ] = resourceType()
9 // if ([ $a ] = stringType()) => [ ++$a ] <: numberType() or stringType()
10 --$a // (pre decrease)
11 // same rules as above apply for --$a

```

Listing 4.15: Pre increment operators

4.1.4 Binary operators

Addition-, subtraction-, multiplication-, division-, modulus-, bitwise-, comparison-, and logical operators are binary operators.

Addition operator The result of an addition operator will always be a number or an array (i). If the left and right hand side are both arrays, the return type will be array (ii). In this case two arrays are merged. In all other cases the result of this operation is a number (iii).

$$\frac{E \equiv (E_1 + E_2)}{\llbracket E \rrbracket <: \{ \text{arrayType}(_), \text{numberType}() \}} \text{ (i)}$$

$$\frac{E \equiv (E_1 + E_2) \quad \llbracket E_1 \rrbracket <: \{ \text{arrayType}(_) \} \wedge \llbracket E_2 \rrbracket <: \{ \text{arrayType}(_) \}}{\llbracket E \rrbracket <: \{ \text{arrayType}(_) \}} \text{ (ii)}$$

$$\frac{E \equiv (E_1 + E_2) \quad \llbracket E_1 \rrbracket ! <: \{ \text{arrayType}(_) \} \vee \llbracket E_2 \rrbracket ! <: \{ \text{arrayType}(_) \}}{\llbracket E \rrbracket <: \{ \text{numberType}() \}} \text{ (iii)}$$

Code sample:

```

1 $a + $b // (addition)
2 // [ $a + $b ] <: arrayType() or numberType()
3 // if (([ $a ] and [ $b ]) <: arrayType(_)) => [ $a + $b ] <: arrayType(_)
4 // if (([ $a ] or [ $b ]) !<: arrayType(_)) => [ $a + $b ] <: numberType()

```

Listing 4.16: Addition operator

Subtraction multiplication division operators The *subtraction* (i), *multiplication* (ii), and *division* (iii) operators are merged together in this paragraph because they have identical behaviour. The result of these operations is always of **number** type. The operations cannot be used if one of the sides is of type **array**. Therefore we can say that the left and right hand side cannot be of **array** type.

$$\frac{E \equiv (E_1 - E_2)}{\llbracket E \rrbracket <: \{ \text{numberType}() \}, \llbracket E_1 \rrbracket ! <: \{ \text{arrayType}(_) \}, \llbracket E_2 \rrbracket ! <: \{ \text{arrayType}(_) \}} \text{ (i)}$$

$$\frac{E \equiv (E_1 * E_2)}{\llbracket E \rrbracket <: \{ \text{numberType}() \}, \llbracket E_1 \rrbracket ! <: \{ \text{arrayType}(_) \}, \llbracket E_2 \rrbracket ! <: \{ \text{arrayType}(_) \}} \text{ (ii)}$$

$$\frac{E \equiv (E_1 / E_2)}{\llbracket E \rrbracket <: \{ \text{numberType}() \},} \text{ (iii)}$$

$$\llbracket E_1 \rrbracket ! <: \{ \text{arrayType}(_) \},$$

$$\llbracket E_2 \rrbracket ! <: \{ \text{arrayType}(_) \}$$

Code sample:

```

1 $a - $b // (subtraction)
2 $a * $b // (multiplication)
3 $a / $b // (division)
4 // [$a - $b] <: numberType()
5 // [$a * $b] <: numberType()
6 // [$a / $b] <: numberType()
7 // [$a] !<: arrayType()
8 // [$b] !<: arrayType()

```

Listing 4.17: Subtraction multiplication division operators

Modulus and bitwise shift operators The merge of *modulus* (i) and *bitwise shift* (ii, iii) operators seems not so obvious at first, but they have the same behaviour. The results of these operations is of integer type.

$$\frac{E \equiv (E_1 \% E_2)}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (i)} \quad \frac{E \equiv (E_1 << E_2)}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (ii)} \quad \frac{E \equiv (E_1 >> E_2)}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (iii)}$$

Code sample:

```

1 $a % $b // [$a % $b] = integerType() // (modulus)
2 $a << $b // [$a << $b] = integerType() // (bitwise shift left)
3 $a >> $b // [$a >> $b] = integerType() // (bitwise shift right)

```

Listing 4.18: Modulus and bitwise shift operators

Bitwise operators The results of the bitwise operators *and* (i, ii, iii), *or*, and *xor* is always of **integer** or **string** type. When the left and right hand side are both strings, the result of the operation is also of type **string**. In all other cases the result of this operation is a number. For readability reasons we omitted the constraints for the operators *bitwise or* (inclusive or, |) and *bitwise xor* (exclusive or, ^), because they provide the same constraints as the *bitwise and* operator.

$$\frac{E \equiv (E_1 \& E_2)}{\llbracket E \rrbracket = \{ \text{stringType}(), \text{integerType}() \}} \text{ (i)}$$

$$\frac{E \equiv (E_1 \& E_2) \quad \llbracket E_1 \rrbracket = \{ \text{stringType}() \} \wedge \llbracket E_2 \rrbracket = \{ \text{stringType}() \}}{\llbracket E \rrbracket = \{ \text{stringType}() \}} \text{ (ii)}$$

$$\frac{E \equiv (E_1 \& E_2) \quad \llbracket E_1 \rrbracket \neq \{ \text{stringType}() \} \vee \llbracket E_2 \rrbracket \neq \{ \text{stringType}() \}}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (iii)}$$

Code sample:

```

1 $a & $b // (bitwise And)
2         // [$a & $b] = stringType() or integerType()
3         // if (($a and $b) = stringType()) => [$a & $b] = stringType()
4         // if (($a or $b) != stringType()) => [$a & $b] = integerType()
5 $a | $b // (bitwise Or)
6         // [$a | $b] = stringType() or integerType()
7         // if (($a and $b) = stringType()) => [$a | $b] = stringType()
8         // if (($a or $b) != stringType()) => [$a | $b] = integerType()
9 $a ^ $b // (bitwise Xor)
10        // [$a ^ $b] = stringType() or integerType()
11        // if (($a and $b) = stringType()) => [$a ^ $b] = stringType()
12        // if (($a or $b) != stringType()) => [$a ^ $b] = integerType()

```

Listing 4.19: Bitwise operators

Comparison operators The result of the comparison operators is always of boolean type. The comparison operators are *equals* (i), *identical* (ii), *not equal* (iii), *not equal* (iv), *not identical* (v), *less than* (vi), *greater than* (vii), *less than or equal to* (viii), and *greater than or equal to* (ix) operators.

$$\begin{array}{ccc}
\frac{E \equiv (E_1 == E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (i)} & \frac{E \equiv (E_1 === E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (ii)} & \frac{E \equiv (E_1 != E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (iii)} \\
\frac{E \equiv (E_1 <> E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (iv)} & \frac{E \equiv (E_1 !== E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (v)} & \frac{E \equiv (E_1 < E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (vi)} \\
\frac{E \equiv (E_1 > E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (vii)} & \frac{E \equiv (E_1 <= E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (viii)} & \\
\frac{E \equiv (E_1 >= E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (ix)} & &
\end{array}$$

Code sample:

```

1 $a == $b // [$a == $b] = booleanType()
2 $a === $b // [$a === $b] = booleanType()
3 $a != $b // [$a != $b] = booleanType()
4 $a <> $b // [$a <> $b] = booleanType()
5 $a !== $b // [$a !== $b] = booleanType()
6 $a < $b // [$a < $b] = booleanType()
7 $a > $b // [$a > $b] = booleanType()
8 $a <= $b // [$a <= $b] = booleanType()
9 $a >= $b // [$a >= $b] = booleanType()

```

Listing 4.20: Comparison operators

Logical operators Just like the comparison operators, the result of the logical operators is always of boolean type. The logical operators are *and* (i), *or* (ii), *xor* (iii), *and* (iv), and *or* (v).

$$\begin{array}{ccc}
\frac{E \equiv (E_1 \text{ and } E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (i)} & \frac{E \equiv (E_1 \text{ or } E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (ii)} & \frac{E \equiv (E_1 \text{ xor } E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (iii)} \\
\frac{E \equiv (E_1 \&\& E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (iv)} & \frac{E \equiv (E_1 \|\| E_2)}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (v)} &
\end{array}$$

Code sample:

