

Bachelor Thesis

Profile Caching for the Java Virtual Machine

Marcel Mohler

Zoltán Majó
Tobias Hartmann
Oracle Cooperation

Prof. Thomas R. Gross
Laboratory for Software Technology
ETH Zurich

August 2015



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



Laboratory for Software Technology

Introduction

Virtual Machines like the Java Virtual Machine (JVM) are used as the execution environment of choice for many modern programming languages. The VMs interpret a suitable intermediate language (e.g., Java Byte Code for the JVM) and provide the runtime system for application programs and usually include a garbage collector, a thread scheduler, interfaces to the host operating system. As interpretation of intermediate code is time-consuming, VMs include usually a Just-in-Time (JIT) compiler that translates frequently-executed functions or methods to “native” code (e.g., x86 instructions).

The JIT compiler executes in parallel to a program’s interpretation by the VM, and as a result, compilation speed is a critical issue in the design of a JIT compiler. Unfortunately, it is difficult to design a compiler such that the compiler produces good (or excellent) code while limiting the resource demands of this compiler (the compiler requires storage and cycles – and even on a multi-core processor, compilation may slow down the execution of the application program). Consequently, most VMs adopt a multi-tier compilation system. At program startup, all methods are interpreted by the VM (execution at Tier-0). The interpreter performs profiling, and if a method is determined to be “hot”, this method is then compiled by the Tier-1 compiler. Methods compiled to Tier 1 are then profiled further and based on these profiling information, some methods are eventually compiled at Tier 2. One of the drawbacks of this setup is that for all programs, all methods start in Tier 0, with interpretation and profiling by the VM. However, for many programs the set of “hot” methods does not change from one execution to another and there is no reason to gather again and again the profiling information.

The main idea of this thesis is to cache these profiles from a prior execution to be used in further runs of the same program. This would allow the JIT compiler to use more sophisticated profiles early in program execution and avoid gathering the same profiling as well as prevent further compilations when more information about the method is available. I present an implementation on top of the Java Hotspot Virtual Machine as well as profound performance analysis using state-of-the-art benchmarks.

Contents

1	Motivation	1
1.1	Tiered Compilation in Hotspot	1
1.2	On Stack Replacement	2
1.3	Deoptimizations	3
1.4	Compile Thresholds	3
1.5	Examples	3
2	Implementation / Design	7
3	Performance	9
3.1	Examples	9
3.2	SPECjvm 2008	9
3.3	Nashorn / Octane	9
4	Conclusion	11
A	Appendix	13
A.1	Tiered Compilation Thresholds	14
	Bibliography	14

1 Motivation

1.1 Tiered Compilation in Hotspot

As mentioned in the introduction, Programming Language Virtual Machines like Java Hotspot feature a multi-tier system when compiling methods during execution. Java VM's typically use Java Bytecode as input, a platform independent intermediate code generated by a Java Compiler like `javac`. The Bytecode is meant to be interpreted by the virtual machine or further compiled into platform dependend machine code. Hotspot includes one interpreter and two different compilers with different profiling levels resulting in a total of 5 different levels. The following Figure 1.1 gives a short overview as well as showing the standard transitions.

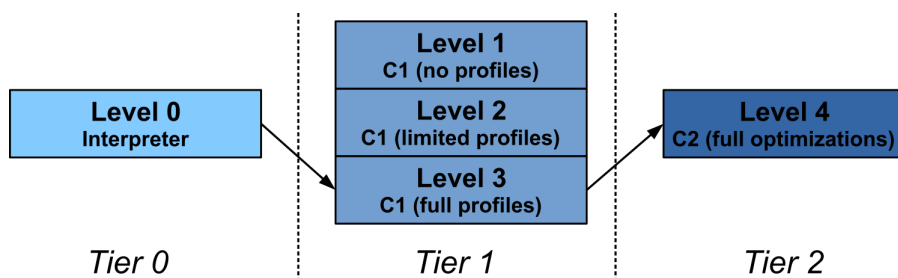


Figure 1.1: Overview over compilation tiers

All methods start being executed by Tier-0 also called the Interpreter. The interpreter is template-based meaning for each bytecode instruction it emits a predefined assembly code snippet. During execution this code is also profiled. This means method execution counters, loop back-branches and additional statistics are counted. More importantly information about the program flow and state are gathered. These information contain for example which branches get taken or the final types of dynamically typed objects. Once one these counters exceed a predefined, constant threshold the method is considered *hot* which usually results in a compilation at a higher tier.

