

Macquarie University
Department of Computing

PROJECT REPORT

Static Type Analysis of Dynamically Typed Programming Language

Author: Stepan Sindelar

Student ID: 43600220

Supervisor of the project: Matthew Roberts

Study programme: exchange student

Contents

1	Introduction	3
1.1	The Problem	3
1.2	Thesis structure	3
2	Analysing PHP Code	5
2.1	PHP semantics	5
2.1.1	Local Variables	5
2.1.2	Global and Local Scope	5
2.1.3	Closures	6
2.1.4	References	6
2.1.5	Arrays	7
2.1.6	Interesting Control Flow Structures	7
2.1.7	Conditional Declarations	7
2.1.8	Auto-loading	8
2.1.9	PHPDoc Annotations	8
2.2	Static Code Analysis	8
2.2.1	Terminology	9
2.2.2	Data Flow Analysis	10
2.2.3	Intraprocedural Analysis	12
2.3	Control Flow for Phalanger Approach	12
3	Existing Software	13
3.1	Phantm[9]	13
3.2	HipHop Type Inference for Hack	13
3.3	Weverca: Web Verification Tool[10]	13
4	Implementation	14
4.1	Implementation Specific Constraints	14
4.2	Overall Design	15
4.3	Control Flow Graph	15
4.4	Data Flow Analysis	16
4.5	Tables	16
4.6	Aliasing and Constant Propagation Analysis	16
4.7	Type Analysis	16
4.8	The Analyser pipeline and AnalysisDriver	16
5	Results	17
5.1	PhpUnit	17
5.2	Zend Framework	17
5.3	Nette	17
5.4	WordPress	17
5.5	Drupal	17
6	Conslusion	18
6.1	Future Work	18

Bibliography	19
List of Tables	20
Attachments	21

1. Introduction

1.1 The Problem

Three out of top ten programming languages in TIOBE index[1] are dynamically typed languages. One of the reasons for their popularity is that they are usually easier to use and suitable for fast prototyping. But at the same time the possibility to omit type information, which might be helpful during the early stages of a software project, can lead to more error prone code, and eventually to problems in later phases of the development and maintenance. Dynamic typing is also challenging for the compilers or interpreters designers. With the type information, a compiler is usually able to emit more efficient code.

Programmers are aware of the possible problems with the maintenance of dynamically typed code and they often document the type information in documentation comments. However, the correspondence of the documentation and the actual code is not checked, and moreover the compiler usually does not take any advantage of having type hints in the comments.

The PHP programming language is one of the mentioned popular dynamic languages and Phalanger [2] is an implementation of a PHP compiler that compiles PHP code into the .NET intermediate code. Phalanger was developed at the Department of Software Engineering of the Charles University in Prague. A part of the Phalanger project is also an implementation of PHP tools for Visual Studio.

Because of its dynamic nature, PHP code is more difficult to analyse than code written in a statically typed language, especially if we want the analysis to be reasonably fast so that it can be used in everyday development. There is an ongoing research of the static analysis methods for many different families of programming languages, including dynamic languages. The problem this project is addressing is to adapt and apply those methods on a real world and widely used programming language PHP. The result is a library that is capable of performing a static analysis of PHP code and can be integrated into the Phalanger project. The library should allow to plug in any kind of analysis, for example constant propagation. However, the main goal is to provide a type analysis in order to discover possible type related errors and mismatches with the type information in the documentation comments. Furthermore, in the future the library can be integrated into the Phalanger as a middle-end to provide optimizations for the compiler.

This project and the library has a code name Control Flow for Phalanger and we refer to it using this name in the following text.

1.2 Thesis structure

The following chapter describes the challenges connected with analysing source code written in the PHP programming language, approaches to static analysis of source code in general and our adoption of those to PHP. In chapter 3, we

briefly discuss existing software of this kind. Chapter 4 provides more detailed description of our implementation. The analysis has been evaluated on several middle to large sized open source PHP projects and the results are presented in chapter 5.