```

1 $a and $b // [$a and $b] = booleanType()
2 $a or $b  // [$a or $b]  = booleanType()
3 $a xor $b // [$a xor $b] = booleanType()
4 $a && $b  // [$a && $b]  = booleanType()
5 $a || $b  // [$a || $b]  = booleanType()

```

Listing 4.21: Logical operators

4.1.5 Array

For arrays we define array declaration and array access. We rely on information in the AST to extract constraints from arrays.

Array declaration From the array declarations (`array()` or `[]`) we can extract the constraint that they should be of any `array` type.

$$\frac{E \text{ is_a array}}{\llbracket E \rrbracket <: \{ \text{arrayType}(_) \}}$$

Code sample:

```

1 array(/*...*/); // [array(/*...*/)] <: arrayType(any())
2 [/*...*/];      // [[(/*...*/)]] <: arrayType(any())

```

Listing 4.22: Array declaration

Array access From the usage of array access syntax you cannot tell what the type of the expression is. The same syntax is used to access strings. We can extract that the type of the base expression should not be of object type (i). If we know that the base type is of `string` type, we know that the result of the expression will also be a string (ii). When the base type is an array, the result type is the type of the elements in there array (iii). For all other cases, when the base type is not an string or array, the result of the expression will be of `null` type (iv).

$$\frac{E_1[E_2] \quad E_1 \text{ is_a arrayAccess}}{\llbracket E_1 \rrbracket \neq \{ \text{objectType}() \}} \quad (\text{i})$$

$$\frac{E \equiv (E_1[E_2]) \quad E_1 \text{ is_a arrayAccess} \quad \llbracket E_1 \rrbracket = \{ \text{stringType}() \}}{\llbracket E \rrbracket = \{ \text{stringType}() \}} \quad (\text{ii})$$

$$\frac{E \equiv (E_1[E_2]) \quad E_1 \text{ is_a arrayAccess} \quad \llbracket E_1 \rrbracket = \{ \text{arrayType}(E_2) \}}{\llbracket E \rrbracket = \llbracket E_2 \rrbracket} \quad (\text{iii})$$

$$\frac{E \equiv (E_1[E_2]) \quad E_1 \text{ is_a arrayAccess} \quad \llbracket E_1 \rrbracket \neq \{ \text{stringType}() \} \quad \llbracket E_1 \rrbracket ! <: \{ \text{arrayType}(_) \}}{\llbracket E \rrbracket = \{ \text{nullType}() \}} \quad (\text{iv})$$

Code sample:

```

1 $a[$b];
2 // [$a] != objectType()
3 // if ([$a] == stringType())      => [$a[$b]] = stringType()
4 // if ([$a] == arrayType(x))      => [$a[$b]] = [x]
5 // if ([$a] != (string or array)) => [$a[$b]] = nullType()

```

Listing 4.23: Array access

4.1.6 Casts

Casts PHP contains syntax to arrays, booleans, integers, floats, objects, strings, and to unset variables. The result of a cast to array is of any `array` type (i). For casting to boolean there are two keywords, `bool` (ii) and `boolean` (iii), and the result will always be of `boolean` type. There are three keywords to cast to floats, `float` (iv), `double` (v), and `real` (vi). Casts to integer integer type, you can use `integer` (vii) or `int` (viii) keywords. Any cast to `object` (ix) will result in any `object` type. A cast to `string` will always result in a `string` type. String casts (x) will always result in a `string` type. If we know that the expression (E_1) is an object, we know that this method needs to have an `__toString()` method (xi). The last cast, `unset`, results in a `null` type.

$$\begin{array}{c}
\frac{E \equiv (\text{array})E_1}{\llbracket E \rrbracket <: \{ \text{arrayType}(_) \}} \text{ (i)} \quad \frac{E \equiv (\text{boolean})E_1}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (ii)} \quad \frac{E \equiv (\text{bool})E_1}{\llbracket E \rrbracket = \{ \text{booleanType}() \}} \text{ (iii)} \\
\frac{E \equiv (\text{float})E_1}{\llbracket E \rrbracket = \{ \text{floatType}() \}} \text{ (iv)} \quad \frac{E \equiv (\text{double})E_1}{\llbracket E \rrbracket = \{ \text{floatType}() \}} \text{ (v)} \quad \frac{E \equiv (\text{real})E_1}{\llbracket E \rrbracket = \{ \text{floatType}() \}} \text{ (vi)} \\
\frac{E \equiv (\text{integer})E_1}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (vii)} \quad \frac{E \equiv (\text{int})E_1}{\llbracket E \rrbracket = \{ \text{integerType}() \}} \text{ (viii)} \quad \frac{E \equiv (\text{object})E_1}{\llbracket E \rrbracket <: \{ \text{objectType}() \}} \text{ (ix)} \\
\frac{E \equiv (\text{string})E_1}{\llbracket E_1 \rrbracket = \{ \text{stringType}() \}} \text{ (x)} \quad \frac{E \equiv (\text{string})E_1 \quad (E_1 <: \{ \text{objectType}() \})}{\llbracket E_1 \rrbracket <: \{ \text{hasMethod} \text{ " __toString" } \}} \text{ (xi)} \\
\frac{E \equiv (\text{unset})E_1}{\llbracket E \rrbracket = \{ \text{nullType}() \}} \text{ (xii)}
\end{array}$$

Code sample:

```

1 (array)$a // [(array)$a] <: arrayType(_)
2 (bool)$a  // [(bool)$a] = booleanType()
3 (float)$a // [(float)$a] = floatType()
4 (int)$a   // [(int)$a] = integerType()
5 (object)$a // [(object)$a] = objectType()
6 (string)$a // [(string)$a] = stringType()
7           // if ($a <: objectType()) => [$a] has method "__toString()"
8 (unset)$a // [(unset)$a] = nullType()

```

Listing 4.24: Casts

4.1.7 Clone

Clone From the PHP function `clone` we can extract the constraint that the type of the given expression and the result must be of any `object` type. We also know that the type will not change, and so the type of the expression will be the same as the type of the variable that is cloned.

$$\frac{E \equiv \text{clone}(E_1)}{\begin{array}{l} \llbracket E \rrbracket <: \{ \text{objectType}() \} \\ \llbracket E_1 \rrbracket <: \{ \text{objectType}() \} \\ \llbracket E \rrbracket = \llbracket E_1 \rrbracket \end{array}}$$

Code sample:

```
1 clone($a) // [$a] <: object
2           // [clone($a)] <: object
3           // [$a] = [clone($a)]
```

Listing 4.25: Clone

4.1.8 Class

This section contains fact extraction rules for class instantiations, special keywords, and method calls.

Class instantiation Classes can be instantiated with the name of the class. The type of the whole expression is then of the specific `class` type (i). When a class is dynamically instantiated, we only know that it should be of some `object` type, and that the type of the expression should be any `object` or `string` type (ii).

$$\frac{E \equiv \text{new } C_1()}{\llbracket E \rrbracket = \{ \text{classType}(C.\text{decl}) \}} \text{ (i)} \quad \frac{E \equiv \text{new } E_1}{\begin{array}{l} \llbracket E \rrbracket <: \{ \text{objectType}() \}, \\ \llbracket E_1 \rrbracket <: \{ \text{objectType}(), \text{stringType}() \} \end{array}} \text{ (ii)}$$

Code sample:

