

Bachelor Thesis

Profile Caching for the Java Virtual Machine

Marcel Mohler
ETH Zurich

Zoltán Majó
Tobias Hartmann
Responsible engineers

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

The JIT compiler executes in parallel to the interpreted program 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 that produces good or excellent code while limiting the resource demands of this compiler. The compiler requires storage, CPU 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 virtual machine (execution at Tier 0). The interpreter gathers execution statistics called *profiles* and if a method is determined to be executed frequently, this method is 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 higher tiers. One of the drawbacks of this setup is that for all programs, all methods start at Tier 0, with interpretation and profiling by the VM. However, for many programs the set of the most used methods does not change from one execution to another and there is no reason to gather profiling information again.

The main idea of this thesis is to cache profiles from a prior execution to be used in further runs of the same program. Having these *cached profiles* available avoids the JIT compiler to gather the same profiling information again. As well as allow the compiler to use more sophisticated profiles early in program execution and prevent recompilations when more information about the method is available. While this in general should not significantly influence the peak performance of the program, the hope is to decrease the time the JVM needs to achieve it, the so called *warmup*.

This thesis proposes a design and an implementation of a profile caching feature for *HotSpot*, an open source Java virtual machine maintained and distributed by Oracle Corporation as well as a profound performance analysis using state-of-the-art benchmarks.

Contents

1 Overview of HotSpot	1
1.1 Tiered compilation	1
1.2 Deoptimizations	3
1.3 On-Stack replacement	4
1.4 Compile thresholds	4
2 Motivation	7
2.1 Example 1: Benefit of early compilation	8
2.2 Example 2: Benefit of fewer deoptimizations	9
2.3 Similar systems	10
3 Implementation / Design	13
3.1 Creating cached profiles	13
3.2 Initializing cached profiles	15
3.3 Using cached profiles	16
3.4 Different usage modes for cached profiles	17
3.4.1 Compile thresholds lowered (Mode 0)	18
3.4.2 Unmodified compile thresholds (Mode 1)	18
3.4.3 Modified C1 stage (Mode 2)	19
3.5 Problems	19
3.6 Debug output	20
4 Performance	21
4.1 Setup	21
4.2 Benchmark performance	22
4.2.1 SPECjvm warmup performance	22
4.2.2 Octane performance	24
4.3 Deoptimizations	24
4.4 Effect on compile queue	30
4.5 Number and type of compilations	36
4.6 Effect of interpreter profiles	39
4.7 Effect of intrinsified methods	41

5 Possible improvements	43
6 Conclusion	45
A Appendix	47
A.1 Tiered compilation thresholds	48
A.2 Cached profile example	49
A.3 Code changes	50
A.4 SPECjvm benchmark	50
A.5 Octane benchmark	52
A.6 Additional graphs	54
Bibliography	54

1 Overview of HotSpot

This chapter will provide the reader with an overview of the relevant parts of Java HotSpot. The chapter explains the core concepts that are needed to understand the motivation of this thesis and the implementation of the system described in this thesis.

1.1 Tiered compilation

As mentioned in the introduction, virtual machines (VMs) 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` [10]. The bytecode is meant to be interpreted by the virtual machine or further compiled into platform dependent machine code (e.g., x86 instructions). HotSpot includes one interpreter and two different just-in-time compilers with different profiling levels resulting in a total of 5 different *compilation tiers*. Since in literature and the JVM source code the *tiers* are also called *compilation levels* the terms will be used synonymously.

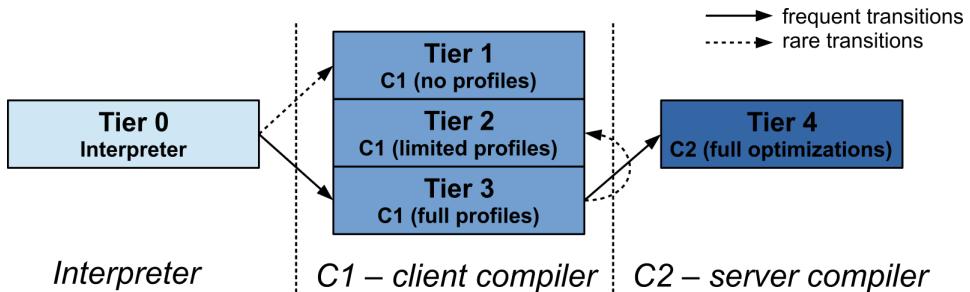


Figure 1.1: Overview of compilation tiers

All methods start being executed at Tier 0, which denotes the interpreter. The interpreter performs a template-based replacement, that is, for each bytecode instruction the interpreter emits a pre-defined assembly code snippet. During execution, the assembly code is also profiled. The snippets also contain structures to gather method information like execution counters or loop back-branches. If a counter exceeds a predefined threshold, the method is considered *hot* and a call back to the JVM is initiated that usually results in a compilation at a higher tier.

The standard behavior of HotSpot is to proceed with Level 3 (Tier 3). The method gets compiled with C1, also referred to as *client compiler*. C1's goal is to provide a fast compilation with a low

Listing 1.1: Example that shows potential compilation based on profiling information

```

1 public static void m(int i) {
2     if (i == 0) { // very common branch (a)
3         Math.sin(0);
4     } else { // very uncommon branch (b)
5         Math.sin(pi + i)
6     }
7 }
8 /**
9 // If the JVM realizes based on profiling information,
10 // that branch (a) is taken all the time:
11 // compiler could compile the method as follows:
12 /**
13 public static void m(int i) {
14     if (i != 0) // very common branch (a)
15         // UNCOMMON TRAP, call to JVM
16     return 0; // result of sin(0)
17 }
```

memory footprint. The client compiler performs simple optimizations such as constant folding, null check elimination, and method inlining based on the information gathered during interpretation. Most of the classes and methods have already been used in the interpreter, allowing C1 to inline them to avoid costly invocations. 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. For example, if certain branches were not taken during execution, further compilations might abstain from compiling these branches and replace them with static code to provide a faster method execution time (see the example in Listing 1.1). The uncommon branch includes an *uncommon trap* which notifies the JVM that an assumption does not hold anymore. This then leads to so called *deoptimizations* which are further explained in the separate Section 1.2.