2. Analysing PHP Code

The PHP programming language first appeared in 1995[3]. Over the years the language has evolved and so have the ways how programmers were using it. This project focuses on PHP version 5.5¹ and the aim for the analysis is to work well on PHP source code written in an object oriented style, using modern PHP patterns and idioms that are described later in this text. The analysis, however, should provide correct results for any valid PHP code of any PHP version. We do not focus only on websites, but also on PHP libraries and frameworks that by themselves do not contain any PHP files that produce HTML or any other output for the user.

2.1 PHP semantics

This section describes some important parts of the semantics of the PHP programming language, especially those that represent a challenge for static analysis.

In PHP, local or global variables, object fields and function or method parameters are dynamically typed, which means that they can hold values of completely different types at different times of the execution.

2.1.1 Local Variables

Local variables in PHP do not need to be declared explicitly. Instead the first usage of a variable is also its declaration. If a variable's value is used before the variable got any value assigned, then the interpreter generates a notice, however the execution continues and value `null` is used instead. A variable can get a value assigned to it when it appears on a left hand side of an assignment or when a reference to that variable is created, in which case it gets value `null`, but no notice is generated. References are discussed in one of the following subsections.

The scope of a local variable is always its parent function not the parent code block as in other languages like C or Java. So in the following example, the usage of variable `$y` at the end of the function can generate uninitialized variable notice, however, if `$x` was equal to 3, `$y` will have a value although it was declared in the nested code block.

```
function foo($x) {  
    if ($x == 3) { $y = 2; }  
    echo $y; // uninitialized variable if x != 3  
}
```

2.1.2 Global and Local Scope

PHP distinguishes two scopes for variables: global scope and local scope. Local scope is a scope of local variables within a user defined routine. Variables that are declared in global scope, that is outside of a user defined routine, are available

¹From this point, if the PHP version is not stated explicitly, it is implicitly 5.5.

anywhere in global scope and are called global variables. Global variables are also available in user defined routines as long as they are imported into the routine's scope using the keyword `global`.

```
$g1 = 1; // global variables g1 and g2
$g2 = 2;
function foo() {
    global $g1;
    echo $g1; // prints the value of global variable g1
    $g2 = 4; // sets the value of local variable g2,
    // because global variable g2 was not imported,
    // its value does not change
}
```

Global variables can be imported and used in any user defined routine. This means that even if we know some type information about a global variable's value at some point in the analysed code (e.g. straight after assignment to that variable), any time another user defined routine is invoked, we have to take into account that the other routine can change the value of the global variable even if we do not pass the global variable to the invoked routine as an argument passed by reference.

2.1.3 Closures

PHP also supports anonymous functions. An anonymous function has its own scope as any other function and its local variables are not visible to the scope where it was declared. Variables from the parent scope are available in the closure scope only if they are explicitly imported to its scope and they can be captured by value or by reference. Only the later represents a challenge for the analysis, because any code that can access the closure can invoke it and thus change the values of variables imported to the closure's scope by reference. By invoking a closure, we can influence the values of variables in a completely different and otherwise inaccessible scope.

2.1.4 References

References in PHP are similar, but not same, as pointers in the C programming language. PHP has a special operator `=&` (assign by reference) that turns the variable on the left hand side into a reference to the variable on right hand side. For example `$a=&$b`, after this, any assignment to `$a` will in fact change the value of `$b` and wherever the value of `$a` is used (e.g. in an expression), the value of `$b` is used instead.

The variable where the reference is pointing to is determined in a transitive fashion, which means that if we assign by reference another reference, the new reference will point to where the other reference was pointing to, but the intermediate link is lost. The following example illustrates this.

```
$a =& $b; // a points to b
$c =& $a; // c points to where a points, that is b
$a =& $d; // a points to d, but c still points to b
```

2.1.5 Arrays

Arrays in PHP do not have to be homogenous and they can be indexed by either integers or strings. In fact, PHP arrays are hash maps rather than arrays in the usual sense and that is also how they are implemented internally.

String indexed heterogenous arrays are often used as flexible ad-hoc structured data type. Instead of defining a class with required fields, one can use what would be a field name as an index into an array. Such arrays are usually indexed only with finite number of constant string values.

