



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

BACHELOR THESIS

Petr Houška

**Compilation of a dynamic language
Generators into MSIL**

Department of Software Engineering

Supervisor of the bachelor thesis: Mgr. Jakub Míšek

Study programme: Computer Science

Study branch: General Computer Science

Prague 2017

I declare that I carried out this bachelor thesis independently, and only with the cited sources, literature and other professional sources.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date

signature of the author

I would like to thank my advisor Mgr. Jakub Míšek for his valuable advice and guidance he has given me for this thesis. I would also like to thank him for starting the Peachpie project in the first place and bringing me on despite the disadvantages mentoring a student who is writing a bachelor thesis inherently brings.

This thesis would also not be possible without the endless support of my parents, friends, and classmates and the endless encouragement of my girlfriend. All of them and practically everyone else who I had the pleasure to meet during my studies deserve an acknowledgment.

Title: Compilation of a dynamic language Generators into MSIL

Author: Petr Houška

Department: Department of Software Engineering

Supervisor: Mgr. Jakub Míšek, Department of Software Engineering

Abstract: The goal of this thesis is to design and implement support for generators within the Peachpie framework, a PHP to CIL compiler. Generators are the simplest form of methods that resume from the same state in which they returned earlier when called repeatedly. The reference PHP interpreter Zend engine supports generators natively. Due to that fact generators in PHP support a number of features not usually seen in other languages. CIL on the other hand does not have a native support for generators. Therefore, languages build on top of CIL (e.g. C#, F#) have to implement them by other means such as by rewriting the original generator methods into state machines. In this thesis we will design and implement support for generators through semantic tree transformations. All that with the intention of keeping the maximum possible compatibility with reference PHP generators. We will also make a comparison to generators in C# whose main implementation also uses CIL as a backend.

Keywords: compiler php msil .net generators roslyn peachpie

Contents

Introduction	3
Thesis structure	4
1 Generators	5
1.1 Iterators	5
1.2 Generators universally	6
1.3 Generators in PHP	7
1.4 Generators in other languages	8
2 .NET platform	10
2.1 Common intermediate language	10
2.1.1 Evaluation stack	11
2.1.2 Exception handling	11
3 Peachpie project	13
3.1 Peachpie architecture	13
3.2 Peachpie compiler	13
3.3 Semantic graph	14
3.3.1 Statements and expressions	14
3.3.2 Graph structure	15
3.3.3 Graph creation	16
3.4 CIL emit phase	17
3.4.1 Code generator	17
3.4.2 Emit	17
3.4.3 Generate methods' invariants	18
3.5 Peachpie runtime library	19
4 Generators implementation in other platforms	20
4.1 CSharp and Roslyn	20
4.1.1 Iterator object and generator methods	20
4.1.2 Rewriter	20
4.1.3 Local variables parameters	20
4.1.4 Execution position	20
4.2 PHP and Zend Engine	20
5 Generators in Peachpie	21
5.1 Basic generators implementation	22
5.1.1 Iterator object	22
5.1.2 Next method implementation and local variables	22
5.1.3 Accessibility of fields on the Generator type	22
5.1.4 Context handling	22
5.1.5 Rewriter	22
5.1.6 Bound yield expression	22
5.1.7 Start block	22
5.1.8 Method symbol	22
5.2 Yield as an expression - theory	22

5.2.1	Possible approaches	22
5.2.2	Branch capture & yield splitting	22
5.2.3	Semantic tree transformation	22
5.2.4	Short circuit evaluation	22
5.3	Yield as an expression - implementation	22
5.3.1	Binding multiple elements	22
5.3.2	Capturing branches with yields	22
5.3.3	Correctness of modified capturing algorithm	22
5.3.4	Creating and keeping the pre-bound graph	22
5.3.5	Path between the root and yields	22
5.3.6	Conditioned branches	22
5.3.7	Implementation remarks	22
5.4	Yield in exception handling blocks	22
5.4.1	Yields and exception handling blocks in PHP	22
5.4.2	Solution in Peachpie	22
5.5	Future work	22
Conclusion		23
Attachments		24
	Compilation	24
	Structure	24
	Manual testing	24
	Automatic testing	25
Bibliography		26
List of Figures		28
List of Abbreviations		29

Introduction

Despite a slow decline in recent years [TIOBE, 2017], PHP is still one of the main languages used for a server side programming on the web [Stack Overflow, 2017]. Its only two relevant implementations are the reference and almost exclusively used Zend engine¹ and slowly emerging HHVM by Facebook². Both of them are standalone virtual machines and neither of them supports easy interfacing with the outside world. Hence, it is quite difficult to share code between a web backend and, for example, a mobile or traditional desktop application.

Fortunately, there is a solution in the form of a Peachpie project³ that is being researched at the Charles University. The project aims to provide a compiler from PHP to “.NET bytecode” CIL⁴ and a reimplement of PHP base class library, thus creating a bridge between PHP and the whole .NET ecosystem. Due to it being a full compiler that takes PHP sources and spits out .NET assemblies indistinguishable from those created by other .NET languages compilers (e.g C#, F# or IronPython) it provides both ways interoperability. It enables both calling normal unmodified .NET libraries from PHP and vice versa. Also, thanks to an extensive compile-time type analysis and proven .NET just in time compiler (RiuJIT) it achieves better performance than reference Zend engine in certain operations [Míšek, 2017, Fistein and Míšek, 2016 - 2017].

PHP, like many other modern languages, has a first class support for generators. Simply put, generators are methods that resume computation from the very place and with the same state they returned at previously when called repeatedly. They are usually used for generating large sequences of data lazily, hence the name generators. Since the execution state gets saved automatically on the special pause and return places (usually called yields) one can write an algorithm as if the sequence was being created at once and only insert yields at appropriate times, e.g. when a new item gets created. The language handles the rest. Each subsequent call to the generator method resumes computation from the last evaluated yield and continues to the next one, e.g. creating a new element each time.