The standard behavior of Hotspot is to proceed with Level 3 (Tier 1). This means the method gets compiled with C1, also referred to as *client* compiler, and continues gathering full profiles. C1's goal is to provide a fast compilation with a low memory footprint. The client compiler performs simple optimizations such as constant folding, null check elimination and method inlining. More information about C1 can be found in [4] and [2]. The levels 1 and 2 include the same optimization but offer no or less profiling information and are used in special cases, for example if the compiler

is already very busy or enough profiles are already available.

Eventually, when further compile thresholds are exceeded, the JVM further compiles the method with C2, also known as *server* compiler. The server compiler makes use of the gathered profiles in Tier 0 and Tier 1 and produces highly optimized code. C2 includes far more optimizations like loop unrolling, common subexpression elimination and elimination of range and null checks. It performs optimistic method inlining, for example by converting some virtual calls to static calls. A more detailed look at the server compiler can be found in [3].

The naming scheme *client/server* comes from back in the days where tiered compilation was not available and one had to choose the compiler via a Hotspot command line flag. The *Client* compiler was meant to be used for interactive client programs with graphical user interfaces where response time is more important than peak performance. For long running server applications, the highly optimized but slower server compiler was used.

Tiered compilation was introduced to improve start-up performance of the JVM. Starting with the interpreter means that there is zero wait time until the method is executed since one does not need to wait until a compilation is finished. C1 allows the JVM to have more optimized code available early which then can be used to create a richer profile to be used when compiling with C2. Ideally this profile already contains most of the program flow so less deoptimizations (see 1.3 occur.

1.2 On Stack Replacement

Since the JVM does not only count method invocations but also loop back branches (see also Section 1.4) it can happen that a method gets compiled while it is still running and the compiled method is ready before the method has finished. Instead of waiting for the next method invocation Hotspot can replace the method directly on the stack (see Figure 1.2) and is called *on stack replacement*. An example would be a simple method consisting of a very long running loop.

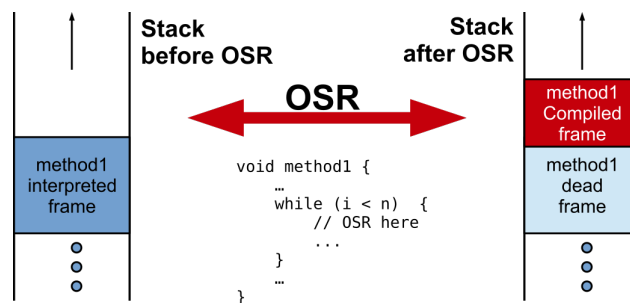


Figure 1.2: Overview over compilation tiers

1.3 Deoptimizations

Ideally we compile a method with as much profiling information as possible. Since the profiling information are usually gathered in levels 0 and 3 it can happen that a method compiled by C2 wants to execute a branch it never used before. In this case the information about this branch are not available in the profile and therefore have not been compiled into the C2-compiled code. This is done to allow further, very optimistic optimization and to keep the compiled code smaller. So instead, the compiler places an uncommon trap at unused branches or unloaded classes which will get triggered in case they actually get used at a later time in execution.

The JVM then stops execution of that method and returns the control back to the interpreter. This process is called *deoptimization* and considered very costly. The previous interpreter state has to be restored and eventually the method might get recompiled with the newly gained information.

1.4 Compile Thresholds

The transitions between the compilation levels (see Fig. 1.1) are chosen based on predefined constants called *compile thresholds*. When running an instance of the JVM one can specify them manually or use the ones provided. A list of thresholds and their default values relevant to this thesis are given in Appendix A.1. The standard transitions from Level 0 to 3 and 3 to 4 happen when the following predicate returns true:

$$i > TierXInvocationThreshold * s \\ || (i > TierXMinInvocationThreshold * s \ \&\& \ i + b > TierXCompileThreshold * s)$$

where X is the next compile level (3 or 4), i the number of method invocations, b the number of backedges and s a scaling coefficient (default = 1). The thresholds are relative and individual for interpreter and compiler.