In this light, it is no surprise that using the subscribe operator `[]` with string index on an object instance will access the field with the same name as the index value.

2.1.6 Interesting Control Flow Structures

The `break` and `continue` statements with optional numeric argument are supported in PHP in a similar way as in other standard imperative programming languages. There are however important differences to be noted.

Firstly, The numeric argument can be an arbitrary expression in some of the older versions of PHP, in which case we cannot statically determine the target of the jump for the control flow resolution.

Secondly, the `switch` statement is considered a loop for the purposes of both `break` and `continue`. The semantics of `break` is intuitive. One of the meaningful use cases is to break a loop from within a `switch` by using `break 2`. The semantics of `continue` statement is perhaps not so intuitive: within a `switch` it works the same way as `break`.

2.1.7 Conditional Declarations

User defined functions, classes, etc. are declared in a global scope in PHP, that is a scope where one can as well place any arbitrary code. Therefore a the declaration can be wrapped in any control structure. It is not allowed to redeclare once declared symbol, however.

A typical use case is to dynamically import a file that may provide some functions and then check using `function_exists` whether the functions where indeed declared and if not, provide default implementation. This is a pre object oriented way of doing overriding and is usually not to be found in modern projects. Nonetheless, popular system WordPress[4] still relies on this pattern in parts of its code base.

Although the mentioned pattern could be deemed as reasonable and useful. This feature allows to write very tricky code as in the following example that may or may not crash on fatal errors “Cannot redeclare foo()” or “Call to undefined function foo()” depending upon the user input.

```
while ($_POST['a'] != 3) {  
    function foo() { return 5; }  
    $_POST['a'] = $_POST['b'];  
}  
echo foo();
```


2.1.8 Auto-loading

Historically, in PHP in order to reference any symbol from a different file, one had to import that file explicitly. Newer versions of PHP allow to define custom auto-loading routine that is invoked by the runtime at any time an undefined class is referenced. The auto-loading routine is then responsible for importing the file(s) that contain the code of the required class.

Auto-loading routine can use arbitrary logic to determine what file(s) to import, in fact, it can execute arbitrary code. Typical pattern used for example in Zend Framework[5] before namespaces were introduced to PHP is to have a file per class and use class names in form of `CodeFolder1_SubFolderName_FileName` for a class placed in file `CodeFolder1\SubFolderName\FileName.php`.

2.1.9 PHPDoc Annotations

Although not part of the official PHP syntax, there is widely recognized format for documentation comments of JavaDoc style called PHPDoc. PHPDoc comments may contain type information that cannot be expressed using PHP syntax. For example, PHP allows “type hints” for routine parameters, but only for class types, not for primitive types like `int`. However, primitive type expectations can be included in the documentation comments. The important difference is that PHP will throw an exception at runtime if a routine is invoked with a parameter of different type than what its “type hint” is. The documentation comments, on the other hand, are of course not checked by the runtime.

The PHPDoc defines a fairly advanced syntax for expressing type information. It allows to include multiple primitive and class types, description of homogenous and heterogenous arrays as well multidimensional arrays, and some constants like `false`. For example, in the following code the documentation comment tells us that function `foo` can return either `null`, or `false` (but should never return `true`), or an array of `integer` values.

```
/**
 * @return null|false|int[]
 */
function foo() { ...
```

2.2 Static Code Analysis

Static analysis of source code is an analysis that is performed without executing the code. This means that we do not have to have a web server for example in order to analyse code of a web application. We can also guarantee some properties that would not be possible to guarantee if we executed the code. Namely the halting property and upper bounds on time and space complexity. Arbitrary code may not halt if executed, but static analysis of such code can still halt and give us some results.

Static analysis can be used to get possible types of an expression in a dynamically typed language, to find out expressions that have constant value through constant propagation and many other problems. Static analyses usually do not

give accurate solution, but its approximation, which can be an over-approximation or an under-approximation and it is up to the designer and user of the analysis which one is acceptable for his or her² purposes.