The Zend engine has a native support for generators. It intrinsically understands yields and is, on their evaluation, able to save the state of current execution [Popov, 2017]. CIL has no such first class support. For that reason languages built on top of CIL have to implement generators through other means [Lippert, 2008]. Usually by rewriting generator methods into state machines with explicit state saving before each yield and state retrieval in the beginning.

That is exactly what this thesis covers. It describes the design and implementation of support for PHP generators within the Peachpie compiler through semantic tree transformations, implementation of new semantic tree nodes, and extensions to Peachpie runtime library. In the implementation parts it tries to not only plainly cover the code but also to depict the decision process that led to choosing certain approaches over others. During the whole work we will compare

¹zend.com/en/community/php

²hhvm.com/

³peachpie.io/

⁴Chapter 2.1

our approach with the one taken by C# team and its compiler Roslyn. C# was chosen as a reference language due to it being the prominent language in .NET platform.

While the goal is to implement support for generators with as much original PHP semantic as possible, due to the scope of this work we will not discuss the specific implementation of all PHP generators features. Namely, we will not cover handling yields in exception control blocks (try, catch, finally) in detail and will leave its implementation for future work.

Thesis structure

This thesis is divided into seven chapters. The first one covers general concepts of generators both in PHP and in other languages, explaining what they are, what features and limitations do they have, and where they stand in regards to iterators.

The second chapter briefly introduces the .NET platform and its intermediate language CIL. The third is all about the Peachpie project. It describes its architecture focusing mainly on the semantic tree data structure and CIL emit phase of the compiler. In fourth chapter we take a look at how generators are implemented in C#'s Roslyn and PHP's Zend engine. Especially the Roslyn's approach is important because it serves as a basis for our own implementation.

Generators within Peachpie is the focus of the fifth chapter which itself is further divided into five more sections. First describes an implementation of generators limited to circa C# generators. It builds on the theoretical basis described in previous section about Roslyn's approach. Second proposes theoretical algorithm to handle yield as expression. Third talks about implementation of said algorithm within peachpie. Fourth briefly mentions possible solutions for yields in exception handling blocks. And the fifth is about possible future work that could be done for generators support within peachpie.

The sixth chapter concludes and summarizes the whole thesis. And last but not least the final chapter provides a lightweight user documentation for the peachpie project and overview of attachments.

1. Generators

1.1 Iterators

Before going into generators we need to define iterators first. Iterators, as their name suggests, represent a state of iteration of some sequence backed by either themselves or some other object. In both C# and PHP they can be arbitrary objects implementing an *Iterator*¹ (for PHP) or an *IEnumerator*² (in C#) interface.

Both of these contain a number of methods (Figure 1.1). However, for now we will focus on only two of them they both share, *current* and *next*. The first one - *current* - should always return the element the iterator is currently pointing at. As such it should never modify the state of the iterator and should generally be side effect free. The second one, *next*, should advance the iterator to the next item effectively changing what element the *current* method returns.

Iterator	IEnumerator<T>
+ current(): PhpValue	+ Current: T
+ key(): PhpValue	+ MoveNext(): bool
+ next(): void	+ Reset(): void
+ rewind(): void	
+ valid(): boolean	

Figure 1.1: Iterator and IEnumerable interfaces.

Iterators serve, among other things, as a useful tool to handle large sequences of items. Instead of generating everything at once and then keeping it, for example, in an array one can simply create an iterator that generates and serves the elements one by one. This enables lower memory footprint because there is no need to ever keep all the elements in memory at once.

Furthermore instead of spending a lot of resources at once in the beginning to create the whole sequence the iterator can now spend just a little bit for each element's creation when the *next* method is called. Though the final price is the same regardless it enables a more even distribution of the performance hit. Also, one does not need to pay for the elements that do not get created - the ones that are not iterated over.

There is one problem with iterators, though. They are quite tedious to write. The main issue stems from the fact that unlike in a normal method in which you would create the whole sequence at once in the iterator's *next* method you need to always explicitly save and then retrieve the current state of the execution. There is a lot of boilerplate associated with them as well. You need to create a new type and implement a number of methods that do not actually do anything useful.

¹[PHP.Net]

²[MSDN]

To give an example (Listing 1), creating an array of numbers 1 to n is a straightforward for loop with an assignment. In the case of an iterator's *next* method you need to first retrieve last used number, increase it by one, and then store it. You also need to implement a *current* method that returns the last stored index and a few other ones that are basically just a busywork.

```
<?php
function by_one_at_once(){
    $result = [];
    for($i = 0; $i < 10; $i++){
        $result[] = $i;
    }
    return $result;
}

class byOneIterator implements Iterator {
    private $position = 0;

    public function rewind() {$this->position = 0;}
    public function key() {return $this->position;}
    public function current() {return $this->position; }
    public function next() {
        $curr_pos = $this->position;
        $curr_pos += 1;
        $this->position = $curr_pos;
    }
    public function valid(){ return ($this->position < 10); }
}
```

Listing 1: Method that creates everything at once and as an Iterator.

While for a monotone sequence of numbers even the iterator is still simple and readable it quickly stops to be the case with higher complexity of the iteration. The amount of code needed for state keeping increases quite quickly and an algorithm which would otherwise be really straightforward when used for the creation of the whole sequence at once gets convoluted. And that is where generators come into play.

1.2 Generators universally

Generators provide an easy way to write methods that return iterators while the code can be almost the same as if they returned whole sequences at once instead. All the transforming of the algorithm into the *next* method with its retrieving of the last state in the beginning and state saving after a new element gets set as *current* is handled transparently for the programmer. There is also no need to create a new type and implement other iterator's busywork methods with them. They automatically return correct and fully implemented iterators.

To achieve this a new keyword *yield* or *yield return* is usually introduced. It serves two purposes. First, it marks the spots where the *next* method of the returned iterator should stop executing and save the current state. Second, much like the *return* keyword it specifies a value being returned, in this case actually a value that should be set as *current* on the iterator.

To continue with our example of creating a sequence of numbers from 1 to n (Listing 2). One would write a generator method the same way a normal method generating the whole array with the only difference that instead of an assignment into a result array the for loop would contain a *yield* of the current index.

```
<?php
function by_one_generator(){
    for($i = 0; $i < 10; $i++){
        yield $i;
    }
}
```

Listing 2: By one sequence as a generator.

