

ECOLE POLYTECHNIQUE FEDERALE DE LAUSANNE

MASTER THESIS

Efficient Deoptimization

Author:
Adrien GHOSN

Supervisor:
Prof. Jan VITEK, Prof. Viktor
KUNCAK

*A thesis submitted in fulfilment of the requirements
for the degree of Master in Computer Science*

in the

LARA
Computer Science

February 1, 2016

Declaration of Authorship

I, Adrien GHOSN, declare that this thesis titled, “Efficient Deoptimization” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

“Thanks to my solid academic training, today I can write hundreds of words on virtually any topic without possessing a shred of information, which is how I got a good job in journalism.”

Dave Barry

ECOLE POLYTECHNIQUE FEDERALE DE LAUSANNE

Abstract

Faculty Name
Computer Science

Master in Computer Science

Efficient Deoptimization

by Adrien GHOSN

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Contents

Declaration of Authorship	iii
Abstract	vii
Acknowledgements	ix
1 Introduction	1
1.1 Proposed Solution	1
1.2 Paper Overview	1
2 On-Stack Replacement	3
2.1 Overview	3
2.1.1 Definition	3
2.1.2 Why is OSR interesting?	4
2.2 On-Stack Replacement Mechanisms	5
2.2.1 The OSR Points	5
Guarded & Unguarded	5
OSR Entries & OSR Exits	6
2.2.2 The Transition Mechanism	6
High-Level requirements and challenges for OSR transitions	6
The OSR transition at a low-level	8
Transition as a Function call	8
2.2.3 Versionning	11
One version	11
Multiple Versions	11
Caching vs. Generating on the fly	11
2.3 The Deoptimization case	12
2.3.1 Why the Deoptimization case is more interesting?	12
2.3.2 Where do we exit?	12
The Interpreter	13
The compiled base version	13
A less optimized version	14
2.4 Constraints on the OSR Mechanism	15
2.4.1 The OSR trade-offs	15
2.4.2 Limits on optimizations	15
3 Related Work	17
3.1 The origins: SELF Debugging	17
3.2 OSR & VMs	20
3.2.1 Java HotSpot	20
3.2.2 Jikes RVM	21
3.2.3 WebKit VM	22
3.3 OSR & LLVM	24

3.3.1	What is LLVM?	24
3.3.2	Why do we want OSR in LLVM?	25
3.3.3	OSR as an LLVM library: the McJIT OSR implementation	26
3.4	A classification summary	30
4	Deoptimization in RJIT	31
4.1	Overview	31
4.1.1	Justification	31
4.1.2	An LLVM Library	32
4.2	Resolved & Open OSR	33
4.2.1	The Open OSR	33
	Similarities with McJIT	34
4.2.2	The Resolved OSR	34
4.3	The API	35
4.3.1	StateMaps	35
4.3.2	OSR Points	35
4.3.3	OSR Conditions	35
4.3.4	The continuation function	35
4.4	Using OSR Kit to deoptimize	35
4.4.1	The unease of deoptimizing	35
4.4.2	SOMETHING	35
4.4.3	SOMETHING ELSE	35
4.5	Extending OSR Kit	35
4.5.1	Handling versions	35
4.5.2	Unguarded OSR points	35
4.5.3	Fusing versions?	35
5	Case Study	37
5.1	The RJIT project	37
5.2	OSR in RJIT: requirements	37
5.3	Implementation	37
5.3.1	The Function Call in RJIT	37
5.3.2	Challenges	37
5.4	OSR Conditions	37
6	Conclusion	39
A	Appendix Title Here	41
	Bibliography	43

List of Figures

2.1	OSR exit options	13
3.1	physical vs. source-level stacks	18
3.2	Recovering the source-level state	19
3.3	The WebKit FTL	23
3.4	SSA example	25
3.5	OSR classification	26
3.6	Retrofitting an existing JIT with OSR support	27
3.7	A CFG of a loop with no OSR point	27
3.8	The CFG of the loop in Figure 3.7, after inserting an OSR point	28
3.9	The transformed CFG of the loop in Figure 3.8 after the OSR Pass	28
3.10	A CFG of a running function before inserting the blocks for state recovery	29
3.11	The CFG of the loop represented in Figure 3.10 after inserting the state recovery blocks, from[15].	30

List of Tables

3.1 Summary of OSR implementations and their features.	30
--	----

List of Abbreviations

OSR	On-Stack Replacement
JIT	Just In Time
VM	Virtual Machine
LLVM	Low Level Virtual Machine

For/Dedicated to/To my...

Chapter 1

Introduction

Put the introduction here.

1.1 Proposed Solution

Description of what I've done.

1.2 Paper Overview

Explain each chapter.

Chapter 2

On-Stack Replacement

2.1 Overview

2.1.1 Definition

On-Stack replacement (OSR) is a set of techniques that allow to interrupt running code during the execution of a function, transform it, e.g., re-optimize it, and resume the interrupted function at the point and state at which it was interrupted. The action of transferring the execution using the OSR mechanism is called an *OSR transition*.

On-Stack replacement can be viewed, at a high level, as a mechanism that allows to transform the currently executing function into another version of itself. This transformation mechanism has been used to allow transitions between different levels of code optimizations. We can therefore reduce it to two main purposes: (1) transforming an executing function into a more optimized version of itself, and (2) undoing transformations that were previously performed. While similar, these two types of transformation have very different goals.

In several virtual machines ([19, 15] MORE), some of which will be presented in Chapter 3, on-stack replacement has been used to improve the performance of long running functions. When the virtual machine (VM) identifies a piece of code as being "hot", i.e., it hogs the execution, the VM suspends the function's execution, recompiles it to a higher level of optimization, and transfers the execution to the newly generated version of the function. This differs from a simple Just-In-Time (JIT) compiler, since the recompilation and transfer of execution both take place during the execution of the function, rather than just before its execution. However, both techniques rely on the same assumption, namely, that run-time profiling data enables to uncover new optimization opportunities.

On-Stack replacement allows a compiler to perform speculative transformations. Some optimizations rely on assumptions that are not bound to hold during the entire execution of a program. A simple example is function inlining in an environment where functions can be redefined at run time. A regular, and correct, compiler would not allow to inline a function that might be modified during the execution. The OSR mechanism, however, enables to perform such an optimization. Whenever the assumption fails, i.e., the function is redefined, the OSR mechanism will enable to transfer the execution to a corresponding piece of code where the inlining has not been performed. In this case, OSR is used to allow potentially erroneous transformations while preserving correctness.

In summary, on-stack replacement is a powerful technique that can be used to either improve performance, or enable speculative transformations of the code while preserving correctness. The next section presents the advantages that OSR uncovers in both cases.

2.1.2 Why is OSR interesting?

This section highlights the benefits of using on-stack replacement, with respect to both optimization and deoptimization.

On-Stack replacement increases the power of dynamic compilation. A function can be recompile while it is executing. This enables more aggressive adaptative compilation than just-in-time (JIT) compilation. In fact, the OSR mechanism, in theory, allows to perform a recompilation at any instruction boundary. As an adaptative compilation mechanism, it enables to delay the compilation of a code artefact and therefore enables to gather more information about the current execution profile. These information can then be used to produce higher quality compiled code, with better execution performances.

REFORMULATE AND USE QIAN CITED PAPERS!! For dynamic languages and VMs, code specialization is the most efficient technique to improve performances (IS THAT TRUE? FIND AND CITE). Code specialization consists in tuning the code to better fit a particular use, hence yielding better performances. Specialization can be viewed as a mechanism relying on the versioning of some piece of code. One of the main challenges is to identify which version better fits the current execution need. This requires to gather enough profiling information, some of which might not be available until some portion of the code is executed multiple times.

OSR, coupled with an efficient compiler to generate and keep track of specialized functions, enables to uncover new opportunities to fine tune a portion of code. While techniques like JIT compilation can generate specialized code at a function level, i.e., before the function's execution, OSR enables to make such tuning while a function is running. For example, in the case of a long running loop inside a function, JIT techniques would need to wait until the function is called anew to improve its run time performance by recompiling it. OSR, on the other hand, gives the compiler the means to make such improvements earlier, hence yielding a better overall performance of program's execution.

OSR is a tool that enables the compiler to recompile and optimize at almost any time during the execution of a program. A clever compiler can perform iterative recompilation of the code in order to improve the quality of the generated compiled artefact. OSR enables these iteration steps to be closer to each other and potentially converge to a better solution faster than other dynamic compilation techniques.

On-Stack replacement's most interesting feature is deoptimization. While optimization enables to increase performance, deoptimization's goal is to preserve the program's correctness in the presence of aggressive speculative optimizations. Virtually any assumption can be used to generate compiled code and, if the assumption fails, OSR enables to revert back to a safe version during the execution.

2.2 On-Stack Replacement Mechanisms

2.2.1 The OSR Points

The OSR mechanism is enabled at specific instruction boundaries in the user's code. Depending on the OSR implementation, these points can be sparse, restricted to special points in the program control flow, or associated with special framework dependent data. In the literature, such boundaries were given various names, e.g., *interrupt points*[11], *OSR points* [9, 11, 14, 15], *OSR entries*, and *OSR exits*[14, 15]. This section presents the terminology for such points that will be used in the remainder of this report.

Definition 2.2.1 *An OSR point is a portion of code that belongs to the OSR instrumentation. It is located at an instruction boundary at which the program can be suspended. OSR points can correspond to several instructions inserted by the OSR framework and enable to switch the execution from one version of the code to another.*

An OSR point has to be located at a point of the code where the state is *safe*, i.e., points where the state is guaranteed to be consistent. For example, as explained in Section 3.1, the SELF debugger[11] considers points at method prologues and at the end of loop bodies. An OSR point must be such that the state of the computation can be extracted and transferred to another version of the code from which we can resume the execution. As a result, different OSR implementations yield different restrictions on the location of OSR points.

Guarded & Unguarded

OSR points can be divided into two categories: *guarded* and *unguarded*.