2.2.1 Terminology

Static analyses are usually described in the form of inference rules. An example of an inference rule can be “if the types of expressions e_1 and e_2 are integers, then the type of expression $e_1 + e_2$ is integer”. Those rules can be more formally described with the following notation

$$\frac{\vdash e_1 : Integer \quad \vdash e_2 : Integer}{\vdash e_1 + e_2 : Integer}$$

where above the horizontal line we have hypothesis and below is the conclusion. The exact notation is not important we will be using it intuitively to illustrate the ideas that we describe.

Flow Sensitivity. The example inference rule is not valid, if we admit that evaluation of expression e_1 can influence the type of expression e_2 or vice versa. In such case, the ordering of the expressions is important, but not captured in the inference rule. Therefore this rule is *flow insensitive*. If we make the hypotheses more complex to capture the ordering it will be *flow sensitive*.

Path Sensitivity. If we admit conditional control flow statements, like if-then-else, we can have more possible paths through the program. In our example, let us say that there is an if statement before the expression e_1 and the expression e_1 can evaluate to different type if the then branch of the if-then-else statement is taken than if the else branch is taken. This is illustrated in the following code listing. If the inference rules do not model this, like our example rule, we say that the inference system is *path insensitive*, otherwise we say that it is *path sensitive*.

```
if (...) $x = 3;
else $x = 'string';
$e_1 = $x;
$e_2 = 4;
$e_1 + $e_2;
```

Abstraction. Another example of problem that can be partly solved with static analysis is the sing of integral variables. We can have inference rules of the following form.

$$\frac{\vdash v_1 : -10(sign : \ominus) \quad \vdash v_2 : 3(sign : \oplus)}{\vdash v_1 + v_2 : -7(sign : \ominus)}$$

However, the implementation would not be very efficient and we sometimes do not have the full information about variables values, but in some cases we can deduce another useful piece of information by other means. For example, variable of type `unsigned integer` will always be positive, we can count on that even if

²“His” or “he” should be read as “his or her” or “he or she” through the rest of the text.

we do not know the actual value. What we can do is to abstract the possible integral values with set $\{0, \ominus, \oplus\}$ with the following meanings

- \ominus represents all negative integers,
- \oplus represents all positive integers,
- 0 represents zero,

and rewrite the inference rules as follows:

$$\frac{\vdash v_1 : \ominus \quad \vdash v_2 : \ominus}{\vdash v_1 + v_2 : \ominus}$$

Nonetheless, there is another problem. What to do when we have \ominus and \oplus in the hypothesis.

$$\frac{\vdash v_1 : \ominus \quad \vdash v_2 : \oplus}{\vdash v_1 + v_2 : ?}$$

The solution is to extend the domain so that it is closed under all operations. We add another element to our set:

- \top represents an unknown value (either positive, negative, or zero).

Then the rule will be:

$$\frac{\vdash v_1 : \ominus \quad \vdash v_2 : \oplus}{\vdash v_1 + v_2 : \top}$$

And for example another rule dealing with \top in hypothesis:

$$\frac{\vdash v_1 : \ominus \quad \vdash v_2 : \top}{\vdash v_1 + v_2 : \top}$$

2.2.2 Data Flow Analysis

Data Flow Analysis (DFA) is a static analysis framework for compiler optimizations and verification that scales to large code bases [6], [7].

Control Flow Graph. DFA is typically performed on a control flow graph, although there are exist approaches to DFA without explicit control flow graph construction [8].

Control flow graph nodes, also called basic blocks, are program statements that are always sequentially executed. Directed edges represent the control flow between basic blocks, i.e. jumps in the control flow due to conditionals, goto statements or any other statements that can change the flow of the program. Control flow graphs usually contain two special nodes: entry node and exit node. Entry node does not have any incoming edges and all the paths lead to the exit node. An example of a control flow graph is given in figure 2.1.