An iterator returned from such generator method (Figure 1.2) would have a *next* method that would always start from the last encountered *yield*, execute update and condition part of the for loop, and then set a new element as the *current* one.

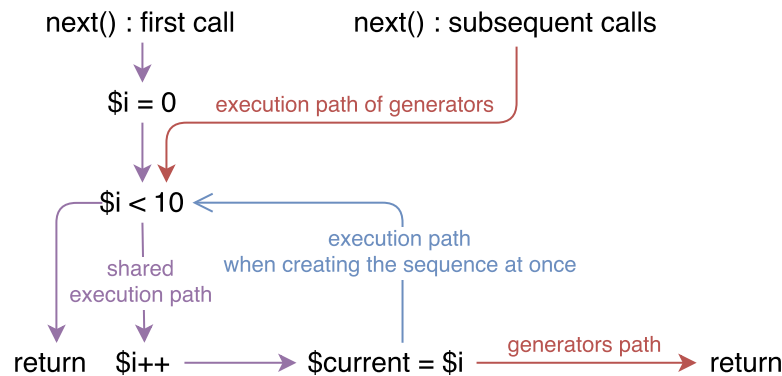


Figure 1.2: Generator's control flow graph.

1.3 Generators in PHP

Generators in PHP are not just about creating data, however. They support consuming data as well (Listing 3). Namely, one can send an arbitrary value into the returned iterator through its method called *send*. This method takes one argument - the value being sent, sets it as a result of the last *yield* expression, and resumes the evaluation the same way as *next* method would.

```

<?php
function logger_generator(){
    $logger = new Logger();
    while(($line = yield) != "END"){
        $logger->log($line);
    }
    $logger->close();
}
class Logger{
    public function close(){ echo "END;";}
    public function log($line){ echo $line; }
}

$gen = logger_generator();
$gen->send("First!");
$gen->send("Second!");
$gen->send("END");

```

Listing 3: Generator method used as a logger.

This means that unlike in C# where *yield* is a statement and therefore can not be a part of a bigger computation or a function call in PHP *yield* can be literally anywhere, even in the place of a function call argument. This further complicates the state saving. In addition to all the local variables when *yield* is part of some bigger expression we need to save its state as well. And it needs to be done the right way to ensure correct order of execution.

Due to various design reasons [Lippert, 2009b] C# also limits where *yields* can happen in regards to exception control blocks. They are not allowed in *catch*³ and *finally* blocks⁴ altogether and can be only in *try* blocks⁵ that do not have any associated catch blocks. PHP, on the other hand, allows yielding everywhere, be it in *try*, *catch*, or *finally* blocks.

1.4 Generators in other languages

These limitations are not unique to C#, however. Both F# and Visual Basic, the only other truly mainstream CIL based languages, also support generators with these restrictions. For them *yield* is just a statement that can not appear in certain exception handling blocks.

In fact, *yield* being an expression holding a send value is not even a feature all dynamic languages share. JavaScript, for example, had just a brief support⁶ for such behavior and as of ECMA 2015 has *yield* only as a statement.

That, nonetheless, does not mean that PHP is completely unique in regards to generators. In python⁷, another mainly interpreted dynamic language, *yields*

³[Lippert, 2009a]

⁴[Lippert, 2009b]

⁵[Lippert, 2009c]

⁶[MDN]

⁷[docs.python.org]

are expressions and are allowed in exception handling blocks almost the same way as in PHP. In Lua⁸ they are implemented as a special version of coroutines and thus have an even wider set of features.

⁸[Lua.org]

2. .NET platform

The .NET platform, or any platform implementing the open CLI standard¹, stands on four pillars (Figure 2.1). The low level intermediate language CIL, the higher level languages such as C# and F# and their compilers to CIL, the base class library known as .NET framework, and - last but not least - the common language runtime, CLR, that actually executes the intermediate code.

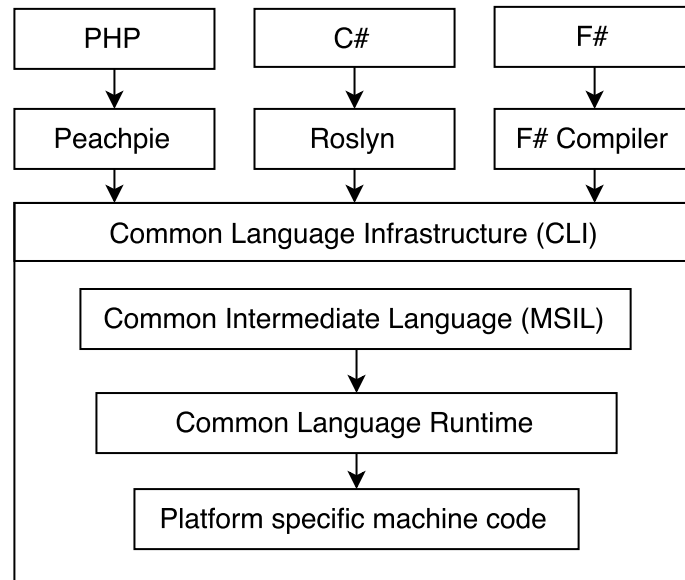


Figure 2.1: Common language infrastructure.

We will talk about the C# and Visual Basic compiler Roslyn later and neither the base class library nor the CLR is particularly interesting to us. The common intermediate language is, however.

2.1 Common intermediate language

CIL is an assembler defined by the common language infrastructure [ECMA-335, 2012] to be a shared basis for all CLI languages (C#, F#, IronPython, ...) and runtime implementations (.NET, Mono, dotGNU, ...). It is platform independent and as such does not natively run on any CPU architecture. Instead it must be either translated to the target platform's native code beforehand or - more commonly - executed by a virtual machine such as CLR.

Despite being an assembler, thus inherently low-level, CIL is actually object oriented and so has a deep understanding of reference types. Its instruction set reflects that with means to create new instances, access their members, and so on. The CLI specification also dictates that by default the CIL should be memory safe.