Definition 2.2.2 *A guarded OSR point is an OSR point which execution is subjected to an explicit boolean test, called the OSR condition. In other words, the framework has enough information locally to test whether or not an OSR transition should be taken.*

Definition 2.2.3 *An unguarded OSR point is an OSR point that does not require to check a boolean condition before being executed.*

The distinction between guarded and unguarded OSR points is present in the literature (e.g., Jikes RVM OSR implementation[9, 21] and the WebKit OSR implementation[14]). Both solutions present different advantages and fit different use cases.

While simple to model at a high-level, a guarded OSR point presents the disadvantage of increasing the code's length. A boolean condition needs to be tested every time the execution comes across an OSR point. As explained previously, one of OSR's most interesting use case is long running loop's optimization. Executing many times the branching test can therefore greatly impact the observed performance.

Unguarded OSR points enable to avoid this performance slow down by letting the framework modify the running function on external/environmental events. In WebKit[14], for example, an unguarded OSR point corresponds to a portion of code that can be overwritten at run time, without recompiling the running function, in order to introduce an unconditional execution of an OSR point. This can be viewed as a lazy triggering of an OSR transition, i.e., the instrumentation to enable the transition can be added just before being executed.

OSR Entries & OSR Exits

There is a further distinction between OSR points that are responsible for optimising the code, and the ones that are used to unoptimize. For that, we will use the same terminology as in WebKit[14] and McJIT[15].

Definition 2.2.4 *An OSR entry is an OSR point that enables to replace the current version of the code with one in which we expect to have better performances. They are used in the optimization process.*

Definition 2.2.5 *An OSR exit is an OSR point that enables to exit the current version of the code when it is invalidated. OSR exits are responsible for preserving the correctness of the program. They are used in the deoptimization process.*

At the implementation level, OSR entries and exits are very similar, if not identical. They both are OSR points and will trigger an OSR transition. Conceptually, however, an OSR exit is part of the mechanism that enables to protect the correctness of the program. Mistakenly not taking an OSR entry is a missed opportunity for better performances. Missing a required OSR exit leads to a possibly incorrect behaviour of the program.

2.2.2 The Transition Mechanism

The OSR mechanism can be narrowed down to its core functionality: allowing a transition at run time from the middle of a running function to a corresponding instruction in the middle of another version of this function. This transition mechanism is therefore the cornerstone of any OSR implementation. This section strives to unify all the existing implementations under a single, general, high-level list of requirements to perform an OSR transition.

In the reminder of this section, we use the terminology introduced by definitions 2.2.6 and 2.2.7, and taken from the OSRKit[6] project. These definitions hold, regardless of the type of OSR transition, i.e., OSR exit or OSR entry, being performed.

Definition 2.2.6 *A function from which we perform an OSR transition is called a **from function** or **source version**.*

Definition 2.2.7 *A function in which we resume the execution after an OSR transition is called a **continuation function** or **destination version**.*

High-Level requirements and challenges for OSR transitions

An OSR transition relies on a mapping between two program points and their respective states. Such a mapping is possible only if all the information available at the OSR point in the from function is sufficient to reconstruct the state expected at the landing site in the continuation function.

Reconstructing the state at the landing site means that all the expected live variables must be assigned the correct values to safely resume the execution in the continuation function. Propagating live values is, itself, a challenge, and requires to keep a mapping

between the live values in the from function and the live values in the continuation function. Compiler optimizations might eliminate variables, fuse them, create new ones, and hence prevent a one-to-one mapping between the two sets of live values. As a result, the OSR transition implementation must provide a function f that takes as input the set of live variables and their assigned values S_s from the source version s and generates a new set of live variables mapped to their values S_d for the destination version d .

$$f : S \rightarrow S$$

$$f(S_s) = S_d$$

An OSR transition is transparent to the user, i.e., it should not be observed by the user running the program. The OSR transition can be viewed as a mechanism that deviates the program from its current execution path in order to either improve the performances (the optimization process) or preserve correctness (the deoptimization process). It is part of the compilation/execution framework and, aside from a performance variation, should not impact the execution of the program, e.g., a correct program should not crash because of an OSR transition.

An OSR transition can be either *unidirectional* or *bidirectional*. In the case of a bidirectional mapping between the source and the destination version, the OSR framework must also provide a function f' that is the inverse of f .

$$f' : S \rightarrow S$$

$$f'(S_d) = S_s$$

If the source is a less optimized version of the destination, f' is used during the deoptimization case. The deoptimization case presents a greater challenge than the optimization one. Consider the unoptimized code in the Snippet ??, and its optimized version in Snippet ??.

```
a <- 5
b <- 6
c <- a+b
```

LISTING 2.1: Unoptimized case

```
c <- 11
```

LISTING 2.2: Optimized case

The function f' must provide a way to re-generate a and b from c . To this end, it must remember the steps that led to the optimized version of the code, and must be able to reverse them. This implies a lot of metadata attached to each version of the optimized code generated by the compiler and might quickly increase the space complexity of the OSR support. The VARMAP structure in the Jikes RVM OSR implementation[21] and the scope descriptors in SELF debugger[11] are examples of such metadata attached to the code that enable to recover the values of a and b from c .

The OSR transition at a low-level

At a machine code level, an ideal OSR transition can be summarised into 5 steps:

1. Save the content of the registers.
2. Save the local live values, i.e., the on-stack values in the current stack frame.
3. Modify the current stack frame to introduce the expected local values, add expected stack frames or remove the stack frames that are not needed.
4. put the correct content inside the registers.
5. modify the program counter and resume execution.

In order to perform such a modification of the runtime stack and registers, the OSR framework must have complete control over the machine executing the program. This requirement is easier to satisfy if the execution is performed inside a virtual machine that enables to access and modify the execution stack and the registers. Section 3.3 expands on the relation between OSR implementations and virtual machines. Furthermore, as explained in 2.2.2, reconstructing a state from a set of live values is challenging. Performing this operation at such a low level is hard. The compiled version of the code loses the high-level semantics specific to the language, translates some high-level operations into several sequential instructions and performs some low level resource management (e.g., register allocation) that requires the OSR mechanism to perfectly understand how the compiled code is generated in order to enable the OSR transition at machine code level.

While tempting for its fine-grained control over the transition, the implementation of the OSR mechanism at a low level seems to unveil great challenges. An alternative solution is to encapsulate the entire process at a higher level, closer to the original program's semantics.

Transition as a Function call

An OSR transition can be modelled as a function call. As explained previously, an OSR transition needs to propagate values and jump to a continuation point outside of the current execution path. A function call allows to "modify" the program counter. The OSR transition therefore correspond to a function call from the from function, to a continuation function responsible for executing the new version of the code.

In order to propagate the current execution state, there are two main solutions:

1. Dumping the state in a shared memory location(e.g., as in McJIT OSR implementation[15] and WebKit[14]), or
2. Passing all the needed values as function call arguments to the continuation function(e.g., as in OSR Kit[6] and Jikes RVM OSR [9]).

The first option implies that the from function must dump its state into a buffer before calling the continuation function. The continuation function starts its execution by loading the values it needs from the buffer. This solution has the drawback of increasing both the from function and the continuation function lengths, i.e., the OSR propagation of the state has a greater impact on the performance of the execution than

a simple function call cost.

```

int f_from(a,b) {
    //Some code.
    //...

    if(OSR_CONDITION)
        goto OSR;
CONT:
    //...

OSR:
    //Saves the state
    //into a shared buffer.
    save_state();
    return f_continuation(a,b);
}

int f_continuation(a,b) {
    //Loads the state
    //from the shared buffer.
    load_state();
CONT:
    //Equivalent for CONT.
    //...
}

```

LISTING 2.3: Figures/ffromBuff.c

For the second option, most of the work needed to propagate the state is performed during the compilation. Values are passed as arguments to the call to the continuation function. If there is not a one-to-one mapping between the from variables and the continuation ones, some compensation code is executed, either before the call or at the beginning of the continuation function, in order to assign the correct values to the continuation local variables. The cost of executing the compensation code is also present in the first option, i.e., both solutions have to perform the same set of transformations in order to generate the expected local values. On the other hand, the second option does not require to access shared memory to store and load propagated values, and is therefore expected to have a smaller cost at execution.

```

int f_from(a,b) {
    //Some code.
    //...
    //[Live values: c,d,e,f]
    if(OSR_Condition)
        goto OSR;
CONT:
    //Some code.
    //...

OSR:
    f_continuation(a, b, c, d, e, f);
}

int f_continuation(a,b, c, d, e, f) {
CONT:
    //Some code.
}

```

}

LISTING 2.4: Figures/funCall.c

```

int fbase (a,b) {
PROLOG:
    if (OSR_LABEL == 0001)
        goto OSR1;
    else if (OSR_LABEL == 0002)
        goto OSR2;
    else
        goto ENTRY;
ENTRY:
    //The regular entry code
    //...
CONT1:
    //Some point in the code
    //that corresponds to osr1.
    //...
CONT2:
    //Some point in the code
    //that corresponds to osr2.
    //...
OSR1:
    load_state1();
    goto CONT1;
OSR2:
    load_state2();
    goto CONT2;
}

```

LISTING 2.5: Figures/MultipleEntry.c

The OSR transition requires to enter the continuation function at a special entry point, that corresponds to the OSR point in the from function. In order to decrease the space complexity, one might want to enable multiple entry points per function. Consider Figure REF, where f_{base} provides three entry points: (1)the regular entry point *ENTRY*, (2)a continuation *CONT1* and (3)a second continuation *CONT2*. The function f_{base} can therefore be used as a continuation function for any from function that has an OSR transition from a point that corresponds to either *CONT1* or *CONT2*, as well as for regular function calls. In order to allow multiple entry points, the f_{base} is instrumented, i.e., a *PROLOG* is executed at its beginning. The *PROLOG* is responsible for identifying the f_{base} 's caller. In FIG, this is done by checking the value of a general flag, called *OSR_LABEL*. According to this value, it then jumps to the correct instruction inside the function, and, if an OSR transition is being performed, loads the state from the scratch buffer, as described before.