Level 1 and Level 2 include the same optimizations but offer no or less profiling information and are used in special cases. Code compiled at these levels is significantly faster than Level 3 because it needs to execute none or little instructions creating and managing the profiles. Since the profiles generated by C1 are further used in C2, HotSpot is usually interested in creating full profiles and therefore uses Level 3. There are, however, rare instances where a compilation of Level 1 or Level 2 is triggered. For example, if enough profiles are available and a method can not be compiled by a higher tier, HotSpot might recompile the method with Tier 2 to get faster code until the higher tier compiler is available again. A compiler can become unavailable if its compilation queue exceeds a certain threshold. Also, if a method is simple, it may be compiled at Level 1.

More information about C1 can be found in [13] and [6].

Eventually, when further compile thresholds are exceeded, the JVM compiles the method with C2, also known as the *server* compiler. The server compiler uses the profiles gathered in Tier 0 and Tier 3 and produces highly optimized code. C2 includes far more and more complex 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. It relies heavily on the profiling information and richer profiles allow the compiler to do more and better optimizations. While the code quality of C2 is a lot better than C1 this comes at the cost of compile time. A more detailed look at the server compiler can be found in [12]. Figure 1.1 gives a short overview as well as showing the standard transitions.

The naming scheme *client/server* is due to historical reasons when tiered compilation was not available and users had to choose the JIT 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 the choice suggested.

Tiered compilation was introduced to improve start-up performance of the JVM. Starting with the interpreter results in instantaneous execution (i.e. a method is executed right away as there is no delay caused by the method's compilation). Also, there are always methods that are executed infrequently. For these, the compilation overhead can exceed the performance gain that results from having a compiled version of the method. C1 allows the JVM to have optimized code available early on. That code can be used to create a richer profile which is then used by the C2 compiler later on. Ideally this profile already contains most of the program flow and the assumptions made by C2 hold. If that is not the case, the JVM might need to go back, gather more profiles and compile the method again. In this case, being able to do quick compilations with C1 decreases the amount of C2 recompilations which are even more costly.

1.2 Deoptimizations

Ideally a method is compiled by making use of as much profiling information as possible. For example, since the profiling information is 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 (again see Figure 1.1). 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 even more 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 during execution.

The JVM then stops execution of that method and returns control back to the interpreter. This process is called *deoptimization* and considered very expensive. The previous interpreter state has to be restored and the method will be executed using the slow interpreter. Eventually the method might be recompiled with the newly gained information.

1.3 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 is 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 program stack. The JVM sets up a new stack frame for the compiled method which replaces the interpreters stack frame and execution will continue using the native method.

This process is called *on-stack replacement* and usually shortened to OSR. The Figure 1.2 presented in a talk by T. Rodriguez and K. Russel [13] gives a graphical representation. The benefits of OSR will become more obvious when looking at the first example in Chapter 2.

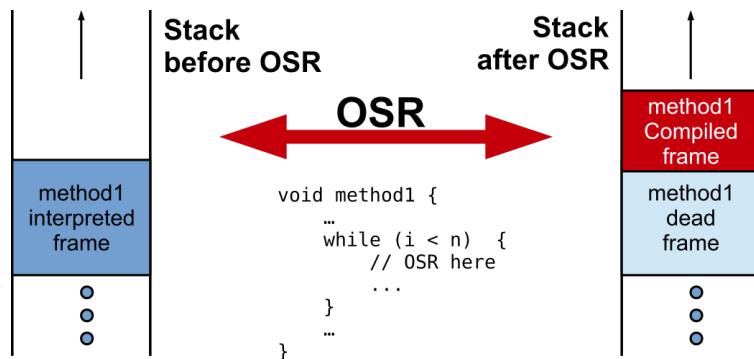


Figure 1.2: Graphical schema of OSR

1.4 Compile thresholds

The transitions between the compilation levels (see Figure 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 Level 3 and Level 3 to Level 4 happen when the following predicate returns true:

$$\begin{aligned}
 & i > \text{Tier}X\text{InvocationThreshold} * s \\
 \parallel & (i > \text{Tier}X\text{MinInvocationThreshold} * s \&\& i + b > \text{Tier}X\text{CompileThreshold} * s)
 \end{aligned}$$

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:

$$b > \text{TierXBackEdgeThreshold} * s$$

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

2 Motivation

We continue with presenting two simple example methods that illustrate the motivation for using cached profiles. The examples should provide an indication of how and why cached profiles can be beneficial for the performance of a Java Virtual Machine. We will omit any implementation details as they will be discussed in Chapter 3 in detail.

Ideally, being able to reuse the profiles from previous runs should result in two main advantages:

1. **Lower start-up time of the JVM:** By having information about program execution, the compiler can avoid gathering profiles and compile methods earlier and directly at higher compilation levels.
2. **Fewer Deoptimizations:** Since cached profiles are dumped at the end of a compilation, when using these profiles the compiler can already include all optimizations for all different method executions. The compiled code includes less uncommon traps and therefore fewer deoptimizations occur.

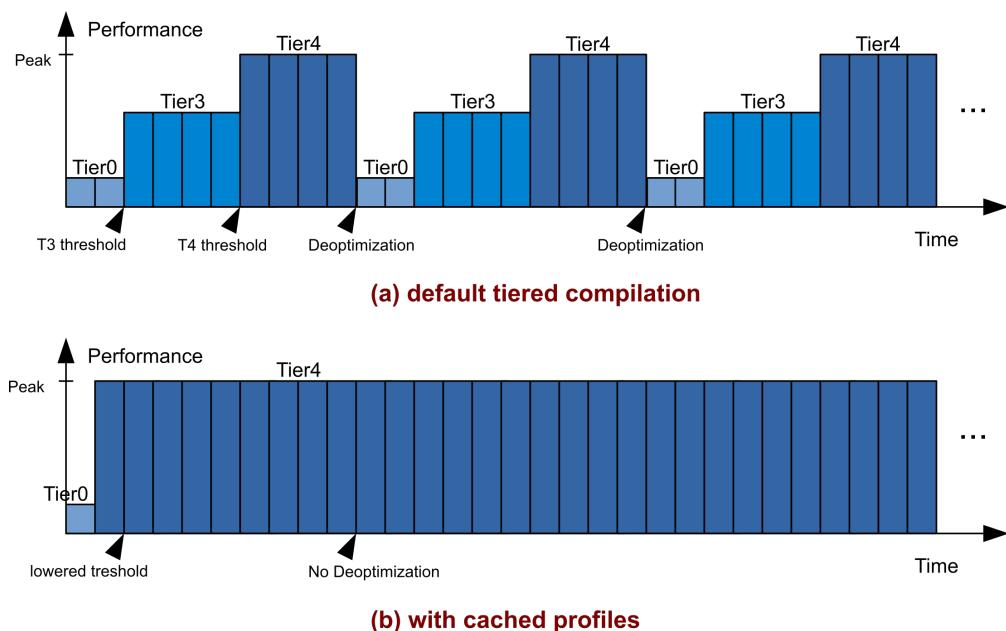


Figure 2.1: Schematic visualization of cached profile benefit

Figure 2.1 gives a schematic visualization of the expected effect on the performance of a single method when using cached profiles compared to the current state without such a system and standard tiered compilation. Each blue bar corresponds to an invocation of the method. Higher bars mean higher compilation levels and therefore higher performance. The x-axis represents time since the start of the JVM. The figure shows the ideal case and abstracts away many details and other possible cases. However, it provides a good visualization for the examples provided in this chapter. A more detailed performance analysis, also considering possible performance regressions is done in Chapter 4.