¹[ECMA-335, 2012]

2.1.1 Evaluation stack

CIL is a stack based assembler, therefore without the notion for registers. Instead it defines a virtual evaluation stack. There are basically two types of instructions in CIL. Memory handling ones that either pop a value from the stack and store it in memory or load a value from memory and push it to the top and the ones that actually do some processing. These pop a few values from the stack, process them somehow, and then store the result on the top of the stack.

There are a few important things to note about the evaluation stack [ECMA-335, 2012, Sec. I.12.4]. Firstly, all parameters and local variables actually live there. They are not ordinary stack values, though. Their place gets reserved and later cleaned automatically and are not accessible through the normal push/pop instructions. Instead there are dedicated ones to work with them.

Secondly, when exiting a function the stack can not contain anything but the returned value. Thirdly, there are instructions only to work with its top. There is no way to query all the elements in the stack, get its height, or to completely save or load it to/from memory.

Lastly, while not a rule the stack is generally used as a store for temporal values instead of proper local variables. For example an expression $2 + 3 * 5$ (Listing 4) would usually result in the load of constants 2, 3, and 5, a multiplication operation ($3 * 5$), at which point the stack would contain 2 and 15, and finally a plus operation that would leave the stack with 17 at its top.

All of these mean that you can not simply pause and save the execution of a method at arbitrary point with just one or even a few CIL instructions. To completely capture the current state you need not only to save all the local variables and parameters somewhere off the stack and but also do the same for every temporal value that might currently live on the stack. And there is no simple way to query what is currently. You need either to construct the information some other way or restrain yourself to saving the state only when the stack is empty.

2.1.2 Exception handling

The last notable thing about CIL is that it has a notion of exceptions and their handling blocks. Try, catch, and finally are all first class citizens in the language and are bound by a number of rules [ECMA-335, 2012, Sec. I.12.4].

CIL does not permit jumping / branching into any exception handling block² unless the source of the jump / branch is within the same block. You can only enter catch and finally regions through the proper exception handling mechanism. And lastly, to leave any of them³ you need to do it via a designed instruction that, in case of try and catch blocks, ensures any potential finally region gets run. Therefore you can neither jump in the middle of a try block nor execute a catch / finally block without throwing a proper exception first.

²[ECMA-335, 2012, Sec. I.12.4.2.8.2.7]

³[ECMA-335, 2012, Sec. I.12.4.2.8.2.8]

```

public void M(int a) {
    int b = 3;
    int c = 5;
    G(a + b * c);
}
public int G(int a){/*Something*/}

.method public hidebysig instance void M (int32 a)
cil managed {
    .maxstack 4
    .locals init (
        [0] int32,
        [1] int32)
    IL_0000: nop          // Do nothing (No operation)
    IL_0001: ldc.i4.3     // Push 3 onto the stack as int32
    IL_0002: stloc.0      // Pop value from stack to local variable 0
    IL_0003: ldc.i4.5     // Push 5 onto the stack as int32
    IL_0004: stloc.1      // Pop value from stack to local variable 1
    IL_0005: ldarg.0      // Load argument 0 (this) onto the stack
    IL_0006: ldarg.1      // Load argument 1 onto the stack
    IL_0007: ldloc.0      // Load local variable 0 onto stack
    IL_0008: ldloc.1      // Load local variable 1 onto stack
    IL_0009: mul          // Multiply values
    IL_000a: add          // Add two values, returning a new value
    IL_000b: call instance int32 C::G(int32) // Call method
        ↪ indicated on the stack with arguments
    IL_0010: pop          // Pop value (returned by G) from the stack
    IL_0011: ret          // Return from method, possibly with a value
} // end of method C::M

```

Listing 4: Simple method in C# and CIL.

3. Peachpie project

The Peachpie project¹ aims to create a bridge between PHP and the .NET ecosystem. While its development started in early 2016 it builds upon the foundations of a much older project, Phalanger², first released in 2004 and also originally developed at the Charles University.

While both projects share the same end goal - bring PHP to .NET, their implementation is quite a bit different. Phalanger, due to being first released before even .NET 2.0 shipped, had to implement almost everything itself. Peachpie, on the other hand, relies heavily on components provided by Roslyn infrastructure in compiler and by DLR at runtime.

Also, while Phalanger supports PHP 5.4 as the highest version, Peachpie was built for PHP 7.1 and beyond from the very beginning. The last major difference is that Peachpie, unlike Phalanger, runs not only on full .NET and Mono but also on a multi platform .NET Core framework.

The Peachpie project is, as of early 2017, still in an active pre-version 1.0 development by the open source community. As such its architecture is not yet finalized and might change in the future rendering following chapter inaccurate.

3.1 Peachpie architecture

As already mentioned Peachpie consists of three more or less separate parts. A compiler that takes PHP sources and produces .NET assemblies, runtime library that provides support for various dynamic features of PHP, and a reimplementa-tion of PHP base class library with its most popular extensions.

Due to the topic of this thesis being an implementation of a compiler feature we will focus mainly on the compiler and only briefly discuss the runtime library. The base class library, while interesting, is completely irrelevant for our work and will be left out.

3.2 Peachpie compiler

The compiler itself is build on the architecture of an open source C# and Visual Basic compiler platform Roslyn. The compilation is logically divided into four main phases (Figure 3.1). First a parser takes a PHP source and creates an abstract syntax tree. In Peachpie this step is actually offloaded to a third party open source PHP parser³. Then Peachpie takes over and binds the AST to a semantic representation in form of a control graph, essentially creating an abstract form of the final final program. The binding phase is also responsible for lowering higher level language constructs.

Next, the semantic graph is used for an extensive data-flow analysis with the intention to resolve dynamic types and generally eliminate as much dynamic behaviour as possible. This step is important mainly for performance reasons.

¹peachpie.io/

²github.com/DEVSENSE/Phalanger

³github.com/DEVSENSE/Parsers

Dynamic dispatch and access at runtime inherently brings a performance hit, especially on .NET CLR that, despite being language agnostic, is still tuned mostly for C# and VB, both of which are statically typed languages.