This instrumentation has an impact on performance as it requires an extra computation at the beginning of the function. On the other hand, the second technique that we described to propagate the state modifies the function's signature. If the set of values to propagate is different at the two points, we need different functions, and hence, should allow only a single entry point for each continuation function. Therefore, we can expect the second solution to have, in general, a higher space complexity (i.e., many versions of the function live in memory at the same time) but a better performance during execution (i.e., smaller execution time).

MAYBE CITE PAPER FROM OSRKIT.

2.2.3 Versionning

On-stack replacement requires, at some point, two have at least two versions of a function: the from function and the continuation function. Furthermore, the continuation function can either be generated on the fly, i.e., the OSR transition includes a transformation step to generate the correct continuation function, or a previously compiled and cached version can be used.

One version

Allowing only one live compiled version of the function enables to simplify the high-level model of the program's execution. Some OSR implementations, such as the McJIT OSR [15], replace the from function with the continuation one in-place, i.e., the continuation function is compiled and registered at the from function's address. This technique enables to propagate the optimized version through out the code without requiring any modification of the callers. One can view it as an eager propagation of optimizations. Therefore, one can assume that the overall performance of the execution should be improved by the global use of this optimized version. On the other hand, this technique requires to ensure that any OSR entry taken is based on a condition that holds at any possible call site of the function. It might therefore prevent the function's specialization in local sub-scopes.

WHAT HAPPENS FOR PREVIOUS CALLS ?

Multiple Versions

Fine-grained specialization of functions enables to improve performances locally. In fact, allowing multiple versions of the same function to be live at the same time enables to optimize a function according to local conditions, regardless of other call sites states. While attractive for the extra flexibility it provides, this solution presents several disadvantages. First, it might increase space complexity. In fact, allowing such fine specialization might lead to having one version per call site. If not carefully handled, this can lead to a dramatic increase of memory usage. Second, this technique is a lazy propagation of optimizations. Assume different calls to the same function. In one call execution, call it C_1 , an OSR condition allows to take a specific OSR entry. Other calls, that will take the same OSR entry, do not benefit from the OSR condition tested by C_1 , i.e., they will have to wait until they reach the same OSR entry, and verify the same condition, before taking the OSR transition.

Caching vs. Generating on the fly

Depending on the reusability of a function's version, and its likelihood to be needed later in the execution, a time vs. space complexity trade-off appears. Generating a continuation function on the fly is time consuming. If the same OSR entry is taken multiple times, and leads to the same version of a function, caching this version might improve run time performances. On the other hand, if OSR entries are rarely taken, i.e., if OSR transitions are rare, generating the continuation function on the fly enables to save memory space. The cost of the compilation is amortised by the improvements allowed by the optimization and the live range of the function. In other words, generating on the fly is viable if the time performance gain of the optimized version over the entire execution is larger than the time required to generate the continuation function. Furthermore, generating the continuation function on the fly enables to take into account

the latest profiling data gathered about the program's execution and use it during the compilation.

2.3 The Deoptimization case

This section presents an overview of the deoptimization case with on-stack replacement. The deoptimization case presents several interesting design decisions and trade-offs.

2.3.1 Why the Deoptimization case is more interesting?

On-stack replacement main interest lies in dynamic deoptimization of code. As explained in 2.1.2, OSR can be used to deoptimize code in order to preserve correctness. A compiler, equipped with an efficient implementation of OSR, is allowed to perform adaptive optimizations, based on some assumptions that might fail at some point, during the execution. Put in another way, OSR enables aggressive speculative optimizations of code by providing a mechanism that preserves correctness.

OSR strives to enable common optimizations performed by static programming language compilers in dynamic language compilers. For example, consider the function inlining optimization, which replaces a function call site with the body of the callee. Inlining in static programming languages is a common optimization (CITE OR GIVE EXAMPLES?). It enables to eliminate the overhead of a function call and return, and enlarges the scope considered for further optimizations. While very common in static languages, inlining proves to be more difficult to implement in compilers for dynamic languages, if functions can be redefined at run time. OSR allows to perform speculative inlining of function calls. The assumption taken by the compiler is that the function call corresponds to a specific function body. A guarded OSR point checks that the inlined function was not redefined before executing the inlined body. If the assumption fails, i.e., the function was redefined, the OSR mechanism allows to exit to a version of the code where the call was not inlined.

The OSR mechanism is only a tool, which efficiency depends on the quality of the assumptions made. This is even more visible in the deoptimization case. The optimization case remains useful for regular optimizations, i.e., an OSR entry can be allowed by some profiling data gathered during the execution, and based on some fact that will hold for the rest of the execution. On the other hand, an aggressive assumption-based optimization exposes a trade-off between the performance improvement that it allows on one side, and the cost it has in the case of an assumption failure on the other. As a result, the OSR deoptimization process allows a greater flexibility but should be coupled with an efficient profiler in order to identify cases where and when the program's execution might benefit from aggressive speculative optimizations.

2.3.2 Where do we exit?

The OSR deoptimization case is interesting in the latitude it leaves with regard to where the optimized code should exit. In the OSR entry case, the continuation function is a more optimized version of the code. In the OSR exit case, there is still a choice to

make: should the continuation function be a fully unoptimized version of the code? Or should we strive to obtain the most optimized version of the code? This section describes several targets for OSR exits, and tries to present the advantages, disadvantages, and requirements for each option.

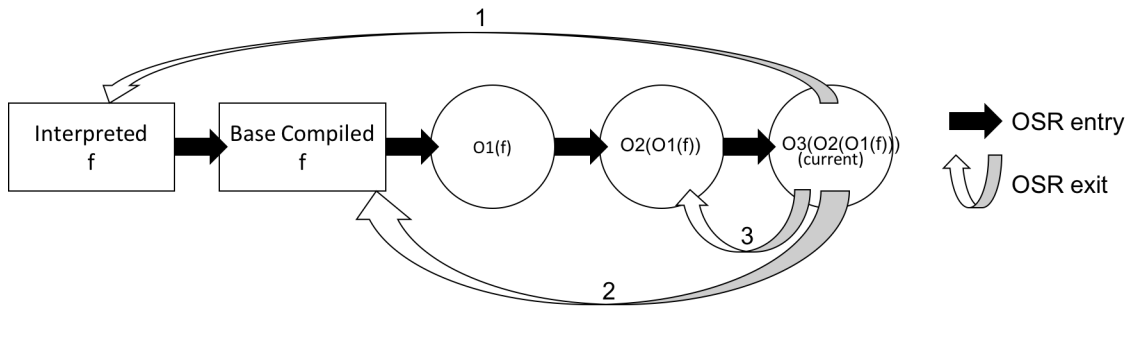


FIGURE 2.1: OSR exit options.

The Interpreter

A first option is to OSR exit back to the interpreter. This approach has been taken by some virtual machines, such as Graal[8] or the Java HotSpot VM[19] (see Section 3.2.1 for more details). This corresponds to the OSR exit 1 in Figure 2.1.

Going back to the interpreter enables to have complete control over the execution environment. The execution framework can therefore easily generate the correct state to resume the execution. Furthermore, this approach gets rid of any assumption taken by the execution framework and goes back to a safe environment. Another way to put it would be to say that, since an assumption failed, the framework takes the conservative decision to go back to an assumption free environment to resume the execution.

Going back to the interpreter is a conservative approach that might impede on the performance of the program. The interpreted version of the function is assumed to be slow. As a result, the function's execution will drastically slow down after exiting to the interpreter. The question is then the following: can we provide the same correctness guarantees while exiting to a more efficient version of the function?

The compiled base version

Another option is to OSR exit to a base, compiled version of the function. It corresponds to the OSR exit 2 in Figure 2.1. This is the approach taken by the WebKit framework[14](see Section 3.2.3). It represents a good intermediary between performance and complexity. If the OSR framework does not have control over the execution environment, it requires careful instrumentation of the compiled code in order to properly propagate the execution state. If the continuation function is generated on the fly, taking this OSR transition implies suffering the cost of compiling a proper continuation

function. On the other hand, the compiled version of the continuation function is expected to run faster than its interpreted equivalent.

On a theoretical point of view, going back to the base version avoids cascading OSR exits. Since an exit was taken, it implies that one of the compiler's assumptions failed. As it is hard to keep track of interferences and relations between different optimization steps, a compiler profiler might not have the means to detect how one assumption's failure impacts other assumptions that were taken. It might be the case that undoing one optimization would yield an incorrect version of the code due to interferences between optimizations (in the worst case), or trigger other OSR exits. An OSR exit also gives some feedback to the compiler's profiler that might improve the decisions that it takes for later speculative optimizations. Restarting the optimization process from scratch might therefore yield a more performant version of the code, compared to the version obtained by undoing the optimization that failed.

A less optimized version

A third option is to exit to an optimized version of the function that does not rely on the assumption that failed. It corresponds to the OSR exit 3 in Figure 2.1. The function to which we exit is a fully optimized one, supposed to yield better performances than both the interpreted and the base versions. This behaviour can be observed in McJIT OSR[15] (see Section 3.3.3) and OSR Kit[6] (see Chapter 4). This option strives to get the best possible performances, even in the presence of an assumption's failure.

As explained previously, there are two options to get the continuation function: either generate on the fly a new fully optimized version that does not rely on the assumption that failed, or "undo" the optimization that failed. Undoing an optimization is hard, especially when it is intertwined with other transformations performed on the code. An alternative is to keep the version of the code that was used to produce the optimized one, and use it as the continuation function for the OSR exit. If the optimization that failed was the last transformation performed on the code, we exit to a version that already underwent all the available transformations and cannot be further optimized. If that is not the case, we exit to a version of the code that is missing several optimizations. For example, consider Figure 2.1. If the assumption used for transformation $O3$ fails, we can easily exit to the version on which we applied the transformation, i.e., $O2(O1(f))$. This corresponds to the OSR exit 3 in Figure 2.1. However, if the assumptions used for $O2$ fails, there is no previous version of f that corresponds to $O3(O1(f))$ and generating one on the fly, from $O3(O2(O1(f)))$, is hard, since it requires to undo $O2$. One possible solution, in that case, is to generate the continuation function on the fly. Starting from f , we apply $O1$ and then $O3$, yielding $O3(O1(f))$. This is a viable solution if and only if it is possible to find a mapping between $O3(O2(O1(f)))$ and $O3(O1(f))$. If such a mapping does not exist, the OSR implementation must ensure that no transformation messes with the OSR instrumentation. In that case, we can always go back to the version on which we applied the transformation that failed, e.g., $O1(f)$ in our example. OSR Kit, presented in Chapter 4, can be used to implement such a solution.