We are using my implementation described in Chapter 3 in CachedProfileMode 0 (see 3.4.1) built into OpenJDK 1.9.0. All measurements in this chapter are done on a Dual-Core machine running at 2 GHz with 8GB of RAM. To measure the method invocation time we use hprof [9] and the average of 10 runs. The evaluation process has been automated using a couple of python scripts. The error bars show the 95% confidence interval.

2.1 Example 1: Benefit of early compilation

For Example 1, on-stack replacement (OSR) has been disabled to keep the system simple and easy to understand. We use a simple class that invokes a method one hundred times. The method consists of a long running loop. The source code is shown in Listing 2.1. Since OSR is disabled and a compilation to level 3 is triggered after 200 invocations this method never leaves the interpreter. We call this run the *baseline*. To show the influence of cached profiles we 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 we use these profiles and achieve a significantly lower time spent executing the cached method as one can see in Figure 2.2. This increase comes mainly from the fact that having a cached profile available allows the JVM to compile highly optimized code for hot methods earlier (at a lower threshold) since there is no need to gather the profiling information first.

Since the example is rather simple neither the baseline nor the profile usage run trigger any deoptimizations. This makes sense because after the first invocation, all the code paths of the method have been taken already and are therefore known to the interpreter and saved in the profile.

When OSR is enabled, the performance difference between using cached profiles and not using them vanishes. That happens because HotSpot realizes the hotness of the method while it is being executed and the simplicity of the method allows the JIT compiler to produce optimized code already. The interpreted version is replaced on the stack by the compiled version during the first method invocation. This example, although rather artificial, shows the influence of early compilation with OSR disabled.

Listing 2.1: Simple method that does not get compiled

```

1 class NoCompile {
2     public static void main() {
3         double result = 0.0;
4         for(int c = 0; c < 100; c++) {
5             result = method1(result);
6         }
7     }
8     public static double method1(double count) {
9         for(int k = 0; k < 10000000; k++) {
10            count = count + 50000;
11        }
12        return count;
13    }
14 }
```

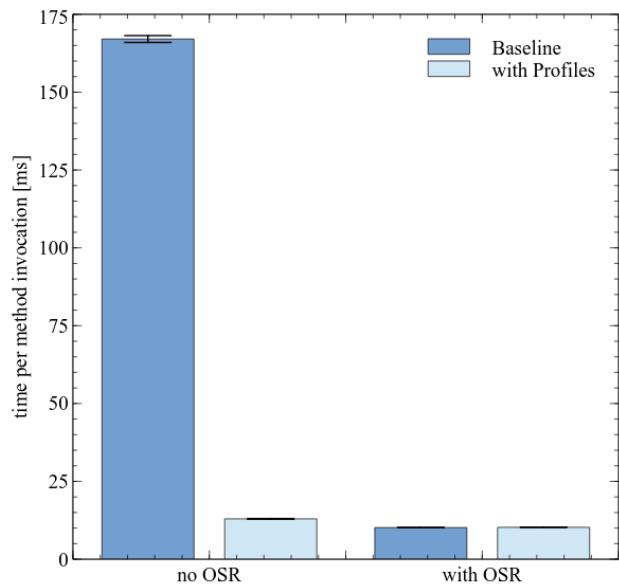


Figure 2.2: NoCompile.method1 - per method invocation time

2.2 Example 2: Benefit of fewer deoptimizations

OSR is one of the core features of HotSpot to improve startup performance of a JVM and disabling it does not give us any practical results. We came up with a second more complex example sketched in Listing 2.2, that demonstrates the influence of cached profiles without disabling any HotSpot functionalities.

The idea is to create a method that takes a different, long running branch on each of its method invocations. Each branch has been constructed in a way that it will trigger an OSR compilation. When compiling this method during its first iteration only the first branch will be included in the compiled code. The same will happen for each of the 100 method invocations. As one can see in Figure 2.3 the baseline indeed averages at around 134 deoptimizations and a time per method

Listing 2.2: Simple method that causes many deoptimizations

```

1 class ManyDeopts {
2     public static void main() {
3         double result = 0.0;
4         for(int c = 0; c < 100; c++) {
5             result = method1(result);
6         }
7     }
8     public static long method1(long count) {
9         for(int k = 0; k < 100000001; k++) {
10             if (count < 100000001) {
11                 count = count + 1;
12             } else if (count < 300000001) {
13                 count = count + 2;
14             .
15             .
16             } else if (count < 505000000001) {
17                 count = count + 100;
18             }
19             count = count + 50000;
20         }
21     }
22     return count;
23 }
24 }
```

invocation of 186 ms.

Now we use a regular execution to dump the profiles and then use these profiles. Theoretically the profiles dumped after a full execution should include knowledge of all branches and therefore the compiled method using these profiles should not run into any deoptimizations. As one can see in Figure 2.3 this is indeed the case. When using the cached profiles no more deoptimizations occur and because less time is spent profiling and compiling the methods the per method execution time is significantly faster with averaging at 169 ms now.

2.3 Similar systems

In commercially available JVMs the idea of caching profiles is not new. The JVM developed and sold by Azul Systems® called Zing® [2] already offers a similar functionality. Zing® includes a feature set they call ReadyNow!™ [1] which aims to increase startup performance of Java applications. Their system has been designed with financial markets in mind and to overcome the issue of slow performance in the beginning and performance drops during execution.