In the last phase the semantic graph is used to emit the final CIL code and produce a complete .NET assembly. While Peachpie controls the emit of each individual CIL instruction their specific bytecode realization, possible CIL level optimizations, and assembly structure creation is done by Roslyn components *ILBuilder* and *PEBuilder* respectively.

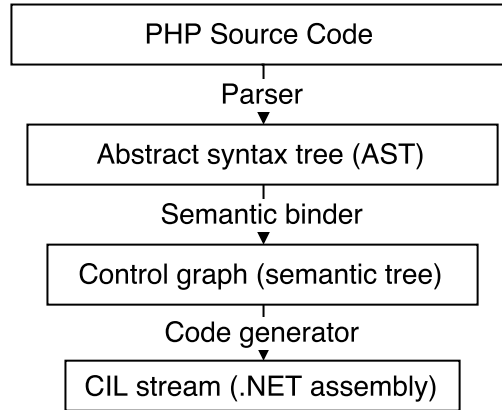


Figure 3.1: Peachpie architecture.

There is obviously more to the compiler part. It also includes a number of code analyzers to catch common PHP bugs, provides an extensive API surface to support projects such as PHP snippets or Visual Studio Code extension, and much more.

3.3 Semantic graph

Since our approach to implementing generator is based on semantic graph transformation let us explore it a bit more. Unlike AST, which is a structured representation of the source code, the semantic graph corresponds more closely to an abstraction of the final program.

It knows the types of all expressions, has all method calls as well as variable/field accesses resolved and bound to specific semantic symbols, and generally contains all the information needed for future compilation.

3.3.1 Statements and expressions

Before going into specific details about the semantic graph itself let us properly define the difference between expressions and statements first. An expression is a combination of values and operations that produces a new value while potentially also having side effects.

For example plus is a binary expression that creates a new value from its two children expressions. Method call, if the method returns some value, is an expression as well. Statement, on the other hand, is an operation that only has some side effects and does not carry a value itself. A good example of a statement is goto jump.

This means that in expression trees, a computation abstraction where each node is an operation that takes the values of its children and produces a new value, statements can only be at the top. Since they do not have a value of their own that could be consumed by their potential parent they simply can not have a parent node.

It is also good to note that while it is simple to transform an expression into a statement from - you simply throw away its value, you can not do it the other way and use a statement in places where expression is expected.

3.3.2 Graph structure

With that out of the way, semantic graph is fundamentally a forest in which every method declared in the source code or synthesized by the compiler corresponds to its control flow graph. Each graph consist of two types of elements. Edges, representing control flow constructs such as loops, branches, and exception handling constructs, and blocks, simple containers holding standalone semantic statements like method calls, assignments, variable declarations, and so on.

These individual statements, can have a form of rather complex graphs themselves (Figure 3.2). They can be arbitrary expression trees with a statement at the top. For example an assignment statement has two expression children, the variable and the one representing the value being assigned. The value expression can also be arbitrary and have children of its own and so on. Semantic expressions and statements are called bound within Peachpie.

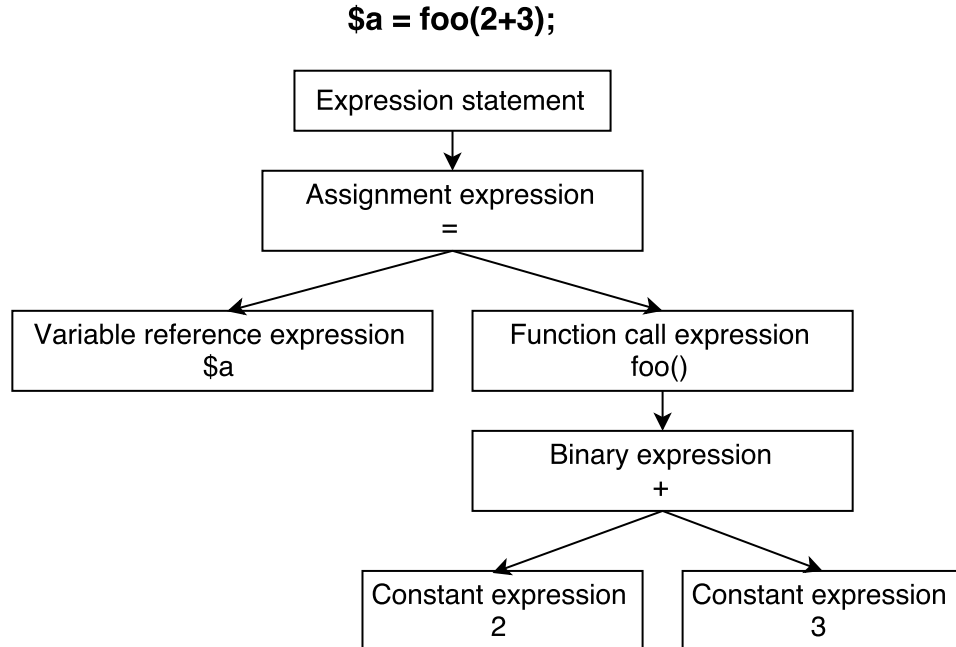


Figure 3.2: Expression tree.

The control flow edges do not have to be simple and connect only two blocks either. In fact most edges connect multiple blocks and some even have references to individual expressions (Figure 3.3). For example a switch edge connects a source block, a switch variable expression, an arbitrary number of case blocks

with their case value expressions, and an end block. Other edges are implemented similarly.