2.4 Constraints on the OSR Mechanism

2.4.1 The OSR trade-offs

On-stack replacement is a powerful technique, that requires to be tailored to its use cases in order to exhibit good results. As seen in this chapter, OSR can be implemented in several different ways, which in turn might lead to very sparse performance results. For example, should the framework cache the compiled versions or should it generate the continuation functions on the fly? This dilemma cannot be answered in the general case. It depends on many conditions, among which the code reusability, i.e., what is the likelihood that the same continuation will be needed later during the program's execution? There is also the trade-off between time and space complexities: caching increases the space usage, generating on the fly increases the time spent in the compiler.

The variety of performances that can be observed does not only depend on the design choices, or the target language, but also on the program's specificities. A program with a certain profile might benefit from a specific OSR implementation, while another might be more amenable to another OSR design. It is hard to know, before hand, which OSR implementation would yield the best results on a specific program, without having enough profiling information about the program's behavior. Furthermore, the program's behaviour might depend on some run time condition, and one OSR implementation is not guaranteed to be the most efficient in every run.

2.4.2 Limits on optimizations

The OSR instrumentation of the code limits the optimizations that the compiler can perform. First, the OSR points must be preserved by the compiler, during the transformations. This has also the disadvantage of impeding on the compiler's performance and more specifically on the quality of the code that it produces. For example, in the Java HotSpot[19], the OSR instrumentation might extend the live range of variables, therefore preventing the compiler from removing them. Another example is the SELF Debugger[11], which deoptimization process does not support tail call elimination and therefore does not implement it in the compiler. Tail call elimination is a compiler optimization that enables, in the case of a tail call, to avoid allocating a new stack frame for the callee. It is commonly used to optimize tail call recursions. The SELF debugger has no way of knowing how many stack frames were eliminated and is therefore unable to reconstruct them properly. However, this limitation does not extend to every deoptimization implementation. A carefully instrumented execution framework could keep track of state information and execution flow in order to enable tail call elimination OSR support.

OSR, and more specifically OSR exits, restricts code motion. In the case of deoptimization, the guarded OSR point protects some portion of code that must not be executed if the OSR exit is to be taken. This implies that the compiler must not move any instruction that depends on this OSR exit condition above the OSR point. The OSR instrumentation therefore limits code motion, which in turn might decrease the number of opportunities for compiler optimizations.

Chapter 3

Related Work

3.1 The origins: SELF Debugging

SELF is a pure object-oriented programming language. It relies on a pure message-based model of computation that, while enabling high expressiveness and rapid prototyping, impedes the languages performances[4, 12]. As a result, the language’s implementation depends on a set of aggressive optimizations to achieve good performances[3, 11]. SELF provides an interactive environment, based on interpreter semantics at compiled-code speed performances.

Providing source level code interactive debugging is hard in the presence of optimizations. Single stepping or obtaining values for certain source level variables might not be possible. For a language such as SELF, that heavily relies on aggressive optimizations, implementing a source code level debugger requires new techniques.

In “Debugging optimized code with dynamic deoptimization”[11], the authors came up with a new mechanism that enables to dynamically de-optimize code at specific interrupt points in order to provide source code level debugging while preserving expected behaviour.

Hölzle, Chambers, and Ungar[11] present the main challenges encountered to provide debugging behaviours, due to the optimizations performed by the SELF compiler. Displaying the stack according to a source-level view is impeded by optimizations such as inlining, register allocation and copy propagation. For example, when a function is inlined at a call site, only a single activation frame is visible, while the source level code expects to see two of them. Figure 3.1, taken from [11], provides another example of activations discordances between physical and source-level stacks. In this figure, the source-level stack contains activations that were inlined by the compiler. For example, the activation B is inlined into A’, and hence disappears from the physical stack.

Single-stepping is another important feature for a debugger. It requires to identify and execute the next machine instruction that corresponds to the source operation. Hölzle, Chambers, and Ungar[11] highlight the impact of code motion and instruction scheduling on the machine instruction layout. Such optimizations re-order, merge, intersperse and sometimes delete source-level operations, therefore preventing a straight forward implementation of source-level single-stepping for the debugger.

Compiler optimizations prevent dynamic changes from being performed in the debugger. Hölzle, Chambers, and Ungar[11] identify two separate issues: changing variable values, and modifying procedures (i.e., functions). To illustrate the first case, the

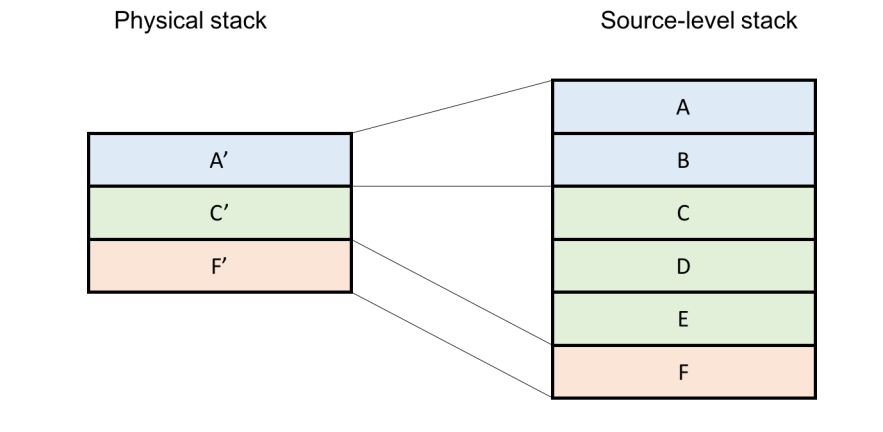


FIGURE 3.1: Displaying the stack, figure from [11].

paper[11] relies on an example where a variable is assigned the sum of two other variables. The compiler identifies the two variables as being constants and replaces the addition by a direct constant assignment. A debugger that allows to change variable values at run time would then yield a non correct behaviour if the user modifies one of the two variables. This problem does not arise in the case of unoptimized code since the addition is still present. For procedures, Hölzle, Chambers, and Ungar[11] provide an example where a function has been inlined by the compiler, but redefined by the user in the debugger.

The paper[11] distinguishes two possible states for compiled code: *optimized*, which can be suspended at widely-spaced interrupt points, from which we can reconstruct source-level state, and *unoptimized*, that can be suspended at any source-level operation and is not subjected to any of the above debugging restrictions.

In order to deoptimize code on demand, SELF debugger needs to recover the unoptimized state that corresponds to the current optimized one. To do so, it relies on a special data structure, called a *scope descriptor*. The scope descriptors are generated during compilation for each source-level scope. This data structure holds the scope place in the virtual call tree of the physical stack frame and records locations and values of its argument and local variables. It further holds locations or values of its subexpressions. Along with the scope descriptor, the compiler generates a mapping between virtual (i.e, scope descriptor and source position within the scope) and physical program counters (PC). Figure 3.2 is taken from[11] and displays a method suspended at two different points. At time t1, the stack trace from the debugger displays frame B, hiding the fact that B was inlined inside of A. At time t2, D is called by C which is called by A, hence, the debugger displays 3 virtual stack frames instead of only one physical frame.

The de-optimization process follows 5 steps described in [11] and summed up here:

1. Save the physical stack frame and remove it from the run time stack.
2. Determine the virtual activations in the physical one, the local variables and the virtual PC.
3. Generate completely unoptimized compiled methods and physical activations for each virtual one.



FIGURE 3.2: Recovering the source-level state (from CITE).

4. Find the unique new physical PC for each virtual activation and initialise (e.g., return addresses and frame pointers) the physical activations created in the previous step.
5. Propagate the values for all elements from the optimized to the unoptimized activations.

Hölzle, Chambers, and Ungar[11] also describe *lazy deoptimization*, a technique to deoptimize a stack frame that is not at the current top of the execution stack. Lazy deoptimization defers the deoptimization transformation until control is about to return into the frame, hence enabling deoptimization for any frame located on the stack.

Deoptimization at any instruction boundary is hard. It requires to be able to recover the state at every single point of the program. The SELF debugger relies on a weaker, relaxed restriction by enabling deoptimization only at certain points called *interrupt points*. At an interrupt point, the program state is guaranteed to be consistent. The paper[11] defines two kinds of interrupt points: method prologues, and backward branches (i.e., end of loop bodies). It estimates the length of the longest code sequence that cannot be interrupted to be a few dozen of instruction, i.e., the average length of a code sequence that does not contain neither a call nor a loop end. Interrupt points are also inserted at all possible run time errors to allow better debugging of synchronous events such as arithmetic overflow. The generated debugging information are needed only at interrupt points, which reduces the space used to support basic debugger operations (as opposed to allowing interrupts at any instruction boundary).

Providing a debugger for SELF limits the set of optimizations that the compiler can support, and decreases the performances of the program when the execution is resumed. Tail recursion elimination saves stack space by replacing a function call with a goto instruction, while fixing the content of registers. SELF debugger is unable to reconstruct the stack frames eliminated by this optimization and hence, it is not implemented in the SELF compiler. More generally, tail call elimination is one important limitation for the SELF debugger.

The debugger slows down the execution when the user decides to resume. The execution should proceed at full speed, but some stack frames might have been un-optimized, hence implying that a few frames might run slowly right after resuming execution.