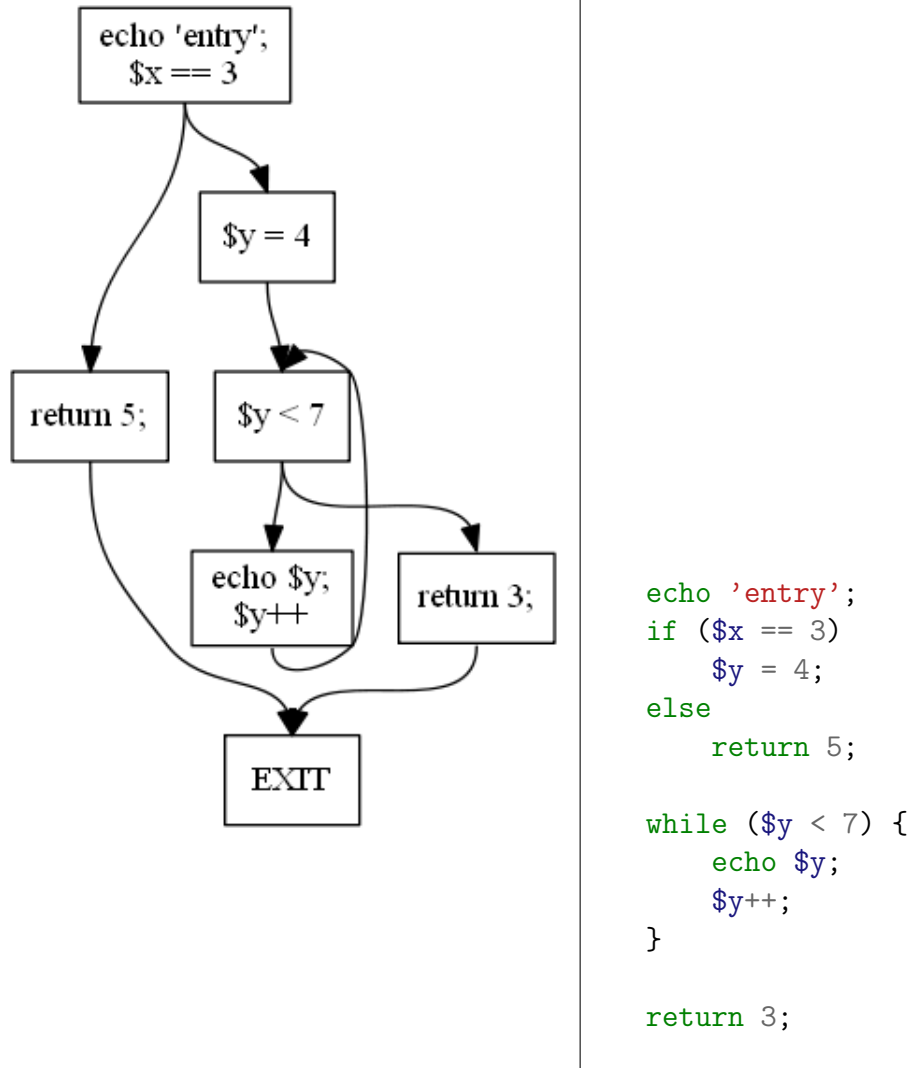


Table 2.1: Control flow graph

Lattices. The reasoning behind correctness and termination of DFA is based on an algebraic structures called lattices. A lattice is a partially ordered set in which every two elements have a least upper bound, called supremum, and a greatest lower bound, called infimum.

Bounded lattice is a lattice that has a greatest element and a least element, usually denoted as \top and \perp . A bounded lattice is depicted in figure 2.1.

What is important about lattices for DFA analysis is that if we have a bounded lattice (S, \leq_s) and a function $f : S \rightarrow S$ that is monotonous, i.e. $\forall a \in S : a \leq_s f(a)$, then $\forall x \in S \exists k \in \mathbb{N}$, such that $f^k(x) = f^{(k+1)}(x)$. $f^k(x)$ is called a fixpoint.

Intuitively, f has to have a fixpoint because for every argument y , it must either return y itself, but then y is a fixpoint, or it returns another element that is strictly greater than y , but this cannot go on forever, because eventually f will be given \top for which it does not have any other option but returning \top and we have a fixpoint again.

The following few paragraphs will finish the description of Data Flow Analysis.