Azul Systems clients reported that their production code usually experiences a significant performance decrease as soon as the market goes live and the clients start trading. The reasons are deoptimizations, that occur for example due to uncommon branch paths being taken or yet unused methods being invoked. In the past, Azul Systems' clients used techniques to warm up the JVM,

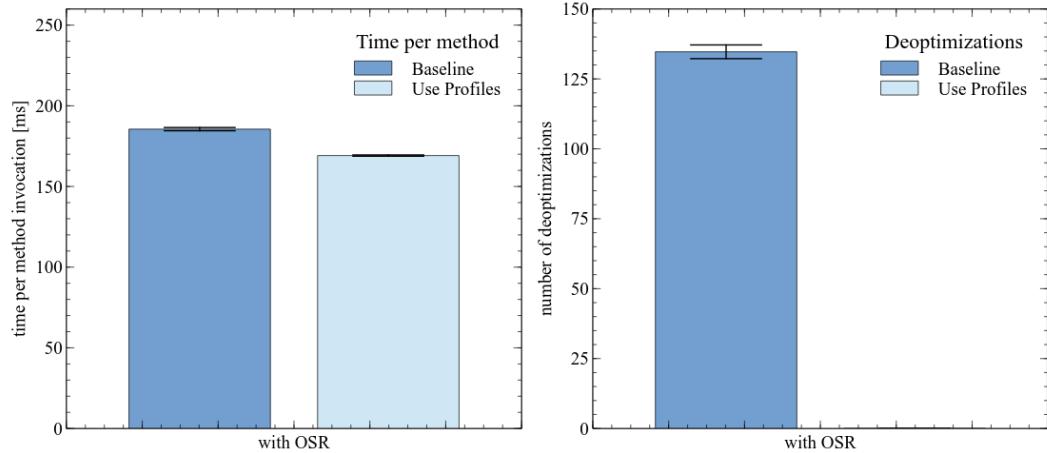


Figure 2.3: ManyDeopts.method1 - per method invocation time and deoptimization count

for example doing fake trades prior to market opening. However this does not solve the problem sufficiently well, since the JVM optimizes for these fake trades and still runs into deoptimizations once actual trades happen, because the code includes methods or specific code snippets that differ between the fake and the real trades.

ReadyNow!TM is a rich set of improvements how a JVM can overcome this issues. It includes attempts to reduce the number of deoptimizations in general and other not further specified optimizations. As one of the core features Azul Systems[®] implemented the ability to log optimization statistics and decisions and reuse this logs in future runs. This is similar to the approach presented in this thesis. However they do not record the actual optimization but the learning and the reasons why certain optimizations happen. This gives them the ability to give feedback to the user of the JVM whether or not certain optimizations have been applied. They also provide APIs for developers to interact with the system and allow further fine-grained custom-designed optimizations.

Unfortunately, Azul Systems does not provide any numbers about how their JVM actually improves performance when executing a software application or any analysis where the speedup originates from in detail.

3 Implementation / Design

This chapter describes the implementation of the cached profiles functionality for HotSpot, written as part of this thesis. HotSpot is a vital part of the open source Java Platform implementation, OpenJDK, and the source code is available at <http://openjdk.java.net/>.

Most of the code additions are included in two new classes `/share/vm/ci/ciCacheProfiles.cpp` and `/share/vm/ci/ciCacheProfilesBroker.cpp` as well as significant changes to `/share/vm/ci/ciEnv.cpp` and `/share/vm/compiler/compileBroker.cpp`.

The core functionality is located in `/share/vm/ci/ciCacheProfiles.cpp`, a class that takes care of setting up the cached profile data structure as well as providing public methods to check if a method is cached or not. The class `/share/vm/ci/ciCacheProfilesBroker.cpp` is used before a cached method is compiled. It is responsible for setting up the compilation environment, so the JIT compiler can use the cached profiles.

A full list of modified files and the changes can be seen in the webrev at <http://mohlerm.ch/b/webrev.01/> or Appendix A.3. The changes are provided in form of a patch for HotSpot version 1aef080fd28d. In the following, the original version is referred to as *baseline*.

We will describe and explain the functionality and the design decisions in the following sections, ordered by their execution order.

3.1 Creating cached profiles

The baseline version of HotSpot already offers a functionality to replay a compilation based on previously saved profiling information. This is mainly used in case the JVM crashes during a JIT compilation to replay the compilation process and allow the JVM developer to further investigate the cause of this incident. Apart from this automatic process, there exists the possibility to invoke the profile saving manually by specifying the `DumpReplay` compile command option per method.

We introduce a new method option called `DumpProfile` as well as a new compiler flag `-XX:+DumpProfiles` that appends profiling information to a file as soon as a method gets compiled. The first option can be specified as part of the `-XX:CompileCommand` or `-XX:CompileCommandFile` flag and allows the user to select single methods to dump their profile. The second command dumps profiles of all compiled methods. The profile is converted to a string and saved in a simple text file

called `cached_profiles.dat`.

The system will only consider compilations of Level 3 or Level 4. Level 1 and Level 2 are rarely used in practice and do only include none or little profiling information. The user can also restrict the profiles to Level 4 by using the compiler flag: `-XX:DumpProfilesMinTier=4`.

The dumped profiling information consists of multiple `ciMethod` entries, `ciMethodData` entries, and one `compile` entry. They are separated by line breaks and keywords to make sure the data can be parsed easily. A shortened example of a cached profile can be found in Appendix A.2. The `ciMethod` entries contain information about the methods used in the compilation (Table 3.1 describes it in more detail). The `ciMethodData` (see Table 3.2) includes all profiling data about the methods itself to be able to redo the compilation. The `compile` entry saves the bytecode index in case of OSR, the level of the compilation and lists all inlining decisions (Table 3.3).

A method can be compiled multiple times and at different tiers, thus compilation information for the same method can be dumped multiple times. This is intentional and is taken care of when loading the profiles (see Section 3.2).

Table 3.1: content of `ciMethod` entry in cached profile

name	description
class_name, method_name, signature	used to identify the method
invocation_counter	number of invocations
backedge_counter	number of counted backedges
interpreter_invocation_count	number of invocations during interpreter phase
interpreter_throwout_count	how many times method was exited via exception while interpreting
instructions_size_name	rough size of method before inlining

Table 3.2: content of `ciMethodData` entry in cached profile

name	description
class_name, method_name, signature	used to identify the method
state	if data is attached and matured
current_mileage	maturity of the oop when snapshot is taken
orig	snapshot of the original header
data	the actual profiling data
oops	ordinary object pointers, JVM managed pointers to object

Table 3.3: content of compile entry in cached profile

name	description
class_name, method_name, signature	used to identify the method
entry_bci	byte code index of method
comp_level	compilation level of record
inline	array of inlining information

3.2 Initializing cached profiles

The information dumped in step 3.1 can now be used in a next run of that particular program. To specify that profiles are available, we introduce a new compiler flag `-XX:+CacheProfiles` that enables the use of previously generated profiles. By default, it reads from a file called `cached_profiles.dat` but a different file can be specified using `-XX:CacheProfilesFile=other_file.dat`.