3.2 OSR & VMs

Virtual machines are privileged environments in which on-stack replacement can be used to its full power. As seen in 2.1.2, OSR is as useful as the compiler's profiler is efficient. A virtual machine (VM) has control over the resources allocation, enables to control the code that is generated by the compiler, maintains important run time data, and state information about the program being executed. Furthermore, modern virtual machines already provide adaptative strategies to recompile hot functions[2, 19, 13, 22], and dispatch calls to the new compiled versions. OSR enables this transition to happen when the function is executing, rather than during a future call.

This section presents several examples of VMs that support on-stack replacement. The section is divided into two parts: we first presents several solutions that provide OSR for various virtual machines, then we briefly introduce LLVM, a VM presenting an interesting framework in which we believe on-stack replacement mechanism should fit.

3.2.1 Java HotSpot

The Java HotSpot Performance Engine is a Java virtual machine developed and maintained by Oracle. The Java HotSpot VM provides features such as a class loader, a bytecode interpreter, Client and Server virtual machines, several garbage collectors, just-in-time compilation and adaptive optimizations.

The Java HotSpot Server Compiler[19] supports on-stack replacement of interpreter frames with compiled-code frames, and deoptimization from compiled-code back to the interpreter. It relies on adaptative optimizations, that focus on performance critical methods[19, 13]. The Java HotSpot Server compiler identifies such methods by using method-entry and backward-branch counters, along with other heuristics related to the method's caller. Thanks to these counters, the run time system is able to identify frequently executed or long running methods and compile them to improve their run time performance. It can further decide to compile one or several of the methods callers. This is called an *upward traversal* and enables to speed the path that leads to a hot method's invocation. On-stack replacement is used when the backward-branch counter exceeds the *OnStackReplacementThreshold* value. The continuation function is

compiled with an *additional* entry point at the target of the backward-branch. A continuation function is therefore able to serve OSR transitions as well as future method invocations. The run time system transfers the execution from the interpreter to this newly compiled code.

In HotSpot, deoptimization might be triggered by two different events: the compilation of a reference to an uninitialised class, and a class loading that invalidates some previous compilation optimization. Class initialization and class loading are performed by the interpreter. Whenever a reference to an uninitialized class is compiled, HotSpot generates an uncommon trap, i.e., a trampoline back to interpreted mode. The optimized code is marked as being unusable, and any thread entering the method are interpreted until the recompilation of the method is achieved properly. The second form of deoptimization happens when a class loading invalidates some optimization performed while generating compiled code (e.g., inlining of methods). A thread that is executing in the invalidated method is "rolled forward"[19] to a safepoint. When the safepoint is reached, their native frame is extracted and converted into an interpreter frame. The thread then continues its execution in the interpreter. For consistency, the class load that invalidated the compiled method becomes visible to the thread only when it reaches its safepoint.

The deoptimization process requires the Java HotSpot to be able to generate an interpreter JVM state at different points of the program. In order to do so, the Java HotSpot Server records the "exact JVM state"[19] at each safepoint and procedure call. This implies that the entire JVM state is considered live at the safepoint, which in turns, might extend the live range of some values. During code emission, the compiler emits a table that maps the JVM state to the resting place of the generated optimized code's JVM state information.

3.2.2 Jikes RVM

The Jikes Research Virtual Machine (RVM)[18] is an open source, self hosted, i.e., it is entirely implemented in Java, virtual machine for Java programs. One specificity of the Jikes RVM is that it exhibits a *compile-only* approach, i.e., the system compiles all bytecode to native code before execution. The RVM provides two compilers: a *base-line* compiler, which generates poor quality code quickly, and an *optimizing* compiler. The optimizing compiler provides a full suite of optimizations, categorized into three different levels. The RVM provides advanced state-of-the-art features such as dynamic compilation, adaptive optimizations, garbage collection, thread scheduling and synchronization[18].

Jikes RVM is an extensible framework in which on-stack replacement techniques have been implemented in two steps, i.e., a first implementation[9], and then an extend version of OSR[21]. Fink and Qian[9] implemented OSR support in Jikes RVM. Their implementation relies on JVM scope descriptor, associated to method activation frames. A scope descriptor contains the thread running the activation, the bytecode index that corresponds to the program counter, values of all local and stack locations, and a reference to the activation's stack frame. The OSR transition is divided into three steps:

1. Extract the compiler-independent state from a suspended thread.

2. Generate the new code for the suspended activation.
3. Transfer the execution in the suspended thread to a new compiled code.

Fink and Qian[9] generate the target function by compiling a specialized version of the method for each activation that is replaced, as well as a new version for future invocations. In other words, instead of allowing multiple entry points per function [15, 19], this JikesRVM OSR implementation generates a specialized version of the target function that has only one entry point, corresponding to the instruction from which the OSR transition was triggered. Each such method contains a special *prologue*, responsible for saving values into locals and loading values on the stack. OSR transitions can be taken at special points, called *OSR points*, introduced by the optimizing compiler and that correspond to points where the running activation may be interrupted. The implementation makes a distinction between unconditional and conditional OSR points. An OSR point is implemented as a call that takes all live variables as arguments. This constraints some optimizations such as dead code elimination, load elimination, store elimination and code motion by extending the liveness scope of some variables. An OSR point transfers control to an exit block, i.e., it can be viewed as a non-return call.

Soman and Krintz[21] proposed a general-purpose OSR mechanism, for the Jikes RVM, that presents less restrictions for compiler optimizations than the previous approach, while decoupling the OSR implementation from the optimization process. They extend the previous OSR implementation[9] to enable OSR transition at, and across points at which the execution can be suspended. In order to do so, they rely on a special data structure called a *variable map* (VARMAP). A VARMAP is associated with each method, and consists in a list of thread-switching points and their live variables. When the compiler performs optimizations, the VARMAP is automatically updated accordingly. Once the compilation completes, the VARMAP is encoded into a compressed map that contains an entry for each OSR point present in the method. The VARMAP is decoupled from the compiled code, i.e., it does not influence the liveness of local variables. Soman and Krintz[21] also propose an alternative lazy triggering of on-stack replacement. Lazy triggering consists in taking an OSR transition due to events in the environment, i.e., events triggered by the runtime. Whenever the runtime deems an assumption invalide, it invokes a helper function called *OSR helper* that either patches the code of the executing methods to call the OSR process, or it modifies return addresses of the callees of the method to point to the OSR helper. When a callee returns, the OSR helper creates a new stack frame with the state extracted from the specialized method's stack and saves all of the specialized method registers into its stack frame. The return address of the OSR helper points to the current instruction in the specialized code, which is then used during OSR to identify the location to resume the execution in the new version of the method. Lazy triggering improves the code efficiency by avoiding the extra cost of guards evaluations.

3.2.3 WebKit VM

WebKit[14] is an open source web browser engine used to improve JavaScript performances. It exhibits a Four-Tier VM architecture. The WebKit run-time compilation flow is described by Figure 3.3.



FIGURE 3.3: The WebKit Four-tier optimization flow.

Forward arrows represent *OSR entries*, i.e., a transformation that yields a more optimized version of the code at run time. Backward arrows correspond to *OSR exits*, i.e., a transformation that yields a less optimized version of the code at run time. The low level interpreter (LLInt) is used for low latency start up. The baseline JIT generates WebKit bytecode with no optimization enabled. The transition from the first tier to the second one happens when a statement is executed more than a hundred times or a function is called more than six times. The data flow graph (DFG) JIT is triggered when a statement executes more than a thousand times or a function is called more than sixty-six times. The Fourth-Tier LLVM(FTL)[20] relies on LLVM machine code optimizations to generate a fast version of portions of the code. In order to hide the costs of the translation to LLVM IR and its compilation time, the FTL is triggered only for long running portions of the code that are currently executing. There are two kinds of transitions in WebKit: the ones contained entirely inside the WebKit framework (i.e., transitions 1, 2 & 4 in Figure 3.3), and the ones that involve LLVM (i.e., 3 & 5 in Figure 3.3).

Transitions to and from LLVM are hard. There is no control over the stack layout or the optimized code produced by LLVM. In the case of transition 3, a different LLVM version is generated for each entry point that the framework desires to have inside this function. In WebKit, such entry points are located at loop headers. This choice makes sense with regard to the condition to enter the FTL, i.e., transition 3 is taken for long running portions of code that could be improved thanks to LLVM low level optimizations. WebKit has to generate a different version of the function for each entry points for two main reasons: LLVM allows only single entry points to functions (going around this limitation would require to modify LLVM IR and implementation), and instrumenting a function with several entry points would impact on the quality and performance of the generated native code by extending the code's length and restricting code motion.

Performing transition 3 requires to get the current state of execution and identify the entry point corresponding to the current instruction being executed. The DFG dumps its state into a scratch buffer. An LLVM function with the correct entry point is then generated, and instrumented such that its first block loads the content of the scratch buffer and correctly reconstructs the state. The mapping between the DFG IR nodes and the LLVM IR values is straight forward since both IR's are in SSA. A special data structure, called a Stackmap, enables to keep the mapping between LLVM values and registers/spill-slots.

Transition 5 is harder as it requires to extract the execution state from LLVM. WebKit

has two different mechanisms to enable OSR exits: the exit thunk and the invalidation points. In the first case, WebKit introduces exit branches at OSR exit points. The branch is guarded by an OSR exit condition and is a no-return tail call to a special function that takes all the live non-constant and not accounted for bytecode values as inputs. The second mechanism enables to remove the guard. Since we assume that the portion of code that is instrumented is executed a lot of times, the cost of testing the condition can have a great impact on the overall execution time. This mechanism relies on special LLVM intrinsics, namely patchpoints and stackmap shadow bytes. A patchpoint enables to reserve some extra space in the code, filled with nop sleds. When the WebKit framework detects that an exit should be taken, it overwrites the nop sleds with the correct function call to perform the OSR exit. This breaks the optimized version of the code which cannot be re-used later on and must be collected. The stackmap shadow bytes improves on this technique by allowing to directly overwrite the code, without having any nop sled generated before hand.