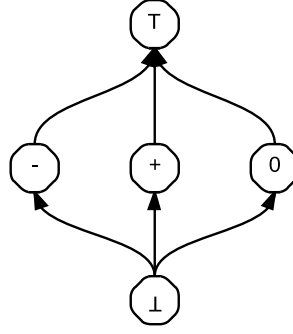


Figure 2.1: Bounded lattice with 5 elements.

2.2.3 Intraprocedural Analysis

So far we have been discussing an analysis of a single function or method³. However, if we want to analyse whole program or a library, the interaction between the routines can be taken into account to make it more precise.

Context Sensitive Intraprocedural Analysis. The most precise solution would be to take into account the calling context when analysing a function. In different contexts, the function can be, for example, given different parameters values which may then influence the result of the analysis. More precise result for a specific call site context can be propagated to that call site, yielding another gain in precision when analysing the function that realises the call.

The following few paragraphs will discuss feasibility of Context Sensitive Intraprocedural Analysis, because call sites are not always known, practical consequences and usual approaches to make Context Sensitive Intraprocedural Analysis more scalable.

2.3 Control Flow for Phalanger Approach

Description of the analysis used in our case using the terminology built up in the previous section. It will be more abstract description, without technical details about implementation.

³We will use term routine to designate a function, static or instance method

3. Existing Software

3.1 Phantm[9]

3.2 HipHop Type Inference for Hack

3.3 Weverca: Web Verification Tool[10]

4. Implementation

4.1 Implementation Specific Constraints

One of the requirements was that the project should be implemented in the context of the Phalanger project. It should be ready to be plugged in between the Phalanger’s front-end and back-end and it should also provide public interface useful for the Phalanger PHP Visual Studio tools.

Abstract Syntax Tree. Phalanger front-end parses PHP code into an Abstract Syntax Tree (AST) [7] structure. This structure is then traversed by the back-end using the visitor design pattern [11]. Phalanger does not use any other intermediate representation than AST.

In order to reduce the memory consumption and provide better modularity, Phalanger code went through small architecture refactoring before this project was started. The classes representing the AST nodes originally contained the code and data needed for emitting the corresponding Microsoft Intermediate Language opcodes. This, however, represents a coupling between the front-end and back-end and the front-end cannot be used just on its own. In the new version, the classes representing AST nodes are capable of storing additional attributes in an extensible way and the back-end have been rewritten to be an ordinary AST visitor that uses the extensible AST nodes attributes. Additional attributes provide a way to annotate AST nodes with any additional information, which in our case will be the results of the analyses.

Integrated Development Environment Integration. The PHP Tools for Visual Studio use Phalanger front-end in order to parse PHP code into AST and then the AST is again traversed to provide code completion and other features. All the AST nodes hold necessary pieces of information, e.g. the position in the source file or documentation comments.

The aim of this project is to provide the results of type analysis and other analyses to the integrated development environment so that the possible errors and warnings can be visualised, e.g. underlined.

The longer term aim of this project, not in the scope of this thesis, is to replace the existing algorithms for code completion, “jump to definition” and “find usages” features. Because with a dynamic language like PHP, it is not trivial to find all the usages of e.g. a class or determine a definition of e.g. a field accessed on some local variable. In order to provide more precise results, type analysis is needed.

One of the challenging parts of integrated development environment integration is also dealing with incomplete code that is being typed in by the user. Therefore one of the requirements was also that the analysis should be capable of performing an ad-hoc re-analysis of once analysed code with a new statement added. This ad-hoc re-analysis should be, if possible, more effective than doing the whole analysis again.

4.2 Overall Design

The project can be divided into several parts:

- Control Flow Graph construction,
- Data Flow Analysis framework,
- Code Tables maintenance,
- The analyses:
 - Dead Code Elimination,
 - Aliasing Analysis,
 - Constant Propagation,
 - Type Analysis.

The Data Flow Analysis is performed on a Control Flow Graph. However, the Data Flow Analysis is an abstract algorithm and in order to implement a concrete Data Flow Analysis, one must provide implementation of required missing operations.