On Stack Replacement uses a simpler predicate:

$$b > TierXBackEdgeThreshold * s$$

Please note that there are further conditions influencing the compilation like the load on the compiler which will not be discussed further.

1.5 Examples

I continue with presenting two very simple examples that illustrate the usage and benefit from using cached profiles. To start I consider a standard Java Hotspot execution with OnStackReplacement disabled.

Listing 1.1: Example usage of the listing package

```

1 class NoCompile {
2     double result = 0.0;
3     for(int c = 0; c < 100; c++) {
4         result = method1(result);
5     }
6     public static double method1(double count) {
7         for(int k = 0; k < 10000000; k++) {
8             count = count + 50000;
9         }
10        return count;
11    }
12 }

```

Example one is a simple class that invokes a method one hundred times. The method itself consists of a long running loop. The source code is shown in Listing 1.1. Since OSR is disabled and a compilation to level 3 is triggered after 200 invocations this method never leaves the interpreter. I call this run the *Baseline*. To show the influence of cached profiles I use a compiler flag to lower the compile threshold explicitly and, using the functionality written for this thesis, tell Hotspot to cache the profile. In a next execution I use these profiles and achieve significantly better performance as you can see in Figure 1.3.

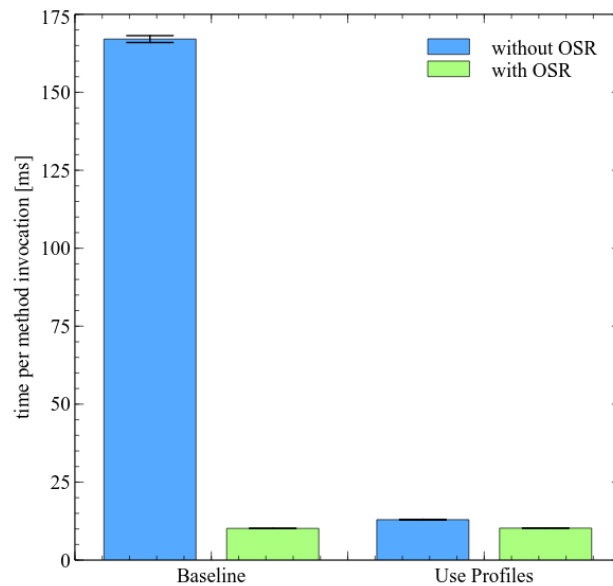


Figure 1.3: NoCompile.method1 - per method invocation time comparison

Enabling OSR again and the difference between with and without cached profiles vanishes. This happens because Hotspot quickly realizes the hotness of the method and compiles it during the first invocation already.

Listing 1.2: Example usage of the listing package

```

1 class ManyDeopts {
2     double result = 0.0;
3     for(int c = 0; c < 100; c++) {
4         result = method1(result);
5     }
6     public static long method1(long count) {
7         for(int k = 0l; k < 100000000l; k++) {
8             if (count < 100000000l) {
9                 count = count + 1;
10            } else if (count < 300000000l) {
11                count = count + 2;
12            }
13            .
14            .
15            .
16            } else if (count < 505000000000l) {
17                count = count + 100;
18            }
19            count = count + 50000;
20        }
21        return count;
22    }