WebKit is a project that heavily, and successfully relies on OSR to improve performances. The web browser engine is used in Apple Web browser Safari and enables a net improvement of performances while proving to be reliable. Although successful, it does not provide a general and reusable framework for OSR in LLVM that other projects could benefit from.

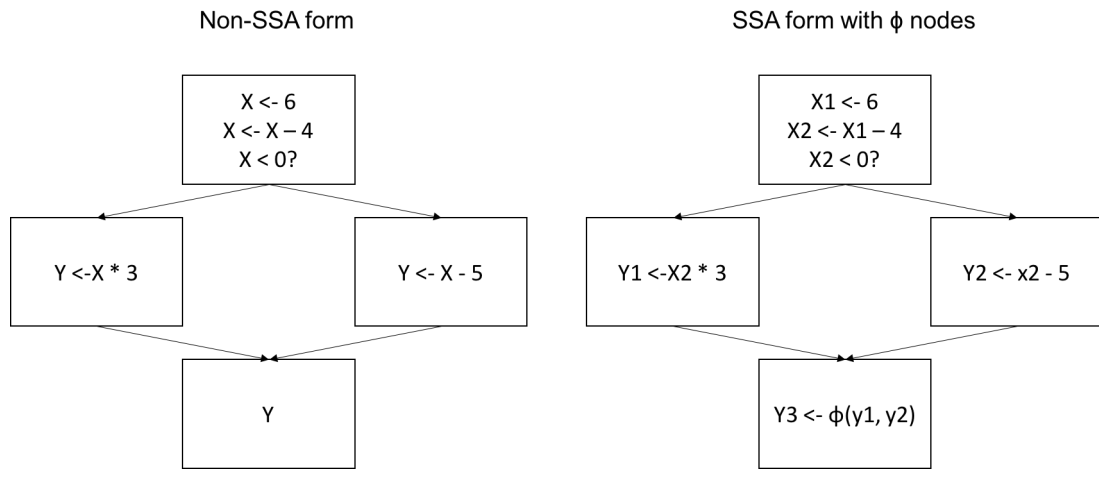
3.3 OSR & LLVM

3.3.1 What is LLVM?

LLVM[10, 16], formerly called Low-Level Virtual Machine, is a compiler infrastructure that provides a set of reusable libraries. LLVM provides the middle layers of a compiler system and a large set of optimizations for compile-time, link-time, run-time, and idle-time for arbitrary programming languages. These optimizations are performed on an intermediate representation (IR) of the code and yield an optimized IR. The LLVM framework also provides tools to convert and link code into machine dependent assembly code for a specific target platform. LLVM supports several instruction sets including ARM, MIPS, AMD TeraScale, and x86/x86-64[10].

The LLVM intermediary representation is a language-independent set of instructions that also provides a type system. The LLVM IR is in static single assignment (SSA) form, which requires every variable to be defined before being used, and assigned exactly once. SSA enables or improves several compiler optimizations among which constant propagation, value range propagation, sparse conditional constant propagation, dead code elimination, global value numbering, partial redundancy elimination, strength reduction, and register allocation[1]. The SSA requirement for variables to be assigned only once requires a special mechanism, called a ϕ -node, when a value depends on which control flow branch was executed before reaching the current variable definition. Figure 3.4 provides an example where we have to choose between two possible values for a variable after the merging of two control flow branches.

The LLVM IR type system provides basic types (e.g., integers, floats), and five derived types: pointers, arrays, vectors, structures, and functions. Any type construct can

FIGURE 3.4: Example of ϕ -node in SSA form.

then be represented as a combination of these types.

The LLVM framework is a versatile tool that enables to implement many programming languages paradigms. LLVM-based compilers exist for several mainstream/popular languages such as Java, Python, Objective-C, and Ruby. Other languages, like Haskell, Scala, Swift, Rust, Ada, and Fortran also have an LLVM compiler implementation. LLVM basic types enable to support object-oriented languages, such as Java and Python, dynamically typed languages like R, or statically typed like Scala. LLVM also enables to model functional languages such as Haskell, as well as imperative ones. Furthermore, it supports reflection and, thanks to dynamic linking, modular languages (e.g., Haskell). The tools provided enable static compilation as well as dynamic compilation techniques such as Just-In-Time compilation (JIT).

3.3.2 Why do we want OSR in LLVM?

On-stack replacement high-level mechanism is language-independent. Therefore, implementing OSR as a clean modular addition to LLVM would enable developers to leverage this feature in many programming languages, without requiring them to write a new compiler from scratch. Furthermore, as explained in 2.1.2, OSR is a useful tool for dynamic and adaptative optimizations. LLVM already provides implementations for many compiler optimizations[10] and tools to allow dynamic recompilation of code. Developers can therefore focus on language specific challenges, such as efficient profilers and new speculative systems, rather than on the optimizations and OSR implementations.

Implementing OSR for LLVM not only serves several languages, but also allows to provide a solution for several target platforms. As explained previously, LLVM supports several instruction sets corresponding to different architectures. By implementing OSR in LLVM, we get portability among these platforms for free.

3.3.3 OSR as an LLVM library: the McJIT OSR implementation

The McJIT OSR support[15] is an attempt at providing an OSR library compatible with the standard LLVM implementation. Lameed & Hendren claim to have come up with a clean modular, and re-usable technique completely defined at the LLVM IR level and compatible with the standard LLVM distribution. The implementation has been tested in the MvVM/McJIT[5, 17], an open source VM and JIT for MATLAB, built on LLVM. Their implementation answers to five challenges, listed in the paper[15], and reproduced here:

1. Identifying correct interrupt points and using the current LLVM IR to represent them.
2. Using the LLVM tools to correctly transform the IR while preserving a correct control flow graph.
3. Making a new version of a function available at the same address as the old one.
4. Providing a clean API for the OSR library, that is compatible with LLVM’s inlining capabilities.
5. Integrating OSR without modifying the LLVM installation.

Lameed and Hendren[15] claim to support optimization and re-optimization, as well as de-optimization by going back to the previous version of the function. Figure 3.5 shows that this feature only allows single-steps to be taken, i.e., the OSR library implemented in[15] does not seem to allow to skip intermediary versions.

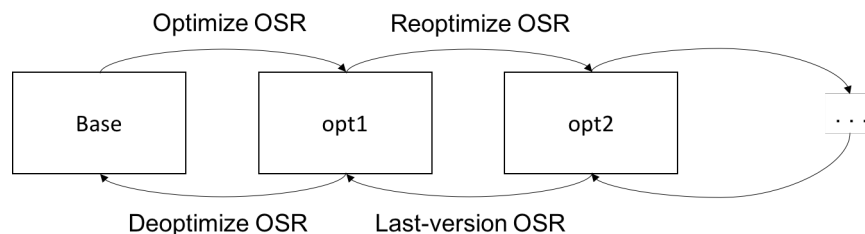


FIGURE 3.5: OSR classification from[15].

The McJIT library that provides OSR functionalities fit into the regular JIT infrastructure provided by LLVM as described in Figure 3.6 taken from[15]. The left figure is a normal JIT in LLVM. The LLVM CodeGen is the front-end of the compiler infrastructure that generates the LLVM IR. The LLVM Optimizer contains a collection of transformations and optimizations that run on LLVM IR. The Target CodeGen outputs the machine code corresponding to the LLVM IR input. The McJIT OSR API instruments the LLVM CodeGen to insert OSR points where the OSR transition can be triggered. A point is a call to the `genOSRSignal` function, which takes as arguments a pointer to a code transformer responsible for generating the new version of the function. The code transformer takes as arguments a pointer to the function that needs to be transformed, and a special OSR label to identify which OSR point triggered the call, if the function contains several ones. The OSR pass is responsible for instrumenting OSR points with



FIGURE 3.6: Retrofitting an existing JIT with OSR support from CITE

the correct OSR machinery. The instrumentation saves the live values and creates a descriptor that contains four elements:

1. A pointer to the current version of the function.
2. A pointer to the control version of the function, i.e., a copy of the old version of the function.
3. A variable mapping between the original version and the control version.
4. The set of live variables at the OSR points.



FIGURE 3.7: A CFG of a loop with no OSR point from [15].

Figures 3.7, 3.8, and 3.9 give an example of OSR instrumentation at a loop header. LH1 is the loop header. LB is the body of the loop and LE the loop exit. Figure 3.7 control flow graph (CFG) is the original CFG. Figure 3.8 is the resulting CFG after the Inserter is executed. Figure 3.9's CFG corresponds to the result of the OSR pass. The

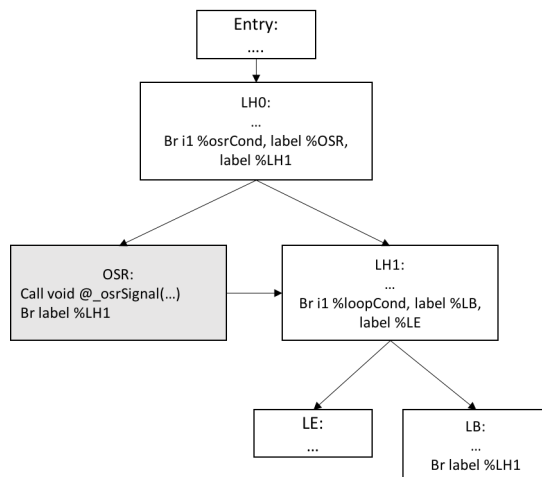


FIGURE 3.8: The CFG of the loop in Figure 3.7, after inserting an OSR point from[15].