Before any cached profiles can be used the virtual machine has to parse that file and organize the profiles and compile information in a data structure. This data structure is completely kept in memory during the whole execution of the JVM to avoid multiple disk accesses. The parsing process is invoked during boot up of the JVM, directly after the `compileBroker` gets initialized. This happens before any methods are executed and blocks the main thread of the JVM until finished.

As mentioned in Section 3.1, the file consists of method information, method profiles, and additional compile information. The parser scans the file once and creates a so called `CompileRecord` for each of the methods that includes compilation information in the file. This compile record also includes the list of method information (`ciMethod`) and their profiling information (`ciMethodData`). As mentioned previously, a method's compile information could have been dumped multiple times, which results in multiple `CompileRecords` for the same method. In this case, HotSpot will only keep the `CompileRecords` based on the latest data written to the file but never overwrite an existing higher level profile. Because a profile dumped by the C1 compiler can not be used by the C2 compiler and the other way around, the level of the profile matters as it influences the compile level transitions described in Section 3.4. And since profiling information only grows, the compilation that happened last contains the richest profile and is considered the best. This is based on the fact that the richer the profile, the more information about the method execution is known and influences the compiled version of that method. For example, a profile for a method might include data for all its branches and can therefore help avoid running into uncommon traps and trigger deoptimizations.

The `CompileRecord` as well as the lists of method information and profiles are implemented as an array located in HotSpot's heap space. They get initialized with a length of 8 and grow when

needed. This choice has been done for simplicity and leaves room for further improvements.

3.3 Using cached profiles

The implementation offers three different modes: Mode 0, Mode 1, and Mode 2, that differ in the way they use the cached profiles. The following paragraph applies to all three modes and I will discuss the differences of the modes in detail in Section 3.4.

The idea is to modify the compiler to use cached profiles if available and continue as usual otherwise. A simplified graphical overview of the program flow for compiling a method with the changes introduced in this thesis can be found in Figure 3.1.

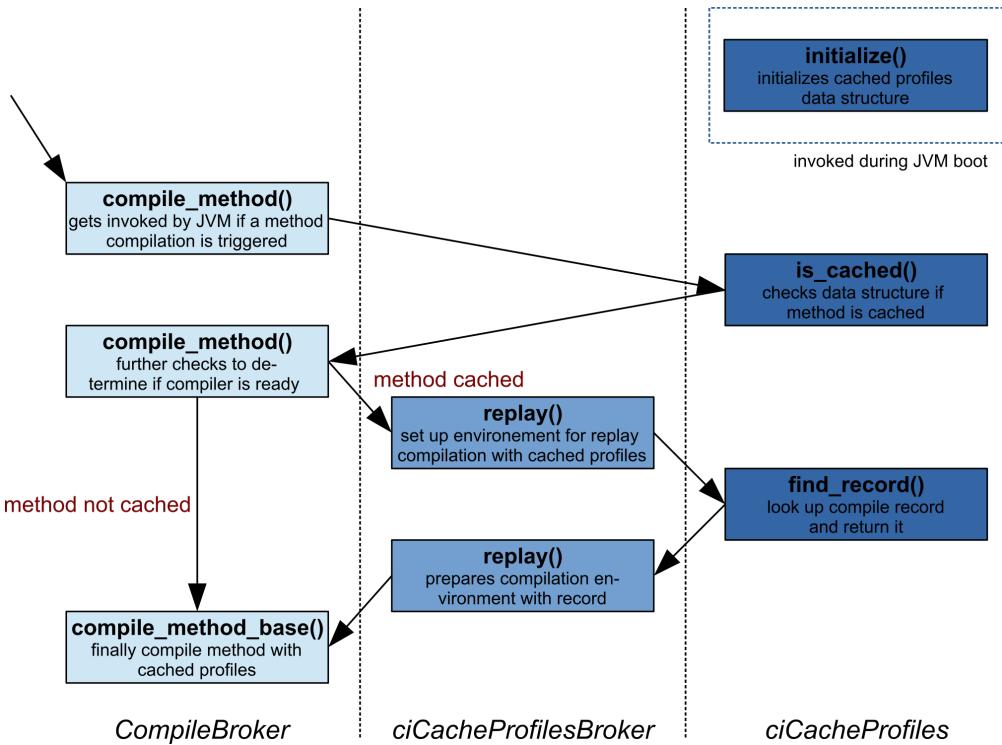


Figure 3.1: program flow for compiling a method

As mentioned before, once compilation thresholds are exceeded a method is scheduled for compilation. This means that the JVM will invoke a method called `compile_method()`, located in the `compileBroker` class. This method tests if certain conditions hold, for example, it checks if the compile queue is not full or if there is already another compilation of that particular method running. We extended this method with a call to `ciCacheProfiles::is_cached(Method* method)` which does a linear scan through the `CompileRecord` array data structure. The method returns either 0 if the method is not cached or returns an integer value, reflecting the compile level, in case a cached profile of this method is available. Because only methods compiled with level 3 or 4 are cached, this call to `is_cached()` only gets executed if the compilation request is also of level 3 or higher.

Depending on the compilation level of the profile, the level of the requested compilation, and the `CacheProfileMode`, the `compileBroker` then schedules either a compilation using freshly gathered profiles or calls into `ciCacheProfilesBroker` to replay the compilation, based on a cached profile. In contrast to cached profiles, fresh profiles describe profiles gathered during the current run of the JVM. Since these decisions are different in each mode, I describe them in detail in the next section. In case the method is not cached, the execution continues like in the baseline version. Otherwise, the `ciCacheProfilesBroker` class then initializes the replay environment and retrieves the compile record from `ciCacheProfiles`. Subsequently, the needed cached profiles get loaded to make sure they are used by the following compilation. `ciCacheProfilesBroker` then returns the execution to the `compileBroker`, which continues with the steps needed to compile the method. Again some constraints are checked (e.g. if there is another compilation of the same method finished in the meantime) and a new compile job is added to the compile queue. Eventually the method is compiled using the cached profiles.