```
1 new C; // [new C] = classType(C)
2
3 $c = "C";
4 new $c; // [new $c] <: objectType()
5         // [$c] <: ( objectType() or stringType() )
```

Listing 4.26: Class instantiation

Special keywords PHP contains a few class related reserved keywords with special behaviour. These keywords can be used inside a class scope ($\in C$). From the usage of the keyword *self* we know that the type of the expression should be the same `class` type as which the keyword is defined in (i). The constraint we can extract from *self* is that the type should be any `object` type and it should be either the contained class or one of the parent classes. The behaviour of *\$this* (ii) and *static* (iii) differs, but the constraints we can extract are equal to the *self* keyword. The *parent* (iv) keyword differs because it must be a super type of the class they keyword is defined in.

$$\frac{(E \equiv \text{self}) \in C}{\begin{array}{l} \llbracket E \rrbracket <: \{ \text{objectType}() \}, \\ \llbracket E \rrbracket >: \{ \text{classType}(C) \} \end{array}} \text{ (i)}$$

$$\frac{(E \equiv \text{static}) \in C}{\begin{array}{l} \llbracket E \rrbracket <: \{ \text{objectType}() \}, \\ \llbracket E \rrbracket >: \{ \text{classType}(C) \} \end{array}} \text{ (ii)}$$

$$\frac{(E \equiv \$this) \in C}{\llbracket E \rrbracket <: \{ \text{objectType}() \}, \llbracket E \rrbracket :> \{ \text{classType}(C) \}} \text{ (iii)}$$

$$\frac{(E \equiv \text{parent}) \in C}{\llbracket E \rrbracket <: \{ \text{objectType}() \}, \llbracket E \rrbracket :> \{ \text{classType}(C) \}, \llbracket E \rrbracket != \{ \text{classType}(C) \}} \text{ (iv)}$$

Code sample:

```

1 self    // in class C -> [self]    = classType(C)
2 parent  // in class C -> [parent] = parentOf(classType(C))
3 static  // in class C -> [static] = classType(C) or <: classType(C)
4 $this   // in class C -> [$this]  = classType(C) or <: classType(C)

```

Listing 4.27: Special keywords

Method calls From the usage of a method call (*expression -> expression*) we can extract the constraint that the type of the left hand side should be an object (i and ii). If the right hand side (E_2) is a name of a method, we can extract the constraint that the left hand side (E_1) must implement this method (ii).

$$\frac{E_1 \rightarrow E_2 \in C \quad E_2 \text{ is_a expression}}{\llbracket E_1 \rrbracket <: \{ \text{objectType}() \}} \text{ (i)}$$

$$\frac{E_1 \rightarrow E_2 \in C \quad E_2 \text{ is_a name}}{\llbracket E_1 \rrbracket <: \{ \text{objectType}() \}} \text{ (ii)}$$

$$\llbracket E_1 \rrbracket = C.\text{hasMethod}(E_2.\text{name}, \text{static} \notin \text{Mfs})$$

Code sample:

```

1 $a->$b() // [$a] <: objectType()
2 $a->b()   // [$a] <: objectType()
3           // [$a] has a method (possible inherited) with the name 'b'

```

Listing 4.28: Method calls

4.1.9 Scope

In the subsection scope we explain which constraints we can extract for variables and return values.

Variables The type of a certain variable is defined by adding an equal constraint on the logical name and the location. The scope of these variables is contained in the logical name. The logical name contains the name of the scope, which is optionally a namespace and the name of a class, method or function, and the name of the variable. An example of a logical name for the variable `$a` in the function `f` in the global namespace is `|php+functionVar:///f/a`.

$$\frac{E \quad E \text{ is_a variable}}{\llbracket E_{\text{definition}} \rrbracket = \llbracket E_{\text{location}} \rrbracket}$$

Code sample:

```

1 function f() {
2     $a = 1; // [|php+functionVar:///f/a|] = [|file:///file.php|(17,2,<2,0>,<2,0>)]
3     $a = "s"; // [|php+functionVar:///f/a|] = [|file:///file.php|(27,2,<3,0>,<3,0>)]
4 }

```

Listing 4.29: Variables

Return types The type of a function or method is defined by the return statements it contains. When there are no return statements declared in a function or method (i), we can extract the constraint that the function will always return `nullType`. The type of a function is also `nullType` when there is a return statement without an expression (ii), like `return;`. When there are multiple return statements, the return type of the function or method is the concatenation of the types of the expressions (iii).

$$\frac{E \text{ is_a return } \not\subseteq f}{\llbracket f \rrbracket = \text{nullType}()} \text{ (i)} \quad \frac{\text{is_a return } E \subseteq f \quad E \text{ is_a noExpr}}{\llbracket f \rrbracket = \text{nullType}()} \text{ (ii)}$$

$$\frac{(\text{return } E_1) \vee \dots \vee (\text{return } E_k) \subseteq f \quad E_{1\dots k} \text{ is_a someExpr}}{\llbracket f \rrbracket <: \llbracket E_1 \rrbracket \vee \dots \vee \llbracket E_k \rrbracket} \text{ (iii)}$$

Code sample:

```

1 function f() {} // [f] = nullType()
2 function f() { return; } // [f] = nullType()
3
4 function f() { // [f] = [$a] or [$b]
5     return mt_rand(0,1) ? $a : $b;
6 }

```

Listing 4.30: Return types

4.2 Annotations

Annotations are pieces of meta data, defined on class, method, function, or statement level. Despite the proposal¹ for official support of annotations, PHP has still no native support for them. PHP has however a `getDocComment`² method in the `ReflectionClass` since version 5.1 in 2005. The `getDocComment` method returns the complete docblock of a certain element as a string. A docblock in php has the format `/**...*/`. Listing 4.31 shows an example of two docblocks in PHP. The first docblock is defined above the class and contains information about the class. The second docblock is related to the method `getSomething`. The block contains a short description of the method, provides type hints for the parameter and the return type, and provides information which possible exceptions can be thrown by the method.

```

1 namespace Thesis;
2
3 /**
4  * Class Example
5  * @package Thesis
6  */
7 class Example

```

¹<https://wiki.php.net/rfc/annotations-in-docblock>

²<http://php.net/manual/en/reflectionclass.getdoccomment.php>

```

8 {
9     /**
10      * This is a description of the method getSomething
11      *
12      * @param SomeTypeHint $someObject
13      * @return string
14      * @throws NoNameException
15      */
16     public function getSomething(SomeTypeHint $someObject)
17     {
18         if (null === $someObject->getName()) {
19             throw new NoNameException();
20         }
21
22         return $someObject->getName();
23     }
24 }

```

Listing 4.31: Examples of PHP DocBlocks

Annotations can be used for type hinting, documentation, and code execution. Software analysis tools and IDE's can use the type hints to aid understanding code and in finding bugs and security issues. Documentation generator are able to generate documentation based on the docblocks. Programs like Symfony2, ZEND Framework, and Doctrine ORM use annotations for controller routing, templating information, ORM mappings, filters, and validation configuration. A standard on using annotations is not in the PHP Standard Recommendations (PSR), but there is a proposal³. For this research we will only focus on the @param, @return, and @var annotations. The annotations @return and @param are only useful for functions, class methods, and closures. Type hints described with @var and can be used on all structures, but mainly occur on variables and class fields. There is no official standard for the use of annotations, but most projects follow the phpDocumentor⁴ syntax. For this research the following annotations are considered:

$$\text{@return} = \left\{ \begin{array}{l} \text{@return type,} \quad \text{unconditionally read @return type.} \end{array} \right. \quad (4.1)$$

$$\text{@param} = \left\{ \begin{array}{ll} \text{@param type \$var,} & \text{if '@param type \$var' occurs at least once.} \\ \text{@param \$var type,} & \text{else if '@param \$var type' occurs at least once.} \\ \text{@param type,} & \text{otherwise try to match '@param type'.} \end{array} \right. \quad (4.2)$$

$$\text{@var} = \left\{ \begin{array}{ll} \text{@var type \$var,} & \text{if '@var type \$var' occurs at least once.} \\ \text{@var \$var type,} & \text{else if '@var \$var type' occurs at least once.} \\ \text{@var type,} & \text{otherwise try to match '@var type'.} \end{array} \right. \quad (4.3)$$

$$\text{type} = \left\{ \begin{array}{ll} \text{type|type,} & \text{if '|' in type.} \\ \text{type,} & \text{otherwise} \end{array} \right. \quad (4.4)$$

In equation 4.1 we show the supported syntax for return annotations. This format starts with the @return annotation followed by the type hint. The return type hint is used on functions and class methods to provide information about the possible return values.

In equation 4.2 we describe the variants of @param annotations. Because the usage is not standardised we support multiple variants. The first variant has a type hint and next provides information for which variable. Next we also support the typehint and variable to be swapped.

³<https://github.com/php-fig/fig-standards/pull/169/files>, July 2014

⁴<http://www.phpdoc.org/>

In equation 4.3 we show the variants for the `@var` annotation. They have the same support as the `@param` annotation.

In equation 4.4 we show the recursive definition of a type. This is a recursive definition because one type hint can have multiple types. When there are multiple return types a `|` is used.

4.3 PHP built-ins

PHP has many built in classes, interfaces, constants, and variables that are always available to end users of the programming language. We include the PHP built-in information to see if we can improve the analysis results. To be able analyse the PHP built-ins we include the PHP representation of the classes, interfaces, constants, and variables. These files, written in PHP, provide annotation with type information. A brief example is shown in listing 4.32. From this example we can fetch the parameter and return type of the `strtoupper` and `strtolower` functions.