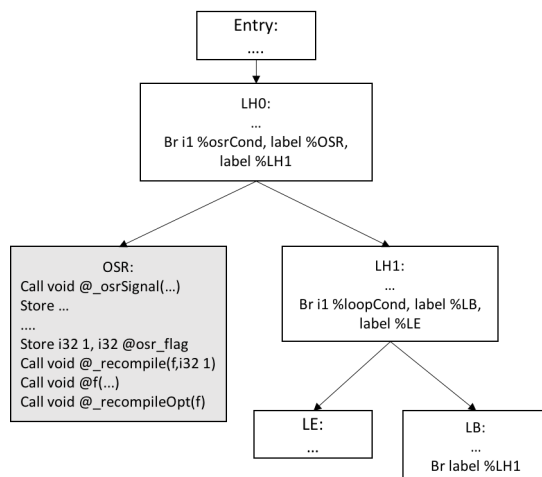


FIGURE 3.9: The transformed CFG of the loop in Figure 3.8 after the OSR Pass from[15].

recompile call in the OSR block recompiles f using the correct code transformer. Then f calls itself, executing the new version of the function. This works since the new version lives at the same address as the previous one and is instrumented to jump to the correct instruction, i.e., the one corresponding to the current point at which OSR was triggered. Figure 3.10 represents the CFG of f before the OSR instrumentation. Figure 3.11 shows the instrumentation of f that enables to jump to the correct instruction in the middle of the function. A prolog entry block is inserted at the function header. This block checks the *OSR flag* to know if an OSR transition is being performed. If that is the case, it branches to the prolog block that restores the state before resuming the execution at the correct instruction.

The OSR implementation proposed in[15] presents interesting features. The implementation is done entirely at the LLVM IR, hence making it language-independent. Furthermore, since the transformation function is provided by the user, any kind of

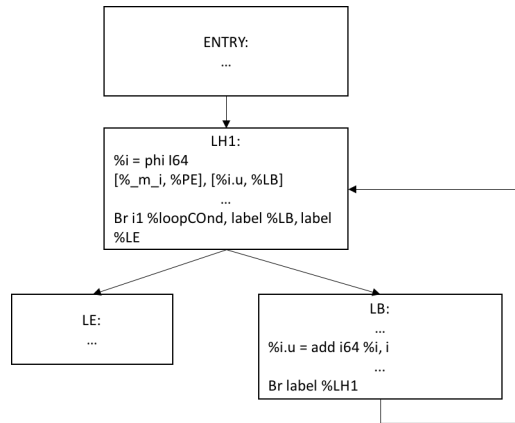


FIGURE 3.10: A CFG of a running function before inserting the blocks for state recovery, from[15].

language specific transformation can be used during the OSR. As a result, McJIT OSR support is a modular, and general purpose OSR library.

On the other hand, this library does not allow to have several versions of the same function live at the same time. This restriction can impede the overall performance of the program by extending the scope in which an assumption on which we base the transformation must hold. For example, a portion of the code, call it *A*, can trigger an OSR, while portion *B* is such that the assumption on which the optimization is based does not hold. The choice that the user has is to either optimize for *A*, and deoptimize for *B*, or to prevent *A* from optimizing by enlarging the scope in which the assumption is supposed to hold. None of these solutions is good if *A* and *B* are executed many times, one after the other. We will either lose a lot of execution time performing OSR, or keep executing *A* without optimization.

In the paper, the deoptimization process is not described. It is implied that it relies on the same set of tools provided by the framework as the optimization process, but no example is provided. As explained in 2.1.2, deoptimization is the most interesting feature of OSR, as it is required to preserve correctness. Being able to identify where a function should exit is a hard task. The McJIT library does not provide tools to ease this process and leaves the responsibility to the developer to instrument his functions correctly in order to exit to the correct landing pad.

To the best of our knowledge, MCJit OSR library is not currently used in production or any important project. As mentioned earlier, WebKit is efficiently integrated in Apple Safari’s web browser, which provides useful feedback on its performances. The lack of usage of McJIT OSR library prevents us from collecting performance results and assess its efficiency. Furthermore, the experimental evaluation in[15] relies on an example that seems artificial for our use of OSR. The case study presents a dynamic inliner that decides to inline a function call if the function is less than 20 basic blocks long, or if it is less than 50 basic blocks long and has an interpreter environment associated with its body. This requirement for inlining is one that can be checked statically and, hence, seems a little artificial.

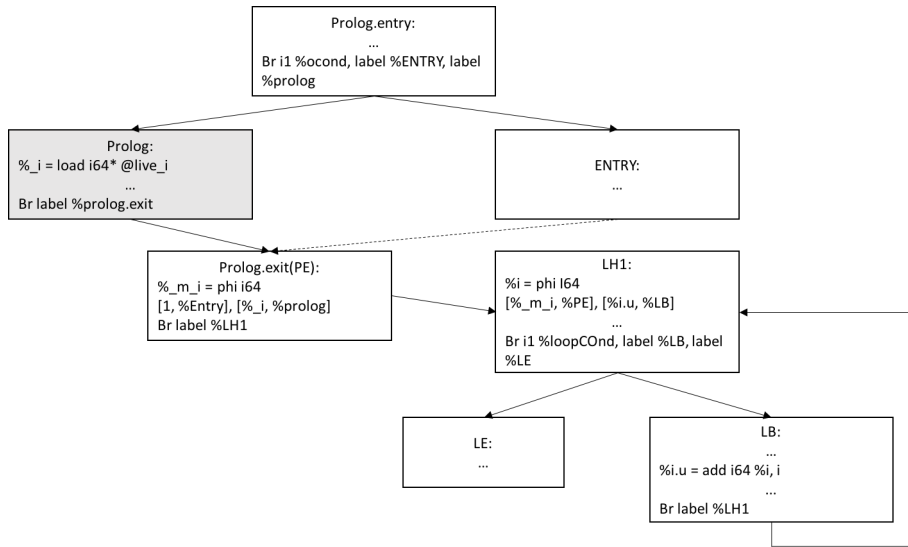


FIGURE 3.11: The CFG of the loop represented in Figure 3.10 after inserting the state recovery blocks, form[15].

	SELF Debugger	HotSpot	Graal	Jikes RVM	WebKit	McJIT OSR	OSR Kit
State Propagation technique							
Transition mechanism							
Modification of function signature							
Multiple entry points per function							
OSR point							
OSR point position							
Guarded							
Unguarded							
Multiple versions live							
Deoptimization support							
Exit target							
Performance results							

TABLE 3.1: Summary of OSR implementations and their features.

3.4 A classification summary

This section presents a summary of the various OSR implementations presented in this Chapter. The summary focuses the state propagation technique, the types of OSR points supported, the implementation of the transition mechanism, the multi-version support, and the deoptimization specifications.

Chapter 4

Deoptimization in RJIT

This chapter presents our design to support OSR transitions and OSR deoptimizations in LLVM, and more particularly in the RJIT project. Our design and implementation is based on OSR Kit[6, 7], an OSR library for LLVM, implemented by D’Elia and Demetrescu[6] at the department of Computer, Control, and Management Engineering of the Sapienza University of Rome. This implementation provides several features of existing techniques, that were never simultaneously supported in a single framework. It is flexible enough to be used to implement our own model for OSR deoptimization. We therefore decided to use and extend it, in order to provide OSR deoptimization for R.

4.1 Overview

This section presents our goals for this Master Thesis Project and justifies the use of the OSR Kit as the base of our OSR support in RJIT.

4.1.1 Justification

The goal of this Master Thesis project is to provide a flexible OSR deoptimization framework in LLVM, and use it to improve performances in RJIT, our LLVM JIT compiler for R. R is a programming language and software environment for the statistical computing and graphics, developed by the R Foundation for Statistical Computing[23]. Due to SAY WHYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYYY OOOOOHHHHH GOOOOOOODDDDDDD WHHHYYYYYYYYYYYY, R exhibits very poor performances CITE SOMETHING. The RJIT project strives to improve these performances by providing a LLVM based JIT compiler for R. SAY MORE. The RJIT compiler is still pretty young, only a few months old. As a result, we lack FEEDBACK; DONT KNOW EXACTLY HOW TO IMPROVE PERF AND NEED TO EXPERIMENT. Therefore, we are looking for a flexible and extensible OSR mechanism that enables us to prototype and experiment various solutions, without trapping ourselves into a single model.

The OSR Kit library[6] is a flexible implementation of on-stack replacement instrumentation in LLVM. The source-code for the library is available on Github[7], and the library can be used in any LLVM project by simply copy-pasting the OSR Kit files inside of it. The simple integration, the availability of the source code, and the flexibility of the framework make it a perfect base implementation upon which we can implement our support for OSR deoptimization mechanisms in RJIT.

This master project thesis focuses on the design and implementation of a prototype for OSR deoptimization support in R_{JIT}. Starting a new OSR transition implementation from scratch requires time. Using a flexible and modifiable OSR transition library therefore seemed like the goto option. We do not waste time reimplementing something that already exists, and can therefore put all our efforts into implementing an interesting OSR deoptimization case, testing it, and extending the OSR Kit library with mechanisms that are specific to our needs (Sections 4.4 and 4.5).

4.1.2 An LLVM Library

The OSR Kit is a general-purpose, target-independent, implementation of on-stack replacement for LLVM. As such, it can be used by any LLVM based compiler. The main goals of the OSR Kit project are:

1. To allow chained OSR transitions, i.e., a continuation function can be instrumented to allow OSR transitions.
2. To support OSR entries and exits via the same instrumentation.
3. To allow transitions at arbitrary function locations.
4. To allow continuation functions to be either generated at run time, or already known at compilation time (i.e., generated on the fly or cached).
5. To encapsulate and hide the OSR implementation details from the front-end.
6. To encode OSR transitions entirely at the LLVM IR level.
7. To limit the intrusiveness of the OSR instrumentation.
8. To allow the LLVM's compilation pipeline to generate the most efficient native code for an instrumented function.

The OSR Kit library's main contribution is the implementation of an efficient and flexible OSR transition mechanism. This transition can be used to implement OSR entries and OSR exits alike. Allowing OSR points to be inserted at any location enables the user to experiment the OSR mechanisms at any instruction boundary. In other words, the OSR library is designed to give the user as much freedom as possible. Section 2.4.1 explained the basic trade-offs between generating the continuation function on the fly and caching already compiled versions. Where other implementations, such as McJIT OSR[15], made a clear choice between the two techniques, the OSR Kit library decided to enable both of them, hence allowing the user to experiment and select the design that better suits specific use cases.