Since the implementation is only invoked by the static class `compileBroker`, `ciCacheProfiles` and `ciCacheProfilesBroker` are static classes as well. The `compileBroker` is solely called by the JVM main thread, therefore there is no need to make the `compileRecord` data structure or any of the new methods thread safe.

3.4 Different usage modes for cached profiles

The implementation of cached profiles offers 3 different modes, which differ from each other in the transitions between the compilation tiers. The motivation as well as the advantages and disadvantages of the modes are described in the following three subsections. While Mode 0 and Mode 1 are similar except for the compile thresholds, mode2 differs significantly. Figure 3.2 provides a graphical overview of the differences in the compilation tier transitions of the modes.

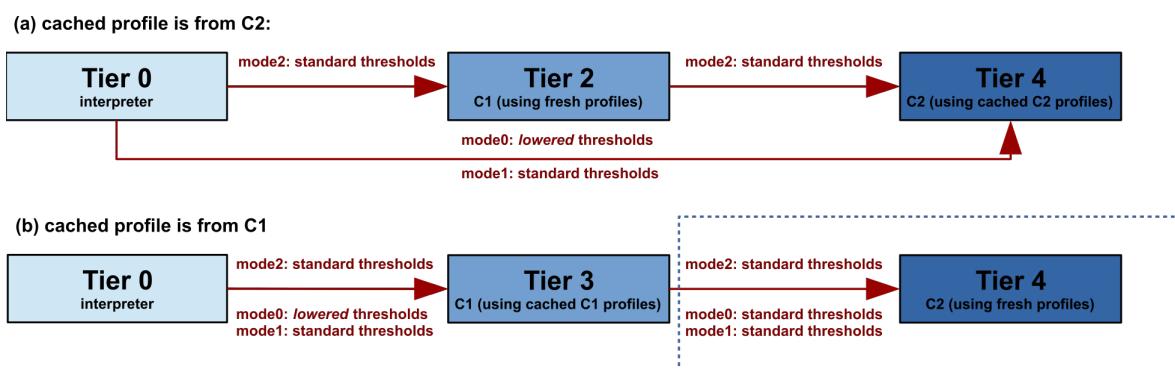


Figure 3.2: Tier transitions of different modes

All three modes have advantages and disadvantages and their performance is evaluated in Chapter 4. Since depending on the methods, different modes might perform best. It is up to the JVM user

to decide which mode he or she wants to use. Because **Mode 2** is the most conservative one, which does not modify thresholds it is considered the default mode and used if not further specified.

3.4.1 Compile thresholds lowered (Mode 0)

The first mode is based on the consideration that a method that has a profile available does not require profiling anymore. Therefore, the compile thresholds (see Section 1.4) of these methods are lowered. If the cached method has a C2 profile, all thresholds are lowered, in case of a C1 profile, only the thresholds that affect the tiers below and equal to Tier 3. This differentiation is done to prevent an increase in early C2 compilations using fresh profiles. Because a lower Tier 4 threshold would mean, that a method which only got compiled at Tier 3 when creating the profile, might now trigger a Tier 4 compilation very early. At this point, the fresh profiles generated by the Tier 3 compiled version are still immature. Using them to compile with Tier 4 results in code that includes only limited profiling information, which causes costly deoptimizations when used.

By default, the thresholds are lowered to 1% of their original values but the threshold scaling can be modified with the JVM parameter: `-XX:CacheProfilesMode0ThresholdScaling=x.xx`. 1% results in the Level 3 invocation counter being reduced from 200 to 2. This means that the method will be interpreted once but then trigger a compilation on the next invocation. Since the interpreter also handles class loading, this decision has been made to avoid the need of loading classes in C1 or C2 which was considered out of the scope for this thesis. HotSpot expects classes of standard libraries to be loaded in a very specific order. Moving class loading to C1 or C2 would mess with that order and can therefore not be done without huge changes to the JVM.

In **Mode 0**, the JIT compiler will always use a cached profile for compilations of Level 3 or Level 4 in case there is a cached profile which has been generated by a compiler of the same level. However, if a method to be compiled on Level 3 has a cached profile available for Level 4, the compiler will skip the C1 compilation and immediately compile with C2. In this case, HotSpot directly uses the highly optimized version generated by C2 and in an ideal case the method will need less time to reach peak performance.

However, since the thresholds of all methods with cached profiles get lowered and some of the C1 compilations are promoted to C2 compilations, the C2 compiler is put under heavy load. Especially during startup of a program, where many compilations happen naturally, C2 might not be able to handle all these requests at the same time and the compile queue fills up. This might negatively affect performance and is analyzed in Section 4.4.

3.4.2 Unmodified compile thresholds (Mode 1)

Mode 1 is doing exactly the same as **Mode 0** but does not lower the compilation thresholds of methods with cached profiles. This is done to reduce the load on C2 as mentioned in Subsection 3.4.1. Apart from this change **Mode 1** has the same behaviour as **Mode 0**.

3.4.3 Modified C1 stage (Mode 2)

Both modes mentioned above use cached profiles as soon as a compilation of Level 3 or 4 are triggered. Since the thresholds for Level 3 are smaller than the Level 4 thresholds (see Appendix A.1) a method reaching a Level 3 threshold could actually trigger a Level 4 compilation, if the cached profile is one of Level 4. So even if Mode 1 is used and the thresholds are untouched, C2 might get overloaded since compilations occur earlier.

Mode 2 has been designed to make as little changes as possible to the tiered compilation and prevent C2 being used more than usual. It does so by keeping the original tiered compilation steps and compilation thresholds and compiles methods with C1 prior to C2. But since there are already profiles available, there is no need to run at Tier 3 to generate full profiles but instead it uses Tier 2. Tier 2 does the same optimizations but offers only limited profiles like method invocation and backbranch counters. They are needed to know when to trigger the C2 compilation and therefore Tier 1 can not be used. By avoiding Tier 3 and using Tier 2 instead methods spend less time in code gathering profiling information and therefore method execution is considered about 30% faster [8]. Eventually, if the Tier 4 thresholds are reached, the method is compiled using C2 and the cached profiles. This still maintains the benefit from having more complete profiles available early but avoids modifying thresholds which could result in a very different load to the compiler.