```
1 /**
2  * Make a string uppercase
3  * @link http://php.net/manual/en/function.strtoupper.php
4  * @param string $string <p>
5  * The input string.
6  * </p>
7  * @return string the uppercased string.
8  * @since 4.0
9  * @since 5.0
10 */
11 function strtoupper ($string) {}
12
13 /**
14  * Make a string lowercase
15  * @link http://php.net/manual/en/function strtolower.php
16  * @param string $str <p>
17  * The input string.
18  * </p>
19  * @return string the lowercased string.
20  * @since 4.0
21  * @since 5.0
22 */
23 function strtolower ($str) {}
```

Listing 4.32: Snippet of PHP file representing built-in functions

Chapter 5

Implementation of PHP type inference

The type inference is executed in 3 main steps. In the first step we create an M^3 model for a PHP program. The creation of an M^3 for PHP is explained in more details in section 5.1. Once the M^3 model is constructed, we extract constraints from the program using ASTs and the M^3 model which is explained in section 5.2. Section 5.3 describes how we solve the extracted constraints to collect type information. The final section 5.4 of this chapter explains how annotations and PHP built-in information are included in the process to improve the type inference.

5.1 M^3 for PHP

As explained in section 2.2, M^3 is a language independent meta model which holds facts about programs. In the next subsections we show the core M^3 elements and the extended PHP specific elements. We briefly describe the declared elements and what characteristics each of them have. Finally we present the required steps and the algorithm to construct a model.

5.1.1 Core elements

The M^3 model has the following core elements: declarations, containment, modifiers, uses, names, and documentation. The rascal code is displayed in Rascal 2. The characteristics of the elements for PHP programs are described in the paragraphs below.

Rascal 2 M^3 core definitions in Rascal

```
anno rel[loc name, loc src] M3@declarations; // maps declarations to file location
anno rel[loc from, loc to] M3@containment; // what is logically contained in what else
anno rel[loc definition, Modifier modifier] M3@modifiers; // associated modifiers
anno rel[loc src, loc name] M3@uses; // maps src locations of usages to the declarations
anno rel[str simpleName, loc qualifiedName] M3@names; // end-user readable names
anno rel[loc definition, loc comments] M3@documentation; // comments and doc-blocks
```

Declarations Declarations defines the declarations of namespaces, classes, interfaces, methods, functions, and variables and holds the relation between the logical name (which is used to refer to the declaration) and the actual file location (which is the physical place in the file system). For example the logical name of a class can be `|php+class://SomeNameSpace/ClassX|` while the actual location might be `|file:///project/SomeNameSpace/ClassX.php|`.

Containment Containment holds information about what elements logically contain other elements. For example, a property or method is contained in a class and a class is contained in a package. When a function is declared in another function, they are both logically contained in the global namespace (the highest level).

Modifiers Modifiers element contains information about the modifiers of classes, fields, and methods. Classes can only be `abstract`, fields can be `public`, `private`, or `protected`, and methods can be all of them. Abstract methods can only be declared in abstract classes.

Uses Uses relation holds information about the usages of certain elements. It is the relation between the usage and the declaration. For example when you instantiate a class, in that case you 'use' that specific class.

Names Names contains the declaration and a simplified version of the name. `names` contains human readable name which can for example be used for presenting the element in a GUI.

Documentation Documentation contains the link to the documentation related to the source code element. The link is mapped to a declaration.

5.1.2 PHP specific elements

Because every programming language differs in syntax and semantics, the M^3 model is extensible to provide language specific elements. The following PHP specific items are added: `extends`, `implements`, `parameters`, `constructors`, and `annotations`. These PHP specific elements are described in more detail in the paragraphs below.

Rascal 3 PHP specific M^3 elements

```
anno rel[loc from, loc to] M3@extends;    // which class extends which class
anno rel[loc from, loc to] M3@implements; // which class implements which interface
anno rel[loc decl, PhpParams params] M3@parameters; // formal functions/methods parameters
anno rel[loc decl, loc to] M3@constructors; // constructor usages
anno rel[loc pos, Annotation annotation] M3@annotations; // annotations from doc blocks
```

Extends Extends contains information about what classes and interface extend other classes or interfaces. Please note that we do not hold information about which class implements which interface, because that information is contained in the `implements` relation. Interface extensions work just like class extensions.

Implements Implements holds information on which class implements which interfaces. One class can implement no, one, or multiple interfaces. Because the information is a relation between the class and the interface, we can easily add multiple interfaces to one class in the model.

Parameters Parameters keeps track of the parameters of a method or function. `PhpParams` is a list relation and contains the optional typehint, if the parameter is required and if it is passed as reference. This information is stored to make it easier to resolve the call to a method or function.

Constructors Constructors lists the constructors for classes. This information is needed because it is not always clear what constructor is used, due to legacy PHP4 way of using class constructors. In PHP4 the constructor was defined as a method with the same name as the class. Since PHP5 the language is provided with a magic method `__construct()`, which results in two ways to have constructors, but only one constructor will be called. The PHP4 constructors will be removed in PHP7.

Annotations Annotations contain a relation between a declaration and the annotations that are known. Annotations are defined in the raw doc blocks and are parsed using regular expressions (regex). Regex was in this case easier then parsing because there is no official grammar or standard defined. For this research we only use `@param`, `@var`, `@returns`.

5.1.3 Constructing a model

An diagram with the required steps to created an M^3 for a PHP program is shown in figure 5.1. The figure shows that independent M^3 models are build for each PHP file in the program. Program facts are collected by traversing the AST representation of the PHP script. Once all individual M^3 s are created, they are combined into one program M^3 . M^3 s are first created for each file because dependencies between files can only be resolved once all files are analysed.

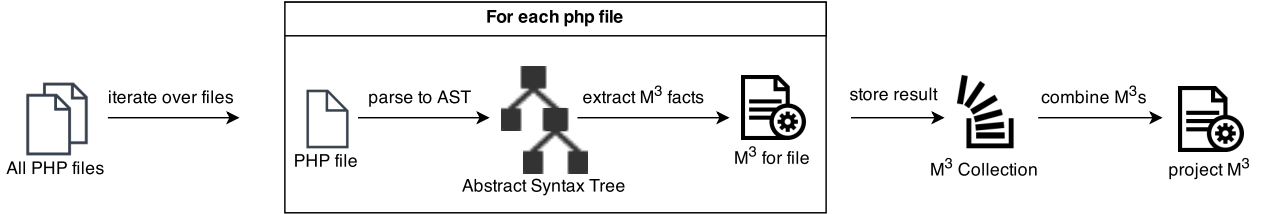


Figure 5.1: M^3 Creation

To show the process in more details provided the algorithm below in Algorithm 1.

Algorithm 1: PHP program to M^3

Input: PHP files of a program

Output: M^3 for the PHP program