```

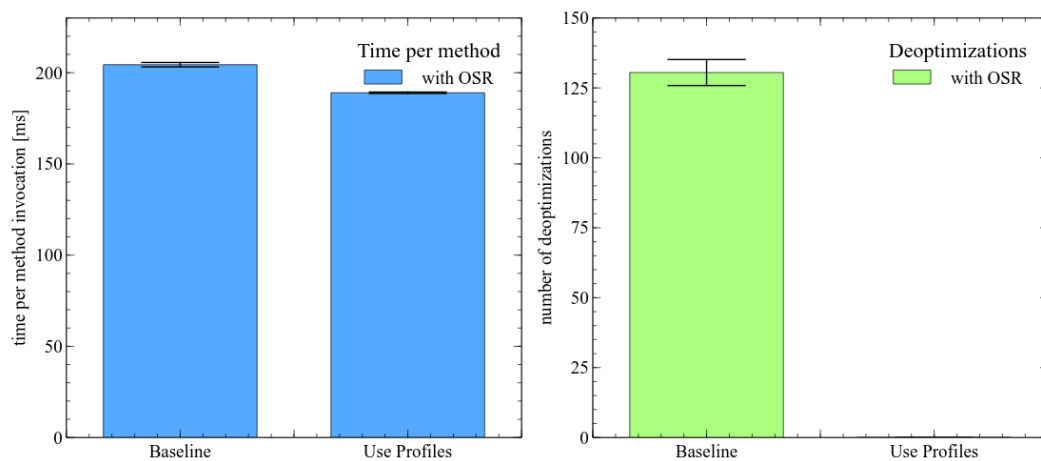


Figure 1.4: ManyDeopt.method1 - per method invocation time and deoptimization count comparison

2 Implementation / Design

3 Performance

3.1 Examples

3.2 SPECjvm 2008

3.3 Nashorn / Octane

4 Conclusion

A Appendix

A.1 Tiered Compilation Thresholds

flag	description	default
CompileThresholdScaling	number of interpreted method invocations before (re-)compiling	1.0
Tier0InvokeNotifyFreqLog	Interpreter (tier 0) invocation notification frequency	7
Tier2InvokeNotifyFreqLog	C1 without MDO (tier 2) invocation notification frequency	11
Tier3InvokeNotifyFreqLog	C1 with MDO profiling (tier 3) invocation notification frequency	10
Tier23InlineeNotifyFreqLog	Inlinee invocation (tiers 2 and 3) notification frequency	20
Tier0BackedgeNotifyFreqLog	Interpreter (tier 0) invocation notification frequency	10
Tier2BackedgeNotifyFreqLog	C1 without MDO (tier 2) invocation notification frequency	14
Tier3BackedgeNotifyFreqLog	C1 with MDO profiling (tier 3) invocation notification frequency	13
Tier2CompileThreshold	threshold at which tier 2 compilation is invoked	0
Tier2BackEdgeThreshold	Back edge threshold at which tier 2 compilation is invoked	0
Tier3InvocationThreshold	Compile if number of method invocations crosses this threshold	200
Tier3MinInvocationThreshold	Minimum invocation to compile at tier 3	100
Tier3CompileThreshold	Threshold at which tier 3 compilation is invoked (invocation minimum must be satisfied)	2000
Tier3BackEdgeThreshold	Back edge threshold at which tier 3 OSR compilation is invoked	60000
Tier4InvocationThreshold	Compile if number of method invocations crosses this threshold	5000
Tier4MinInvocationThreshold	Minimum invocation to compile at tier 4	600
Tier4CompileThreshold	Threshold at which tier 4 compilation is invoked (invocation minimum must be satisfied)	15000
Tier4BackEdgeThreshold	Back edge threshold at which tier 4 OSR compilation is invoked	40000

Bibliography

- [1] T. Hartmann, A. Noll, and T. R. Gross. Efficient code management for dynamic multi-tiered compilation systems. In *2014 International Conference on Principles and Practices of Programming on the Java Platform Virtual Machines, Languages and Tools, PPPJ '14, Cracow, Poland, September 23-26, 2014*, pages 51–62, 2014.
- [2] V. Ivanov. JIT-compiler in JVM seen by a Java developer. <http://www.stanford.edu/class/cs343/resources/java-hotspot.pdf>, 2013.
- [3] M. Paleczny, C. Vick, and C. Click. The Java HotSpotTMServer Compiler. https://www.usenix.org/legacy/events/jvm01/full_papers/paleczny/paleczny.pdf, 2001. Paper from JVM '01.
- [4] T. Rodriguez and K. Russell. Client Compiler for the Java HotSpotTMVirtual Machine: Technology and Application. <http://www.oracle.com/technetwork/java/javase/tech/3198-d1-150056.pdf>, 2002. Talk from JavaOne 2002.