LLVM IR is a stable representation that combines the advantages of both the high-level representation, i.e., it still contains some semantic constructions particular to the language being compiled, and the advantages of a lower level representation, closer to the execution engine. In the case of this master thesis project, i.e., providing OSR mechanisms in the R_{JIT} project, the LLVM IR is the exact middle layer representation that we need. The R AST representation is too high level and prevents us from introducing efficient OSR instrumentation. At the LLVM IR level, the R semantics are still visible and it therefore allows us to efficiently implement our optimizations.

Point 8 is important in order to leverage the full power of the LLVM transformation and optimization passes. One of the main advantages of using the LLVM framework is that LLVM transformation passes and optimizations can be automatically added to the compilation pipeline in order to improve the quality of the code generated. For example, in the case study in Chapter 5, one of the OSR inlining mechanism's goals is to extend the scope in which the LLVM passes run in order to improve the code quality.

4.2 Resolved & Open OSR

The OSR Kit library[6] provides the choice between generating the continuation function of an OSR transition on the fly, or using an already compiled function. This section describes both techniques as they are implemented in the OSR Kit library.

4.2.1 The Open OSR

The *open OSR* scenario corresponds to the case where the continuation function is generated when the OSR transition is fired. Deferring the compilation of the continuation function allows profiled-guided compilation strategies to gather as much information as possible about the current state of the execution, and therefore generate the best possible code.

The open OSR scenario is implemented as a call to a stub function, call it f_{stub} . The f_{stub} function is responsible for generating the continuation function, and then transferring the execution to it. The continuation function is generated by a special *gen* function, that takes as inputs the base function f , i.e., the from function, and the OSR point that triggered the transition.

CODE API

Open OSR are hard in the deoptimization case. Being able to generate a continuation function on the fly for an OSR exit requires the framework to be able to undo an optimization. The RJIT framework, does not, for the moment, provide means to keep track of transformations performed on a function, and efficiently revert them. For example, the VARMAP data structure in the Jikes RVM OSR implementation[21] is automatically updated by the framework when a transformation is performed on the function. This enables to reverse the steps that lead to the optimized version. While the RJIT compiler does not provide such mechanism, we could imagine keeping the LLVM IR of the original function in memory, without compiling it. When an open OSR exit is triggered, the unoptimized LLVM IR for the function is retrieved, the continuation point is identified, and RJIT instruments and compiles the correct continuation function. This solution enables to avoid the cost of going through the entire compilation of the continuation function until it is truly needed. On the other hand, it increases the amount of work performed during the OSR transition.

Similarities with McJIT

The open OSR is similar, in theory, to the McJIT OSR implementation[15]. Both of them rely on a *gen* function to generate the continuation function when the OSR transition is fired. However, D’Elia and Demetrescu[6] made the choice to perform these operations, i.e., generating the continuation function and transferring the execution, inside a separated function, rather than instrumenting the from function directly. This choice was made in order to minimize the code injected inside the from function. In fact, the instrumentation, even when inserted in a separated branch, might interfere with compiler optimizations. Examples of such interferences are the increase of register pressure, the alteration of code layout and instruction cache behavior.

4.2.2 The Resolved OSR

The *resolved OSR* scenario corresponds to the case where the continuation function is known prior to fire the OSR transition, and the from function has been instrumented to transfer the execution to this specific continuation function. In the resolved OSR scenario, the transition is expected to be faster than in the open scenario, since there is no compilation to be performed. This technique also enables to reuse functions that were compiled previously, e.g., either because several from functions have the same continuation function, or because an OSR point is fired several times during the execution of the program.

The resolved OSR is implemented as a call to the continuation function, which takes as inputs all the live variables at the OSR point as inputs. The continuation function is instrumented to jump at the correct location inside the code to resume the execution. Figure REF illustrates the resolved OSR scenario.

FIGURE

CODE API

The resolved OSR scenario does not require to inverse on the fly the transformations performed on the function and can therefore easily be used to implement the deoptimization case. In order to obtain the optimized version of a function, one has to perform transformations on a base function. This base function can therefore be used as the continuation function in the case of an OSR exit. In RJIT, we rely on the resolved OSR to provide deoptimization of functions where call sites were inlined. While generating the optimized version by speculatively inlining calls, the RJIT compiler introduces the OSR exit instrumentation with the base function as continuation function. When several call sites are inlined, each of them is instrumented with an OSR exit that points to its own continuation function, which is a copy of the base function with the proper OSR instrumentation (see Section 4.3.4).

4.3 The API

4.3.1 StateMaps

4.3.2 OSR Points

4.3.3 OSR Conditions

4.3.4 The continuation function

4.4 Using OSR Kit to deoptimize

4.4.1 The unease of deoptimizing

4.4.2 SOMETHING

4.4.3 SOMETHING ELSE

4.5 Extending OSR Kit

4.5.1 Handling versions

4.5.2 Unguarded OSR points

4.5.3 Fusing versions?

Chapter 5

Case Study

5.1 The RJIT project

5.2 OSR in RJIT: requirements

5.3 Implementation

5.3.1 The Function Call in RJIT

5.3.2 Challenges

5.4 OSR Conditions

Chapter 6

Conclusion

Appendix A

Appendix Title Here

Write your Appendix content here.

Bibliography

- [1] Andrew W. Appel. “SSA is functional programming”. In: *SIGPLAN notices* 33.4 (1998), pp. 17–20.
- [2] Matthew Arnold et al. “Adaptive optimization in the Jalapeno JVM”. In: *ACM SIGPLAN Notices*. Vol. 35. 10. ACM. 2000, pp. 47–65.
- [3] Craig Chambers. “The design and implementation of the self compiler, an optimizing compiler for object-oriented programming languages”. PhD thesis. Stanford University, 1992.
- [4] Craig Chambers and David Ungar. “Making pure object-oriented languages practical”. In: *ACM SIGPLAN Notices*. Vol. 26. 11. ACM. 1991, pp. 1–15.
- [5] Maxime Chevalier-Boisvert, Laurie Hendren, and Clark Verbrugge. “Optimizing MATLAB through just-in-time specialization”. In: *Compiler Construction*. Springer. 2010, pp. 46–65.
- [6] Daniele Cono D’Elia and Camil Demetrescu. “Flexible On-Stack Replacement in LLVM”. In: *Proceedings of the 2016 International Symposium on Code Generation and Optimization (CGO 2016)*. 2016.
- [7] Daniele Cono D’Elia and Camil Demetrescu. *OSRKit OSRKit*. 2016. URL: <https://github.com/dcdelia/tinyvm> (visited on 02/01/2016).
- [8] Gilles Duboscq, Thomas Würthinger, and Hanspeter Mössenböck. “Speculation without regret: reducing deoptimization meta-data in the Graal compiler”. In: *Proceedings of the 2014 International Conference on Principles and Practices of Programming on the Java platform: Virtual machines, Languages, and Tools*. ACM. 2014, pp. 187–193.
- [9] Stephen J Fink and Feng Qian. “Design, implementation and evaluation of adaptive recompilation with on-stack replacement”. In: *Proceedings of the international symposium on Code generation and optimization: feedback-directed and runtime optimization*. IEEE Computer Society. 2003, pp. 241–252.
- [10] LLVM Developer Group. *LLVM web site LLVM web site*. 2003. URL: <http://llvm.org/> (visited on 02/01/2016).
- [11] Urs Hölzle, Craig Chambers, and David Ungar. “Debugging optimized code with dynamic deoptimization”. In: *ACM Sigplan Notices*. Vol. 27. 7. ACM. 1992, pp. 32–43.
- [12] Urs Hölzle, Craig Chambers, and David Ungar. “Optimizing dynamically-typed object-oriented languages with polymorphic inline caches”. In: *ECOOP’91 European Conference on Object-Oriented Programming*. Springer. 1991, pp. 21–38.
- [13] Urs Hölzle and David Ungar. “A third-generation SELF implementation: reconciling responsiveness with performance”. In: *ACM SIGPLAN Notices*. Vol. 29. 10. ACM. 1994, pp. 229–243.
- [14] Apple Inc. *WebKit The WebKit Open Source Web Browser Engine*. 1998. URL: <https://webkit.org/> (visited on 02/01/2015).

- [15] Nurudeen A Lameed and Laurie J Hendren. "A modular approach to on-stack replacement in LLVM". In: *ACM SIGPLAN Notices*. Vol. 48. 7. ACM. 2013, pp. 143–154.
- [16] Chris Lattner and Vikram Adve. "LLVM: A compilation framework for lifelong program analysis & transformation". In: *Code Generation and Optimization, 2004. CGO 2004. International Symposium on*. IEEE. 2004, pp. 75–86.
- [17] McLAB. *The McVM Virtual Machine and its JIT Compiler The McVM Virtual Machine and its JIT Compiler*. 2012. URL: http://www.sable.mcgill.ca/mclab/mcvm_mcjit.html (visited on 02/01/2016).
- [18] Jikes RVM Project Organization. *Jikes RVM Jikes RVM*. 2013. URL: www.jikesrvm.org (visited on 02/01/2016).
- [19] Michael Paleczny, Christopher Vick, and Cliff Click. "The java hotspot TM server compiler". In: *Proceedings of the 2001 Symposium on Java TM Virtual Machine Research and Technology Symposium-Volume 1*. USENIX Association. 2001, pp. 1–1.
- [20] Filip Pizlo. *WebKit FTL WebKit FTL*. 2014. URL: <https://webkit.org/blog/3362/introducing-the-webkit-ftl-jit/> (visited on 02/01/2016).
- [21] Sunil Soman and Chandra Krintz. "Efficient and General On-Stack Replacement for Aggressive Program Specialization." In: *Software Engineering Research and Practice*. 2006, pp. 925–932.
- [22] Toshio Sukanuma et al. "A dynamic optimization framework for a Java just-in-time compiler". In: *ACM SIGPLAN Notices*. Vol. 36. 11. ACM. 2001, pp. 180–195.
- [23] R Core Team. *The R project The R project*. 1993. URL: www.r-project.org (visited on 02/01/2016).