```

1 m3Collection = [];
2 forall the file ∈ program do
3     ast = parseUsingPhpParserAndReturnRascalAST(file);
4     m3 = createEmptyM3(ast);
5     m3 = addDeclarations(m3, ast);
6     ast = addScopeInformation(m3, ast);
7     m3 = addContainment(m3, ast);
8     m3 = addExtendsAndImplements(m3, ast);
9     m3 = addModifiers(m3, ast);
10    m3 = addRawDocBlocksAndAnnotations(m3, ast);
11    m3 = calculateUsesFlowInsensitive(m3, ast);
12    m3Collection += m3;
13 end
14 return projectM3 = composePhpM3(m3Collection);

```

The input of the algorithm is all PHP files of a program, which are all files ending on .php within a program. **The output** is an M^3 -model with facts about the provided program. The algorithm starts with initialising an empty M^3 collection in line 1 which will be filled with the result of each individual file. In the big loop on line 2 to 13 we iterate over all the PHP files of the program. The first thing that needs to be done is to create an **Abstract Syntax Tree** (AST) of the program using an external PHP parser¹ which returns an Rascal AST in `parseUsingPhpParserAndReturnRascalAST` on line 3. From this AST we create an empty M^3 in `createEmptyM3` on line 4. In order to be able to refer to any source code element, we need to have the declarations of the elements. These elements are extracted in `addDeclarations` on line 5 can be namespace, class, interface, method, function, variable, field, or constant. For all functions we need to add scope information in case functions are declared inside another function or method. In line 6 we add the scope information to all nodes to the AST by visiting the AST. This is

¹<https://github.com/ruudvanderweijde/PHP-Parser>

needed in order to extract the right facts about the logical containment which is done in line 7. The next three lines (8-10) extract facts about class inheritance and the interface implementations, modifiers, PHP doc blocks, and annotations. Once all the basic information is collected, `calculateUsesFlowInsensitive` on line 11 tries to find the declarations of the used objects, methods, functions and variables. This is flow and context insensitive and only tries to resolve non-dynamic language constructs. The last step of the iterator is to add the constructed M^3 to the M^3 collection. Finally when an M^3 is constructed for all files, all facts of the individual M^3 's are merged into one M^3 model in line 14.

5.2 Constraint extraction

In section 4.1 we've defined the constraint extraction rules. These constraint rules are implemented in Rascal with the help of visiting² and pattern matching³. Visiting is a concept implemented in the core of Rascal and is based on the visitor design pattern. This makes it very easy to traverse trees like AST's, without coupling the design pattern to your code. During the visiting of the AST we use pattern matching on function arguments to find to correct functions to handle the constraint extraction.

5.2.1 Constraint definitions in Rascal

The constraints are implemented in Rascal with the `TypeOf` and `Constraint` definitions, which are shown in Rascal 4.

The `TypeOf` definition represents a literal type or the type of an expression. To define a literal type, we use `typeSymbol` which takes a `TypeSymbol`. `TypeSymbols` are defined in Rascal 1. We use the `typeOf` definition when we define the type of an expression or statement and takes a Rascal location (`loc`) as argument. The location represents the physical location of the expression in the source code.

The `Constraint` definition consists of recursive and non-recursive definitions. `eq` states that the type of the left hand side should be equal to the type of the right hand side. `subtyp` states that the type of the left hand side should be a subtype of the right hand size. This subtype relation is self inclusive, which means that the type can also be the same. `supertyp` states that the type of the left hand side should be a supertype of the right hand side. This relation is also self inclusive. `isMethodOfClass` takes as arguments two expressions and a string. The two expressions represent the possible types for the left and right hand side of the caller expression. The string contains the name of the called method.

The recursive definitions allow support for negation, conjunction, disjunction. The `negation` is added to make it easy invert complex constraints, for example when you can say that the type of an expression can be anything but a string. `conjunction` and `disjunction` are used for and and or relations.

Rascal 4 TypeOf and Constraint definitions

```
data TypeOf
  = typeSymbol(TypeSymbol ts)
  | typeOf(loc ident)
  ;

data Constraint
  = eq(TypeOf lhs, TypeOf rhs)
  | subtyp(TypeOf lhs, TypeOf rhs)
  | supertyp(TypeOf lhs, TypeOf rhs)
  | isMethodOfClass(TypeOf expr, TypeOf classVariable, str name)

  | disjunction(set[Constraint] constraints)
  | conjunction(set[Constraint] constraints)
  | negation(Constraint constraint)
  ;
```

²<http://tutor.rascal-mpl.org/Rascal/Concepts/Visiting/Visiting.html>

³<http://tutor.rascal-mpl.org/Rascal/Concepts/PatternMatching/PatternMatching.html>

5.2.2 Constraint extraction in Rascal

A code snippet of the constraint extraction implementation⁴ is shown in Rascal 5. In this snippet we show a small subset of the code used in the constraint extraction process.

The code starts with a local declaration of `constraints`, an empty set. This set of constraints will be filled in the rest of the functions on this Rascal module.

You can get all the constraints of a program by calling the only public function `getConstraints` with a `System` and `M3` as parameters. `System` represents a collection of all AST's of a program. The provided `M3` is the M^3 -model of the whole program, containing facts about the program. In this function we loop over all `Scripts` of a system and call `addConstraints` for each script.

Using argument pattern matching the function `addConstraints(Script script, M3 m3)` is triggered. The first thing this function executes is `addConstraintsOnAllVarsForScript`, where we add the constraints that the types of a variable should be of one type within a scope. Next we loop over all the statements of the body of a `script`.

Because of pattern matching, the function `private void addConstraints Stmt statement, M3 m3)` will be called for each statement. In this function we execute a *top-down-break visit*. *Top-down* means that the visit will start at the root node and then moves down towards the leaves. *Break* means that the in dept visit will stop as soon as there is a pattern match. We need to include *break* in this visit because for some nodes we need context. For example if a variable names `$this` would match, we need to know in which class this item was declared. In the code snippet we've only included `exprstmt`, a match with a PHP expression.

The last function of the snippet, `private void addConstraints(Expr e, M3 m3)`, handles expressions like assign with operator, a string scalar, and PHP variables that has a name. With this small code snippet we can show the constraint extraction with a tiny example.

Tiny example With this part of the code, we could extract constraints for the file with the content: `<?php $a += '1';`. When we extract the constraints from this file we would come to the following constraints:

- the type of `$a` is a subtype of `numberType()`
- the type of `$a` is a subtype of `any()`
- the type of `'1'` is equal to `stringType()`

See section 4.1 to see the full list of constraints that we extract from the source. The full implementation can be found on github.

⁴The complete source code can be found in <https://github.com/ruudvanderweijde/php-analysis>

Rascal 5 Constraint extraction snippet