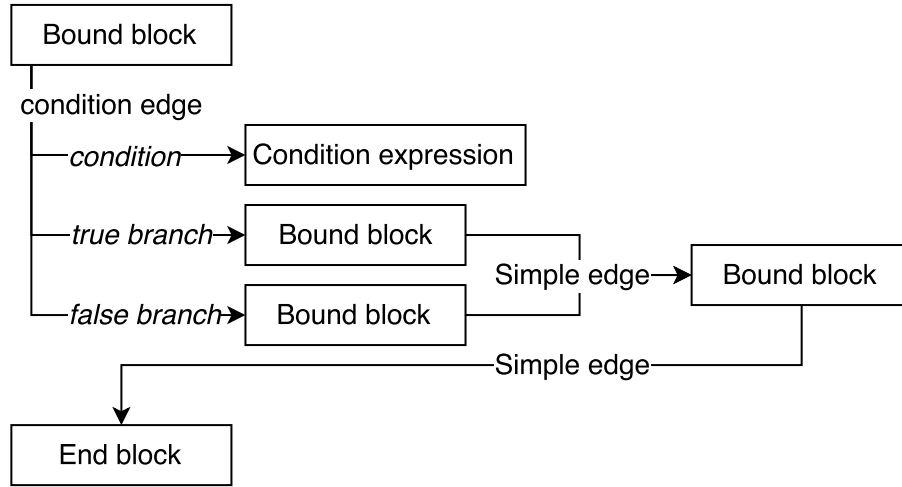


Figure 3.3: Condition edge.

The last type of objects to note in regards to semantic graph are symbols. They represent declared objects. These include types, namespaces, methods, fields, variables, and parameters. Read and write accesses to fields, variables, and parameters inside methods are still represented as bound expressions, however. Bound expressions that hold a reference to these symbols and use them as an identifier of the actual place where their value is.

3.3.3 Graph creation

There are two components in Peachpie compiler responsible for the individual method's control flow graphs creation, *BuilderVisitor* and *SemanticBinder*. *BuilderVisitor* is a higher level component traversing the top level of method's abstract syntax tree and creating aforementioned control flow edges and bound blocks in the process.

It uses the *SemanticBinder* to fill these bound blocks with bound statements and to create sporadic bound expressions needed for the edges, e.g. a switch value expression for a switch edge. Specifically it goes through statements and control flow construct in the method's AST and does two things. It either adds a newly bound statement into the current bound block or, on control flow constructs, creates new bound blocks. When it does that it fills them with statements, connects them to the previous current block, and sets the last of them as the new current block.

The *SemanticBinder* is a component that takes a statement or an expression in form of an abstract syntax tree and creates its semantic representation, either a bound statement or a bound expression, that can be used in the resulting semantic graph.

It also handles the full complexity of the statements/expressions. When it gets asked by the *BuilderVisitor* to bind an assignment it creates not only the bound assignment statement itself but it binds its children, and their children, and so on as well, returning the full bound subgraph.

That is the reason why the *BuilderVisitor* goes only through top level statements and does care about individual expressions. With the exception of some within edges, as noted, each expression is part of a larger expression tree under some statement. A statement that gets bind as a whole tree by the *SemanticBinder*.

3.4 CIL emit phase

In Peachpie the CIL code generation is based solely on the semantic graph. Each of its elements, be it an edge, a statement, or a symbol, has a method, *Generate* or *Emit*, that can create the element's complete CIL code representation. These methods do not produce the bytecode themselves. Instead they use a component called *CodeGenerator*, a thin wrapper around Roslyn's *ILBuilder*.

3.4.1 Code generator

CodeGenerator is fundamentally an abstraction of CIL code stream. It can do two things. Append either individual CIL instructions or their short sequences to the current code stream and realize the stream into actual bytecode. The only higher level service it provides is the ability to change where it should look for certain important items.

One can, for example, set an arbitrary variable as the current this object or specify that local variables should live in some PHP array instead of on the evaluation stack. Appending a load from a local variable then results in the correct CIL sequence being emitted. That means either a load from the locals part of the evaluation stack or, when the place of locals was changed on the current *CodeGenerator*, some *PhpArray*, which itself might need to get loaded from some field first.

This is used, among other places, when there are indirect variable accesses in a method⁴. Because of their indirect and thus dynamic nature it is impossible to resolve and bind them at compilation time. Also there are no CIL instructions to create a new local variable nor to access a variable by its name. Afterall local variables in CIL are just slots in the evaluation stack. Therefore for methods with indirect variables its locals must be moved from the evaluation stack to a *PhpArray* that supports both of these operations.

3.4.2 Emit

The act of emit is very similar to the binding phase. When *Generate* is called on a semantic item it emits CIL representation of not only the item itself but also of all that under it in the semantic graph. A code generation for a method symbol (Figure 3.4), for example, causes emit of all of its bound blocks and edges, each of which triggers emit of their bound expressions and statements, effectively ending up with CIL for the whole method's body being generated.

As such the emit is effectively a mixture of pre and post-order traversal of the semantic tree. First emit of the current item starts, then its children from

⁴php.net/manual/en/language.variables.variable.php

left to right get fully emitted, and then it finishes. Due to that and the fact that execution follows the emitted CIL code there is an invariant. Left children represent code that is executed before right children and nodes lower in the tree finish before their parents.

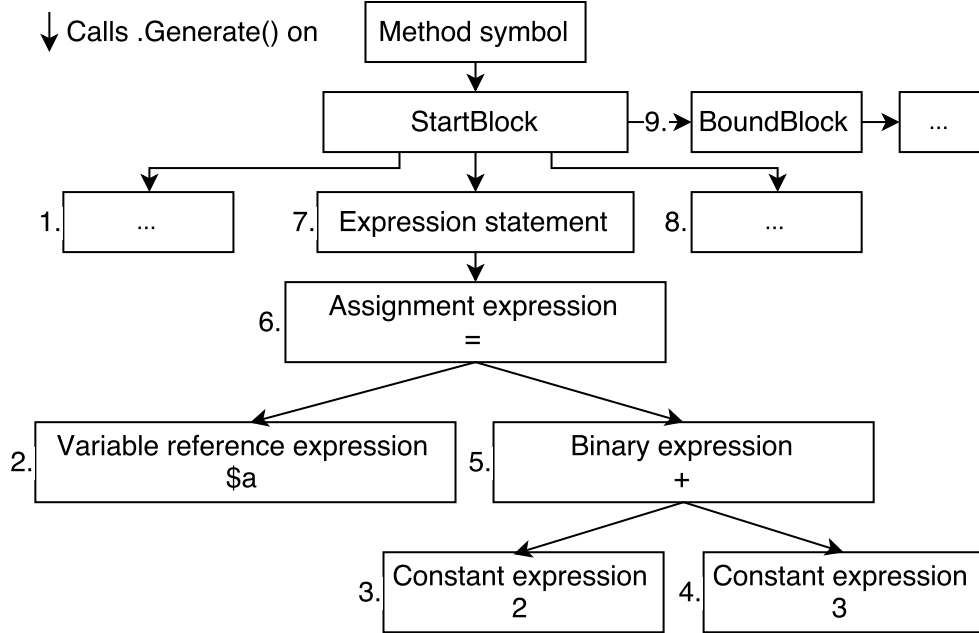


Figure 3.4: Emit order of a method.