An analysis is a concrete Data Flow Analysis or any other analysis, because the analysis can be performed, for example, only on the Control Flow Graph (Dead Code Elimination) or only on the function body (Aliasing Analysis) without using the Data Flow Analysis algorithm. The results of an analysis can be annotations added into the AST node objects, or annotations of the Control Flow Graph nodes, which support extensible attributes the way AST nodes do.

The results of an analysis can also be pushed into the Code Tables. Code Tables gather relevant information about code elements like classes, functions and others. This information is mainly type related information, for example, a return type of a function. Code Tables can be also queried by an analysis, for example, when the Type Analysis reaches a function call and needs to find out what types the given function may return.

The whole project is designed as a class library with many extension points where some of the functionality can be used independently. However, in order to provide better usability, it also contains a Facade class **AnalysisDriver** that plugs in together all the necessary objects and provides a simple interface to perform a defined analysis of a file or a given piece of code.

4.3 Control Flow Graph

Description of the Control Flow graph construction. Discussion of the choice of intermediate representation: in fact, why we do not use any intermediate representation, because we wanted to keep simple correspondence between the AST nodes and analyses results and because the back-end emits MSIL code from AST nodes.

4.4 Data Flow Analysis

Architecture of the generic algorithm implementation using generics in C# in a way that the graph the analysis is performed upon can be arbitrary graph (not only CFG), which can be useful if in the future some of the analyses will be performed, for example, on a definition-use graph.

4.5 Tables

Intra-procedural dependencies handling.

4.6 Aliasing and Constant Propagation Analysis

4.7 Type Analysis

- *Type Information representation*
- *Data Flow representation*
- *AST annotations*

4.8 The Analyser pipeline and AnalysisDriver

Description of the high level public interface and the phases of the full analysis process.

5. Results

The project has been evaluated on the following open source PHP frameworks and websites.

- PHPUnit: a port of JUnit unit testing framework for PHP.
- Zend Framework: popular general purpose PHP framework.
- Nette: another popular PHP framework for building websites.
- WordPress: one of the most popular content management systems.
- Drupal: another popular content management system.

An evaluation always started with downloading a git repository with the latest source code of given PHP project. Then the analysis was run and all the discovered errors were collected and classified. Actual errors were rectified and recorded as commits in the git repository.

5.1 PHPUnit

5.2 Zend Framework

5.3 Nette

5.4 WordPress

5.5 Drupal

6. Conclusion

6.1 Future Work

Integer interval analysis, integration with the compiler back-end.

Bibliography

- [1] “Tiobe index for march 2014.” <http://www.tiobe.com/index.php/content/paperinfo/tpci/index.html>. Accessed: 2014-09-03.
- [2] J. Benda, T. Matousek, and L. Prosek, “Phalanger: Compiling and running php applications on the microsoft .net platform,” *.NET Technologies 2006*, 2006.
- [3] “Php: History of php.” <http://www.php.net/manual/en/history.php.php>. Accessed: 2014-02-04.
- [4] “Wordpress is blog tool, publishing platform, and cms.” <http://wordpress.org/>. Accessed: 2014-07-05.
- [5] “Zend framework.” <http://framework.zend.com/>. Accessed: 2014-07-05.
- [6] F. Nielson, H. R. Nielson, and C. Hankin, *Principles of program analysis*. Springer, 1999.
- [7] A. V. Aho, R. Sethi, and J. D. Ullman, *Compilers: Principles, Techniques, and Tools*. Addison Wesley, 1985.
- [8] M. Mohnen, “A graph—free approach to data—flow analysis,” in *Compiler Construction*, pp. 46–61, Springer, 2002.
- [9] E. Kneuss, P. Suter, and V. Kuncak, “On using static analysis to detect type errors in php applications,” tech. rep., 2010.
- [10] D. Hauzar and J. Kofron, “Hunting bugs inside web applications,”
- [11] E. Gamma, R. Helm, R. Johnson, and J. Vlissides, *Design patterns: elements of reusable object-oriented software*. Pearson Education, 1994.

List of Tables

Attachments