```
private set[Constraint] constraints = {};

public set[Constraint] getConstraints(System system, M3 m3) {
    for (s <- system)
        addConstraints(system[s], m3);
    return constraints;
}

private void addConstraints(Script script, M3 m3) {
    addConstraintsOnAllVarsForScript(script, m3);
    for (stmt <- script.body)
        addConstraints(stmt, m3);
}

private void addConstraints Stmt statement, M3 m3) {
    top-down-break visit (statement)
    case exprstmt(Expr expr): addConstraints(expr, m3);
}

private void addConstraints(Expr e, M3 m3) {
    top-down-break visit (e) {
        case a:assignWOp(Expr assignTo, Expr assignExpr, Op operation): {
            addConstraints(assignTo, m3);
            addConstraints(assignExpr, m3);

            switch(operation)
            case plus(): constraints += {subtyp(typeOf(assignTo@at), numberType())};
        }
        case v:var(name(name(name))):
            constraints += {subtyp(typeOf(v@at), \any())};

        case s:s:scalar(Scalar scalarVal):
            switch(scalarVal)
            case string(_): constraints += {eq(typeOf(s@at), stringType())};
    }
}
```

5.3 Constraint solving

When all constraints are extracted from the source code, we will solve the constraints until we reach a fixed-point. The result will be a list of possible types for each class, method, fields, functions, variable and expression.

The first step is to initialise all type-able objects. In this initial phase the types of each object is of type any, except for the literal types which can be resolved already. When we solve the constraints step by step, we will be able to limit the number of possible types for a certain object. We do this by taking the intersection of the constraint result and the possible types we have. This way there should be less and less possible types for each variable.

When we take the intersection of none overlapping types, we have a type error. There is no possible type for this object, resulting in an empty set of possible types.

5.3.1 The algorithm

In algorithm 2 we show the algorithm to solve the constraints.

Algorithm 2: Constraint solving algorithm

Input: set[Constraint] constraints

Output: map[TypeOf expression, set[Type] possibleTypes] solutions

```
1 map[TypeOf, set[Type]] solutions = init(constraints);
2 while changes in constraints or solutions do
3   constraints = constraints + deriveMore(constraints, solutions);
4   solutions = propagateSolutions(constraints, solutions);
5 end
6 return solutions;
```

The algorithms input is a set of constraints and the output is a map with solutions. The set of constraints are the constraints defined in constraint extracting section, and the result map is the location of all type-able objects with their possible solutions. The init function on line 1 adds all type-able objects to the set of possible solutions. The initial possible types of the objects will be the collection of all possible types, called *Universe()*. Only the literal types can be resolved already, like strings, numbers and boolean constants.

Line starts a loop until a fixed-point is reached. In Rascal this method is called `Solve()` and will loop until there are no more changes in the given data. The methods on line 3-4 may alter the values of constraints or solutions.

The first function of the loop on line 3, `deriveMore`, visits all constraints and checks if there can be new constraints extracted based on the latest constraints and solutions. Example of these new constraints are conditional constraints, which could only be added if more information of the conditional was available. Function `propagateSolutions` on line 4 propagates resolved estimates. Here we take the intersection of the known solutions with the given constraints.

5.4 Additional constraints

We've added the option to add context to the type inference process. These are the ability to include type annotations from phpdoc in the constraint analysis, and to include constraints on php built-ins.

Annotations After we gather all the facts from the source code, when constructing an M^3 model, we can add type information from annotations. We match `@return` and `@param` annotations for functions and class methods, and read `@var` annotations for variables and class attributes. In our first analysis we do not include the facts we gathered from the annotations. In the second analysis we do include the facts. This way we can compare the end results.

PHP Built-ins In order to see if we are able resolve more types, we add type information of PHP built-in features. The built-in features are represented by PHP files and can also be parsed and analysed. There files contain information about internal behaviour of PHP. One special thing that we do is that we not add constraints for the empty body of the method, because this would result in an empty return values which is invalid. We also filter the results type resolutions of the php built-in types out of the results we provide in the next chapter.

Chapter 6

Evaluation

In this chapter we evaluate the type inference algorithm. We first describe the setup of the experiment in section 6.1. Next we describe the taken steps to reproduce the analysis in section 6.2. Then we present the results in section 6.3. In section 6.4 we analyse the measured results and explain the differences between the performed analysis.

6.1 Experiment setup

For this experiment we want:

- Source code of programs to perform the type inference on
- Type inference results
- Variations of the type inference
- Expectations and hypotheses

In order to run our analysis we need actual source code, preferably projects or libraries that are used in real world programs. Next we perform the type inference analysis without adding program context. In the next runs we add type hint annotations and php built-in information to the analysis and compare the results. Finally we want to test our hypotheses. All these items are described in more details in the next paragraphs.

Project selection In order to validate this research we tested our type inference algorithm on a selection of the 30 most popular packages of Packagist¹. Packagist is a repository for Composer² projects.

Composer is a dependency manager for PHP. All projects have between 2 and 6 million downloads, so they should be suitable for our research because they are frequently used in live applications. The selection of projects, sorted on total lines of code, is listed in table 6.1. The statistics are generated using phploc³.

¹<https://packagist.org/explore/popular>, July 2014

²<https://getcomposer.org/>, July 2014

³<https://github.com/sebastianbergmann/phploc>, July 2014

Product			Files		Objects			Lines of code	
Vendor	Project*	Version	D ¹	F ²	C ³	I ⁴	T ⁵	Total ↑	Logical
doctrine	lexer	v1.0	2	7	3	0	0	733	128
phpunit	php-timer	1.0.5	5	11	5	0	0	740	117
phpunit	php-text-template	1.2.2	5	11	5	0	0	768	125
doctrine	inflector	v1.0	2	7	3	0	0	853	130
psr-fig	log	1.0.0	3	15	8	2	2	1 039	155
phpunit	php-file-iterator	1.3.4	5	13	7	0	0	1 071	176

*Selection of the 30 most popular packagist packages ordered by LOC, in July 2014.

¹ = Directories, ² = Files, ³ = Classes, ⁴ = Interfaces, ⁵ = Traits

Table 6.1: Statistics of most popular composer projects

Type inference In this research we are interested in resolving types for expressions. As explained in section 5.3, the result of the constraint solving is a set of types for each expression. We group these resulting typesets in two groups, resolved and unresolved types. Resolved type means that we can deduce the typeset to one type. All other types are unresolved. In unresolved we include typesets that are not resolved and have only type `any()`, empty typesets, and multiple types that cannot be reduced to one unique type for an expression.

Annotations We first run the type inference without reading annotations from PHPDocs. As described in section 4.2, we can use type annotations to read types of variables, functions, methods, and class attributes. Our goal is to see if it is possible to resolve more types when we take the type annotations from the PHPDocs into account.

Built-ins After we run the analysis with and without PHPDocs, we run the analysis also with built-in information. Modern editors stub files written in PHP which represent the internal behaviour of constants, variables, functions, classes, methods and class attributes. Using the type information of the built-in language constructs we want to see if we are able to resolve more types.

Expectations The goals of these three analysis is to find out if taking more context into account will increase type resolution. We expect that when we take annotations into account that we are able to resolve more types. When we add PHP built-in information to the analysis we expect to resolve slightly more information. We’ve seen in the research on PHP features usage [HKV13] that many projects use many built-in features. Because there is no type information available by default we expect to be able to resolve more types when adding type annotations of PHP classes, functions, and variables.

Hypotheses Based on our expectations we’ve constructed following null (H_0) and alternative (H_A) hypotheses:

1. Adding type annotations to the analysis
 - $1-H_0$: adding type annotations has no effect on the number of inferred types
 - $1-H_A$: adding type annotations has an effect on the number of inferred types
2. Adding php built-in information to the analysis
 - $2-H_0$: adding php built-in information has no effect on the number of inferred types
 - $2-H_A$: adding php built-in information has an effect on the number of inferred types

6.2 Steps to reproduce

In order to be able to understand and reproduce the analyses we describe the taken steps to get to the results. To reproduce the experiment you need to run the analysis for the three variants and then compare the results.

Running the analysis The source code of the type inference is open source can be found on github⁴.

Clone the repository and open the project in Eclipse with Rascal installed. Analysis results can be generated by running the “lang::php::experiments::mscse2014::mscse2014” Rascal module. Execute `runAll(loc projectLocation)` for all wanted projects of libraries, which results in writing a the analysis results to the file system. The analysis results file contains type information of all expressions of the program. An example of the output of the Rascal console is shown in figure 6.1.