The individual *Generate* or *Emit* methods are completely independent. They only append CIL instructions representing their semantic node to the *CodeGenerator* instance they got passed as a parameter. In case their respective nodes have any children within the semantic graph they also call *Generate* on them, always passing the *CodeGenerator* instance they themselves got. This ensures that in the end the one *CodeGenerator* instance contains CIL code for all the method's statements, expressions, and edges.

3.4.3 Generate methods' invariants

While individual nodes can emit themselves however they want there are two rules that must hold true. The code emitted by a bound expression must load and leave its value on the top the evaluation stack. Other than that it can not leave there anything else nor can it remove something. Their *Generate* methods also have to return a symbol representing the expression's IL type.

For bound blocks, statements, and edges the rule is similar with the difference that, since they do not have an inherent value, they can not leave anything on the evaluation stack at all. Their methods do not return anything.

These rules ensure a number of things. First, all expressions are basically interchangeable. An emit of a binary expression does not need to care about the operands' types. It can simply call *Generate* on them and know the evaluation stack will contain their values, independent on whether they are constant expressions, method calls, or something else.

Also, since neither statements nor edges can leave anything on the evaluation stack and statements can only be at the top level of semantic tree the evaluation stack is guaranteed to be empty in the beginning of each bound statement's execution. This is important because there are statements, such as return or goto, that transfer execution and therefore need the evaluation stack to be empty.

3.5 Peachpie runtime library

The runtime library consists of a number of important types needed to support the dynamic nature of PHP. Probably the most important one is *PhpValue*, a managed counterpart of Zend Engine's *ZVal*⁵. It is a type used everywhere the data-flow analysis can not ensure a more specific type. It is a lightweight struct that can hold a primitive value (number, string, or boolean), *PhpArray* (essentially a hashtable), reference to a proper class instance, or - in case of a reference variable - a link to another *PhpValue*. It also supports all the required type conversions and serves as a true representation for an arbitrary PHP value.

Another core type present in the runtime library is *Context*. As hinted by its name it holds the context of the current execution. Be it defined global methods, variables, and constants, the value of static fields and much more. It is a rather special type because it gets silently passed to every PHP method as their first argument throughout the whole execution.

When a library originally written in PHP is used from another .NET language the context has to be explicitly created and passed into it as part of an invocation. If the original PHP app is used on its own the context gets created automatically in the beginning before any user code is run.

Other than that the runtime library also contains function call and variable access resolution logic needed to support dynamic behaviors at runtime, variety of types needed for certain PHP features such as base class for lambdas, and full reflection support.

⁵php.net/manual/en/internals2.variables.intro.php

4. Generators implementation in other platforms

4.1 CSharp and Roslyn

4.1.1 Iterator object and generator methods

4.1.2 Rewriter

4.1.3 Local variables parameters

4.1.4 Execution position

4.2 PHP and Zend Engine

5. Generators in Peachpie

5.1 Basic generators implementation

5.1.1 Iterator object

5.1.2 Next method implementation and local variables

5.1.3 Accessibility of fields on the Generator type

5.1.4 Context handling

5.1.5 Rewriter

5.1.6 Bound yield expression

5.1.7 Start block

5.1.8 Method symbol

5.2 Yield as an expression - theory

5.2.1 Possible approaches

5.2.2 Branch capture & yield splitting

5.2.3 Semantic tree transformation

5.2.4 Short circuit evaluation

5.3 Yield as an expression - implementation

5.3.1 Binding multiple elements

5.3.2 Capturing branches with yields

5.3.3 Correctness of modified capturing algorithm

5.3.4 Creating and keeping the pre-bound graph

5.3.5 Path between the root and yields

5.3.6 Conditioned branches

5.3.7 Implementation remarks

5.4 Yield in exception handling blocks

5.4.1 Yields and exception handling blocks in PHP

5.4.2 Solution in Peachpie

5.5 Future work

Conclusion

In previous four sections we have first described the fundamental concepts required for this thesis to even make sense, then designed an algorithm to support our feature, provided an overview of said algorithm's implementation, and, in the end, proposed possible expansions.

While the work on generators support within the Peachpie compiler is by no means done the shipped implementation provides a good foundation that can stand on its own. It brings support for all generator's features except for yields in exception handling blocks. And while that is an useful feature it is a more of an extension of generators than its fundamental building block. Other than that our implementation mimics the reference semantics faithfully while expanding upon the featureset usual in other CLI based languages such as in C#.

The goal of using as much existing architecture as possible and not creating unnecessary abstractions just for generators was also achieved. While there is still room for improvement all generators specific code is either cleanly separated or repurposed for use by other compiler components.

Lastly, while not an explicitly stated goal the compilation of generators is efficient. It does not introduce any new semantic tree or syntax tree traversals and only slightly increases the memory required for the binding phase. Due to separating all specific logic to a special binder it has absolutely no impact on binding and thus compiling non-generator methods.

In conclusion this thesis and the attached implementation fulfill all goals set by both the thesis assignment and us in the introduction section. On top of that it brings a self-contained functionality to a popular open source project.

Attachments

Attached to this thesis is a snapshot of Peachpie project's git repository. It contains not only the implementation that was done as the practical part of this thesis but also the rest of the complete project. A more up to date version can be found on github¹.

To query only commits done by the author of this thesis please filter out author *Petr Houška* or email *houskape@gmail.com*.

Compilation

The project's only implicit dependency is .NET Core runtime and optionally its CLI SDK. If you want to compile the project yourself you can download both of them from the official site², for Linux, Windows, or MacOSX.