The above only makes sense if the cached profile is a C2 profile. If only a C1 profile is available, C1 should gather fresh, full profiles since they might be needed in C2 later. HotSpot will then only use the cached profile during the C1 compilation and then use the generated profiles for possible C2 compilations. In theory this transition is considered rare, because if a method has not been compiled with C2 when creating the profile it is unlikely to get compiled with C2 in the future.

To summarize, we expect a performance benefit C1 compiled code by using Level 2 instead of Level 3 and a benefit for C2 compiled code by using more comprehensive profiles.

3.5 Problems

If the profiles generated by multiple runs of the program deviate sharply it is likely that a cached profile does not fit to the current execution. In this case the compiled version would still trigger many deoptimizations and the method could end up having even worse performance since it's going to use the profile over and over again. For each method, the JVM maintains a deoptimization counter. Cache profiles are used if the counter is below a certain limit. If they are above that limit a standard compilation using freshly gathered profiles will be used instead. The limit is 10 to allow a small number of recompilations. This could for example be useful when the method is deoptimized due to classes not being loaded. The value of 10 seems reasonable for all executed measurements.

3.6 Debug output

For debugging and benchmarking purposes four debug flags are implemented, that can be used along with `-XX:+CacheProfiles`.

flag	description
<code>-XX:+PrintCacheProfiles</code>	enable command line debug output for cached profiles
<code>-XX:+PrintDeoptimizationCount</code>	prints amount of deoptimizations when the JVM gets shut down
<code>-XX:+PrintDeoptimizationCountVerbose</code>	prints total the amount of deoptimizations on each deoptimization
<code>-XX:+PrintCompileQueueSize</code>	prints the total amount of methods in the compile queue each time a method gets added

4 Performance

This section evaluates the performance of the cached profile implementation using two modern benchmark suites. The goal is to provide indicators of the performance influence and analyze where this performance influence comes from.

4.1 Setup

To provide reliable and comparable results, all benchmarks were executed on single nodes of the Data Center Observatory provided by ETH [3]. A node features 2 8-Core AMD Opteron 6212 CPUs running at 2600 MHz with 128 GB of DDR3 RAM and a solid state disk for storage. The nodes are running Fedora 19 and GCC 4.8.3. All JDK builds were compiled on the nodes itself. The benchmarks were each executed on a dedicated node to prevent any influence of inter node performance differences.

To compare performance the following benchmark suites were used:

1. **SPECjvm 2008:** A benchmark suite developed by Standard Performance Evaluation Corporation for measuring the performance of the Java Runtime Environment [14]. We use version 2008 and we run a subset of 17 out of a total of 21 benchmarks. 4 are omitted due to incompatibility with openJDK 1.9.0.

SPECjvm reports the number of operations per minute (ops/m). This is used to compare the performance and more ops/m equals better performance.

2. **Octane 2.0:** A benchmark developed by Google to measure the performance of JavaScript code found in large, real-world applications [4]. Octane runs on Nashorn, a JavaScript Engine of Hotspot. The version used is 2.0 and consists of 17 individual benchmarks of which 16 are used. One benchmark is not compatible with openJDK 1.9.0.

Octane gives each benchmark execution a score reflecting the performance. the higher the score, the better the performance.

A more detailed description of the benchmark suites used in this performance evaluation can be found in Appendix A.4.

The benchmark process was automated using a number of self-written python scripts. Unless specified otherwise, the graphs in this chapter always show the arithmetic mean of 50 runs and the error bars display the 95% confidence intervals.

4.2 Benchmark performance

The main goal of cached profiles is to improve the warmup performance of the JVM. Having a rich profile from an earlier execution will allow the JIT compiler to use a highly optimized version early in method execution.

The different CacheProfilesModes were implemented to being able to compare the performance influence of certain design decisions. We expect the modes to produce different results and the following list suggests reasons for these performance differences:

- If the cached profiles fit nicely to the current method execution compiling these methods earlier than in the baseline will save executions at lower tiers and decrease the time needed for warmup. Benchmarks with these properties should achieve a performance improvement and favor **Mode 0**.
- In case many methods are compiled early this could overload the compile queue and delay compilation of these methods. Also, if the cached profile does not fit the current method execution (the limit of 10 deoptimization was reached) the JVM will eventually use freshly generated profiles. These profiles could be very sparse when using lowered compilation thresholds. In these scenarios we expect **Mode 1** to outperfom **Mode 0**.
- **Mode 2** keeps the steps of the original tiered compilation and is considered the most conservative mode. **Mode 2** also does not modify the compilation thresholds and therefore puts the same load on the compile queue than the baseline version. When using this mode we can exclude the effect of earlier compilation and the performance only originates from the faster C1 phase and the better code quality in C2.

4.2.1 SPECjvm warmup performance

The longer a program is running, the less impact a faster warmup has. Considering most benchmarks include a warmup phase which does not count towards the final score simply running the complete benchmark suite is not an option. Instead, we limited SPECjvm to 1 single operation which, depending on the benchmark take around 6 to 40 seconds. Additionally, the JVM gets restarted between each single benchmark to prevent methods shared between benchmarks being compiled already.

We run each benchmark with all cached profiling features disabled. This run is called the *baseline* and displays the current openJDK 1.9.0 performance.

We then use a single benchmark run where we configure the JVM to dump the profiles to disk. This run is not limited to a single operation and instead uses the default values of the benchmark. By default the benchmark is limited by time and runs for about 6 minutes. The idea is that these profiles include information that are usually not available during warmup and result in less deoptimizations and better code quality.

These profiles are then used in 3 individual runs using the introduced `-XX:CacheProfiles` flag. Each run is using one of the 3 different CacheProfileModes.