```
rascal> import lang::php::experiments::mscse2014::mscse2014;
rascal> runAll(|home:///location/to/project|);
rascal> ...
rascal> Writing output to file |home:///location/to/project/result_file.txt|.
```

Figure 6.1: Rascal console: running the type inference

Variants of analysis For the next two analysis, with type annotation and PHP built-in information we run the same steps as shown in 6.1, but then we change the configuration. In the type annotation analysis run we set `useAnnotations = true`; and execute `runAll(loc projectLocation)` again for all projects. For the PHP built-in analysis we also set `includePhpInternals = true`; and run all the analysis again. These two variant will write the output to different files.

Compare results When all the projects are analysed you will have result files for all projects. Using the “lang::php::experiments::mscse2014::RascalResultAnalysis” Rascal module you can compare the results of the different projects. The output is a table generated in L^AT_EX with the number of resolved items for the three different analysis. The first table compares the default analysis with the analysis where type annotations are taken into account. The second table shows the difference between the annotations analysis and the analysis adding PHP built-in information. An example of the Rascal output is shown in figure 6.2 This output of the generated tables can be used into your L^AT_EX document.

```
rascal>import lang::php::experiments::mscse2014::RascalResultAnalysis;
rascal>main()
\paddmissingzero
\npfourdigitsep
\begin{table}[H]
... content of the table ...
\end{table}
\npfourdigitnosep
\npnoaddmissingzero
```

Figure 6.2: Rascal console: result analysis

6.3 Results

In this section we present the results of the analysis. We start with the results of the analysis with and without taking PHPDocs into account. Next we show the results of the analysis with and without PHP built-ins typehints. Finally we evaluate the hypotheses.

Annotations The results of the analysis without and without taking type annotations from PHPDocs into account are shown in table 6.2. In this table on the left side we present the name of the project and the total number of expressions that were attempted to be inverted. On the right side of the table we present the percentage of resolved types for the analysis with and without taking the PHPDoc context into account.

⁴<https://github.com/ruudvanderweijde/php-analysis/tree/mscse2014>

Project	Total	Resolved types	
		w/o PHPDoc	w/ PHPDocs
php-timer	14	3	9
log	66	27	27
php-text-template	23	8	10
lexer	42	19	24
inflector	31	3	14
php-file-iterator	63	19	28

Table 6.2: Type inference results, with and without PHPDocs

Built-ins The results of type inference where we take type information of PHP built-in functions and variables into account are shown in table 6.3. In this table we present the analysed projects and the total amount of expressions we attempted to infer. On the right side of the table we present the percentage of resolved types for the analysis with and without built-in information.

Project	Total	Resolved types	
		w/o Built-ins	w/ Built-ins
php-timer	14	9	9
log	66	27	27
php-text-template	23	10	10
lexer	42	24	24
inflector	31	14	14
php-file-iterator	63	28	29

Table 6.3: Results of type usage, with and without PHP built-ins

6.4 Analysis

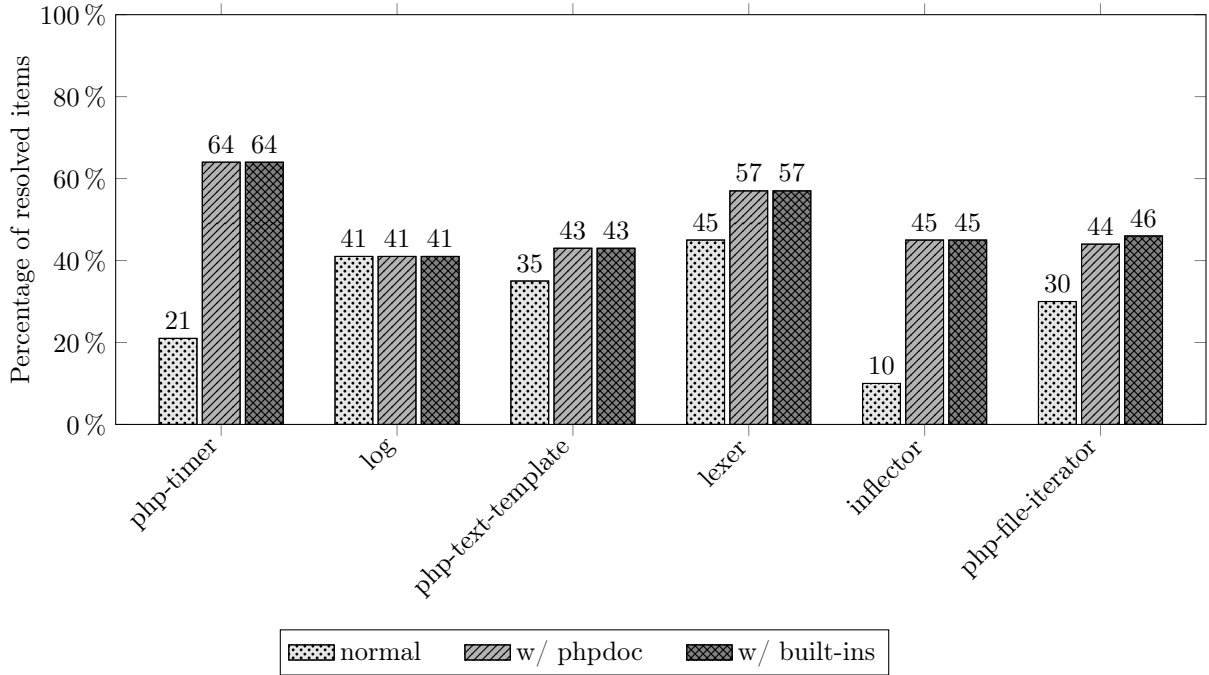


Figure 6.3: Comparison of the type inference results

To get better insight on the difference between the results, we list the percentages of resolved expressions next to each other in the bar chart in figure 6.3. The three bars represent the three variants of the analysis. The bars respectively present the results of the analysis where we do not take extra context into account, the results with using annotations from PHPDocs, and the results where we also take

PHP built-in information into account. The percentages are rounded to the nearest integer to improve readability.

Normal vs PHPDocs In the bar chart we can see that, except for one project, all projects have an increased number of resolved items. There are two projects that show a big difference and three projects that have slightly more resolved items when we take annotations into account. In general we could say that adding annotations aids in having more resolved results.

PHPDocs vs built-ins When we compare the PHPDocs analysis with the built-in analysis we only see a small increase in the number of resolved items for one of the projects. This project is the `php-file-iterator`. All the other projects that we analysed show no difference in the number of resolved item. Because the wide usage of built in functions, classes, and variables we were expecting a bigger difference in the number of resolved items. This might not be the case for these smaller library packages that we used for this analysis. The projects are very small and abstract.

Hypotheses evaluation We constructed the following hypotheses in section 6.1.

1. Adding type annotations to the analysis
 - $1-H_0$: adding type annotations has no effect on the number of inferred types
 - $1-H_A$: adding type annotations has an effect on the number of inferred types
2. Adding php built-in information to the analysis
 - $2-H_0$: adding php built-in information has no effect on the number of inferred types
 - $2-H_A$: adding php built-in information has an effect on the number of inferred types

Hypothesis 1 We've seen that all the analysed projects have a neutral or positive effect on the number of inferred items when adding type annotations. Based on the results of the we can reject the first null hypothesis ($1-H_0$), because there are projects with an increase in the number of inferred types. And therefor we can accept the alternative hypothesis ($1-H_A$): adding type annotations has an effect on the number of inferred types. Almost all projects show a positive effect by adding the type annotations.

Hypothesis 2 : Based on the results of the we can reject the second null hypothesis ($2-H_0$) because we've seen a difference in the number of changed. Therefor we can accept the alternative hypothesis ($2-H_A$): adding php built-in information has an effect on the number of inferred types. This effect is in most cases neutral, but one of the projects shows an increase in the number of inferred types.

Chapter 7

Conclusion

In this section we present the conclusion in section 7.1. Future work is described in section 7.2. Finally we list threats to validity in section 7.3.

7.1 Conclusion

In this work we have presented a type inference implementation written in Rascal. We presented a theoretical constraint based type inference algorithm. This type constraint algorithm is implemented with a PHP parser written in PHP and constraint extractor written in Rascal. To resolve the constraints we also created a constraint solver in Rascal. The results of the constrain solver is a list of code expressions and their possible set of types.

In a small evaluation we have performed the type inference on a few small projects. In this evaluation we found that adding annotations resulted in inferring more types for almost all projects. When we also added type information about PHP built-in classes, functions, and variables, we only saw a small improvement in one of the projects. Hereby we could conclude, although the number of projects is too small to generalise the results, we saw a good improvement when adding type annotations from PHPDocs, but no extra improvement when also adding PHP built-in information.