After obtaining the .NET Core SDK please navigate to the folder with the Peachpie repository in your favourite terminal and:

```
dotnet restore //download all external packages required
dotnet build   //build the complete solution
```

Structure

There are three components relevant for this thesis within the repository. The compiler binaries, the compiler implementation, and the generators tests. Below are listed paths to them and in case of the compiler's implementation also to some files containing the majority of our work to support generators.

1. src/Compiler/peach
2. src/CodeAnalysis
 - (a) ./Semantics/SemanticsBinder.cs
 - (b) ./Semantics/Graph/BuilderVisitor.cs
 - (c) src/Peachpie.Runtime/std/Generator.cs
3. tests/generators

Manual testing

To compile an arbitrary PHP file into a .NET assembly with Peachpie invoke the compiler with a path to the PHP file as its first argument. The compiler assembly resides at aforementioned path and is called peach.exe or peach.dll depending of whether it was compiled for full .NET framework or .NET Core.

¹ github.com/peachpiecompiler/peachpie

² microsoft.com/net/download/core

```
$\src\Compiler\peach> dotnet run .\test.php
```

Please do note that an assembly compiled this way will require PHP runtime libraries to run. These libraries can be found, for example, in the bin output of the compiler (peach) project.

Alternatively it is possible to use a Peachpie console application sample³. It includes a .msbuildproj file that configures the .NET Core CLI to download and use both the Peachpie compiler toolchain and required runtime libraries automatically. More about that approach can be found on a peachpie blog⁴.

Automatic testing

The Peachpie project includes a comprehensive set of automatic tests. These consist of PHP files that get compiled by the Peachpie compiler and run by a .NET runtime. If there is a PHP runtime present in the current path environment variable they get run by it as well. The results are then compared to ensure Peachpie compilation keeps the original PHP semantics and is, in terms of runtime behaviour, indistinguishable from the reference implementation.

There is a number tests created as part of this thesis that ensure the implementation of generators support works correctly. They are located in a subfolder tests/generators. While they are in no particular order it is generally true that the higher their number the more complex aspect of generators they test. Below is a command that invokes all peachpie tests, including generator ones.

```
$\src\Tests\Peachpie.ScriptTests> dotnet test
```

Please do note that two tests usually fail on some machines because of encoding issues.

³github.com/iolevel/peachpie-samples/tree/master/console-application

⁴ peachpie.io/2017/04/tutorial-vs2017.html

Bibliography

- docs.python.org. 6.2.9. Yield expressions. URL <https://docs.python.org/3/reference/expressions.html#yield-expressions>. Accessed: 2017-06-25.
- ECMA-335. Common Language Infrastructure (cli), June 2012. URL <https://www.ecma.ch/publications/standards/Ecma-335.htm>. 6th edition.
- Benjamin Fistein and Jakub Míšek. Peachpie benchmarks, 2016 - 2017. URL <http://www.peachpie.io/2017/06/optimizing-php-code-1.html/>. Accessed: 2017-06-25.
- Eric Lippert. Iterator block implementation details: auto-generated state machines, 2008. URL <http://csharpindepth.com/Articles/Chapter6/IteratorBlockImplementation.aspx>. Accessed: 2017-06-25.
- Eric Lippert. Iterator blocks part four: Why no yield in catch?, July 2009a. URL <https://blogs.msdn.microsoft.com/ericlippert/2009/07/20/iterator-blocks-part-four-why-no-yield-in-catch/>. Accessed: 2017-06-25.
- Eric Lippert. Iterator blocks, part three: Why no yield in finally?, July 2009b. URL <https://blogs.msdn.microsoft.com/ericlippert/2009/07/16/iterator-blocks-part-three-why-no-yield-in-finally/>. Accessed: 2017-06-25.
- Eric Lippert. Iterator blocks, part five: Push vs pull, July 2009c. URL <https://blogs.msdn.microsoft.com/ericlippert/2009/07/23/iterator-blocks-part-five-push-vs-pull/>. Accessed: 2017-06-25.
- Lua.org. 9.1 – coroutine basics. URL <https://www.lua.org/pil/9.1.html>. Accessed: 2017-06-25.
- MDN. Generator. URL https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Global_Objects/Generator. Accessed: 2017-06-25.
- Jakub Míšek. Optimizing php code with peachpie – part 1, June 2017. URL <http://www.peachpie.io/2017/06/optimizing-php-code-1.html/>. Accessed: 2017-06-25.
- MSDN. IEnumerator Interface. URL <http://php.net/manual/en/class.iterator.php>. Accessed: 2017-06-25.
- PHP.Net. The Iterator interface. URL <http://php.net/manual/en/class.iterator.php>. Accessed: 2017-06-25.
- Nikita Popov. PHP 7 virtual machine, April 2017. URL <http://nikic.github.io/2017/04/14/PHP-7-Virtual-machine.html#generators/>. Accessed: 2017-06-25.
- Stack Overflow. Stack overflow developer survey results 2017, 2017. URL <https://insights.stackoverflow.com/survey/2017#technology/>. Accessed: 2017-06-13.

TIOBE. Tiobe index, 2017. URL <https://www.tiobe.com/tiobe-index>. Accessed: 2017-06-13.

List of Figures

1.1	Iterator and IEnumerable interfaces.	5
1.2	Generator's control flow graph.	7
2.1	Common language infrastructure.	10
3.1	Peachpie architecture.	14
3.2	Expression tree.	15
3.3	Condition edge.	16
3.4	Emit order of a method.	18

List of Abbreviations

CLI Common language infrastructure, open standard for runtime environment implemented by .NET, Mono, and others.

CIL Common intermediate language, object oriented assembler defined by *CLI* (also known as *MSIL* or *IL*).

CLR Common language runtime, virtual machine implementing the execution engine specified by *CLI*.

DLR Dynamic language runtime, set of libraries providing compiler and runtime services for dynamic languages build on top of *CLR*.

AST Abstract syntax tree, structured representation of the source code.

CFG Control flow graph, a semantic graph representing a method.