7.2 Future work

There are various extensions of this research possible. The constraint based type inference could be combined with other analysis, like file include resolutions, dead code elimination, reference/alias analysis, and constant folding/propagation. By including these other analyses it should be possible to get more precise results.

The performance of the constraint solver could be optimised. In our version we did not take any effort to optimise the algorithm. But we have seen that the algorithm is too slow to be able to run on bigger code bases. If the performance of the type inference is improved it is also easier to validate the results of this analysis by adding more projects to the evaluation.

At the moment of writing the implementation, PHP7 was not released yet. The benefit we have seen on adding type annotation is optionally implemented in PHP7. It would be interesting to use the new language constructs and compare the results with PHP5 programs.

7.3 Threats to validity

Because we perform an overapproximation of run time values at compile time, the results of the type inference are not 100% accurate. With the dynamic nature of the programming language and the usage of dynamic features it is a challenge to be as accurate as possible.

Missing constraints The constraints in this research are supposed to model the dynamic semantics of PHP soundly, and they only be manually checked by reading the PHP specs or documentation¹. By adding a constraint we limit the number of possible types for that expression. It could be possible that we made a mistake by missing a constraint or adding a wrong constraint. This could result in false type inference results.

Generalisability Because the small sample size in our evaluation, we are not able to generalise the results of this research. The used projects are slim and do not represent the usage of other PHP program. In order to gain more precise results we should increase our sample size. This can be achieved by running analysis on more software projects in the future.

Object references One thing we do not take into account are the side-effects that can be caused by passing arguments as reference. When an referenced argument is modified, the referenced variable is also modified. Full-fledged alias analysis can be added to mitigate this effect.

PHPDoc correctness In this research we also assume that the provided type annotations in the PHPDocs are correct. These annotations are ignored during program execution and therefore could be incorrect without noticing. Using PHP7 new type hint language constructs could help to ensure that we extract the correct constraints for the source code.

¹<http://php.net/docs.php>, December 2016

Glossary

PSR

PHP Standard Recommendation (PSR) is project which provides rules for commonalities between PHP projects. Autoloading, coding style guide, logging, and HTTP Message interface are the first few accepted standards. The PHPDoc standard describes how and what to use in doc blocks and is currently in draft phase. See <http://www.php-fig.org/> for more information.

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.

static code analysis

In a static code analysis the targeted program is analysed without executing the program.

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

- [Age95] Ole Agesen. “The Cartesian Product Algorithm: Simple and Precise Type Inference Of Parametric Polymorphism”. In: *Proceedings of the 9th European Conference on Object-Oriented Programming*. ECOOP ’95. London, UK, UK: Springer-Verlag, 1995, pp. 2–26. ISBN: 3-540-60160-0. URL: <http://dl.acm.org/citation.cfm?id=646153.679533>.
- [ASU86] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 1986. ISBN: 0-201-10088-6.
- [Bas+15] Bas Basten et al. “ M^3 : a General Model for Code Analytics in Rascal”. In: *Proceedings of the first International Workshop on Software Analytics, SWAN*. 2015.
- [Big10] Paul Biggar. “Design and Implementation of an Ahead-of-Time Compiler for PHP”. In: (2010).
- [Cam07] Patrick Camphuijsen. “Soft typing and analyses on PHP programs”. In: (2007).
- [CHH09] Patrick Camphuijsen, Jurriaan Hage, and Stefan Holdermans. “Soft Typing PHP with PHP-validator”. In: (2009).
- [DM82] Luis Damas and Robin Milner. “Principal Type-schemes for Functional Programs”. In: *Proceedings of the 9th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. POPL ’82. Albuquerque, New Mexico: ACM, 1982, pp. 207–212. ISBN: 0-89791-065-6. DOI: [10.1145/582153.582176](https://doi.org/10.1145/582153.582176). URL: <http://doi.acm.org.proxy.uba.uva.nl:2048/10.1145/582153.582176>.
- [Hin69] R. Hindley. “The Principal Type-Scheme of an Object in Combinatory Logic”. English. In: *Transactions of the American Mathematical Society* 146 (1969), pages. ISSN: 00029947. URL: <http://www.jstor.org/stable/1995158>.
- [HKV13] Mark Hills, Paul Klint, and Jurgen J. Vinju. “An Empirical Study of PHP feature usage: a static analysis perspective”. In: *ISSTA*. Ed. by Mauro Pezzè and Mark Harman. ACM, 2013, pp. 325–335.
- [HKV14] Mark Hills, Paul Klint, and Jurgen J. Vinju. “Static, Lightweight Includes Resolution for PHP”. In: *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering*. ASE ’14. Vasteras, Sweden: ACM, 2014, pp. 503–514. ISBN: 978-1-4503-3013-8. DOI: [10.1145/2642937.2643017](https://doi.org/10.1145/2642937.2643017). URL: <http://doi.acm.org/10.1145/2642937.2643017>.
- [Izm+13] Anastasia Izmaylova et al. “M3: An Open Model for Measuring Code Artifacts”. In: *CoRR* abs/1312.1188 (2013).
- [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.
- [KSV09] Paul Klint, Tijs van der Storm, and Jurgen J. Vinju. “RASCAL: A Domain Specific Language for Source Code Analysis and Manipulation”. In: *SCAM*. 2009, pp. 168–177.
- [Mil78] Robin Milner. “A theory of type polymorphism in programming”. In: *Journal of Computer and System Sciences* 17.3 (1978), pp. 348–375. ISSN: 0022-0000. DOI: [10.1016/0022-0000\(78\)90014-4](https://doi.org/10.1016/0022-0000(78)90014-4). URL: <http://www.sciencedirect.com/science/article/pii/0022000078900144>.
- [NNH99] Flemming Nielson, Hanne R. Nielson, and Chris Hankin. *Principles of Program Analysis*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 1999. ISBN: 3540654100.
- [PS91] Jens Palsberg and Michael I. Schwartzbach. “Object-oriented Type Inference”. In: *SIGPLAN Not.* 26.11 (Nov. 1991), pp. 146–161. ISSN: 0362-1340. DOI: [10.1145/118014.117965](https://doi.org/10.1145/118014.117965). URL: <http://doi.acm.org.proxy.uba.uva.nl:2048/10.1145/118014.117965>.
- [PS94] Jens Palsberg and Michael I. Schwartzbach. *Object-oriented type systems*. Wiley professional computing. Wiley, 1994, pp. I–VIII, 1–180. ISBN: 978-0-471-94128-6.
- [Shi88] O. Shivers. “Control Flow Analysis in Scheme”. In: *SIGPLAN Not.* 23.7 (June 1988), pp. 164–174. ISSN: 0362-1340. DOI: [10.1145/960116.54007](https://doi.org/10.1145/960116.54007). URL: <http://doi.acm.org.proxy.uba.uva.nl:2048/10.1145/960116.54007>.
- [Shi91] Olin Shivers. *Control-Flow Analysis of Higher-Order Languages*. Tech. rep. 1991.
- [VH15] Henk Erik Van der Hoek and Jurriaan Hage. “Object-sensitive Type Analysis of PHP”. In: *Proceedings of the 2015 Workshop on Partial Evaluation and Program Manipulation*. PEPM ’15. Mumbai, India: ACM, 2015, pp. 9–20. ISBN: 978-1-4503-3297-2. DOI: [10.1145/2678015.2682535](https://doi.org/10.1145/2678015.2682535). URL: <http://doi.acm.org.proxy.uba.uva.nl:2048/10.1145/2678015.2682535>.
- [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>.