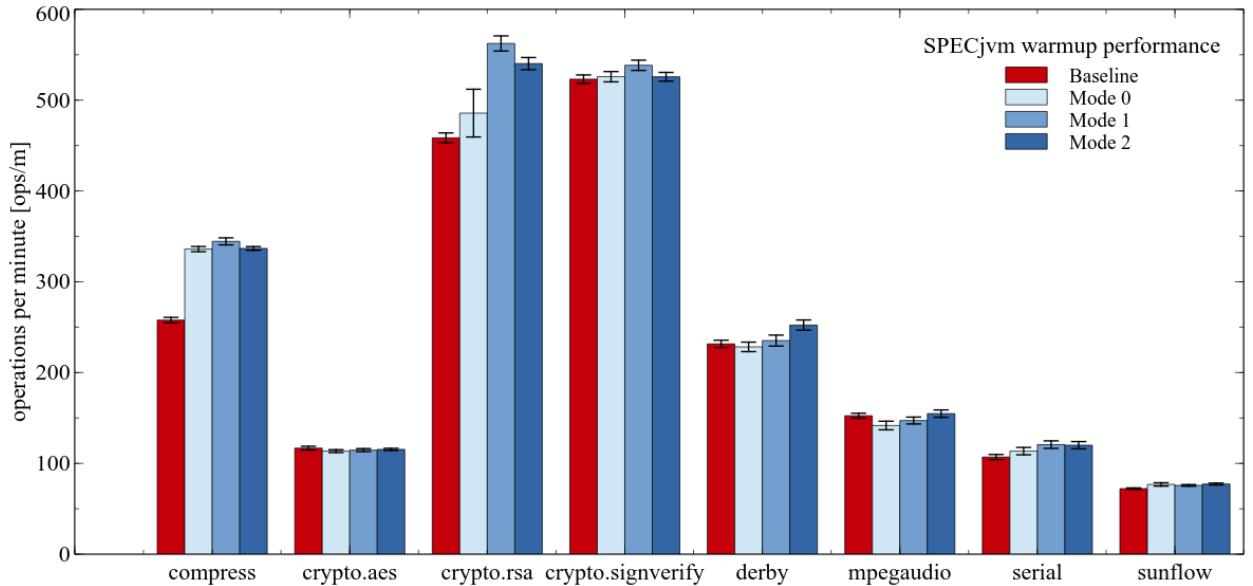


Figure 4.1: SPECjvm benchmarks on all different modes

Figures 4.1 and 4.2 shows the number of operations per minute, measured for each benchmark individually. Note, that the operations per minute is not to be confused with the *1 operation* of the benchmark itself. Figure 4.3 summarizes the results by showing the relative performance compared to the baseline.

The individual benchmarks show different effects on performance. Taking the average of all modes, we see a performance increase up to around 34% in the compress benchmark (Mode 1) and a performance decrease of down to 20% in scimark.sparse.large (Mode 0).

Interestingly, the performance differences between the modes is not the same when comparing the individual benchmarks. For example in crypto.rsa Mode 0 clearly performs worst but in scimark.sparse.small it performs best. JVM performance is known to be very hard to predict and it seems not to be different when cached profiles are used. On average the performance of the benchmark warmup is improved by 2.64%, 3.37%, and 2.67% for Mode 0, Mode 1 and Mode 2.

Between the three different modes there is no clear *winner*. Each mode wins and loses in certain benchmarks against the others in terms of performance. However, in 12 out of 17 benchmarks at least one of the CacheProfileModes improves performance.

We will take a more detailed look at single benchmarks later in this chapter.

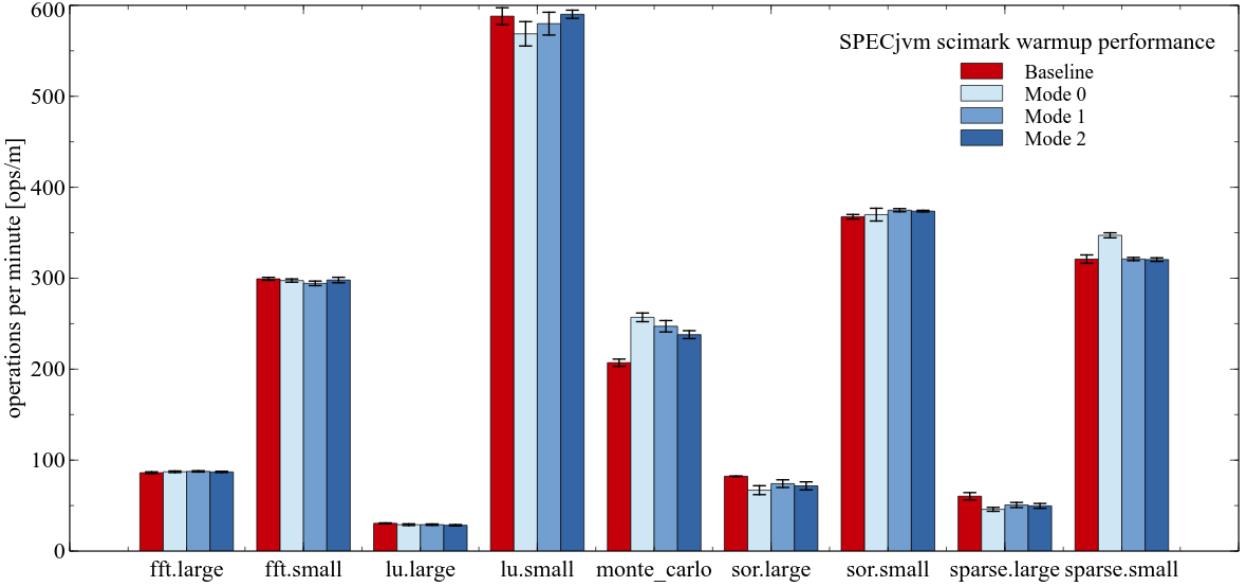


Figure 4.2: SPECjvm scimark benchmarks on all different modes

4.2.2 Octane performance

Since the individual Octane benchmarks are rather short (most of them run for between 4 and 30 seconds) and there is no way to run a fixed number of iterations (without modifying the Octane source) we run the Octane benchmarks completely. We still split up the execution in the individual benchmarks to achieve many JVM restarts. The rest of the setup is identical to the SPECjvm run in Section 4.2.1.

The absolute results are shown in Figure 4.4 and a relative comparison with the baseline in Figure 4.5. Compared to SPECjvm the Octane performance is more scattered. The richards benchmark increases by around 50% in Mode 0 while navierstockes decreases by around 25% in Mode 1. In most benchmarks (9 out of 14) Mode 0 performs worst. We assume this is related to the increased load of the compile queue and will therefore take a more detailed look at this in Section 4.4. The performance of the two other modes is better in most benchmarks, but in total only 6 out of 14 benchmarks result in a performance improvement in at least one mode.

4.3 Deoptimizations

We are still eager to figure out where the performance increase and decrease come from. We aim to lower the time needed for warmup by compiling methods earlier and or at lower tiers but also expect to decrease the number of deoptimizations by having more complete profiles early, which ideally results in better compiled code quality. To measure the total amount of deoptimizations we added a new compiler flag `-XX:+PrintDeoptimizationCount`. The total number of deoptimizations of the SPECjvm benchmarks is shown in Figure 4.6 and Figure 4.7. The Octane numbers are drawn

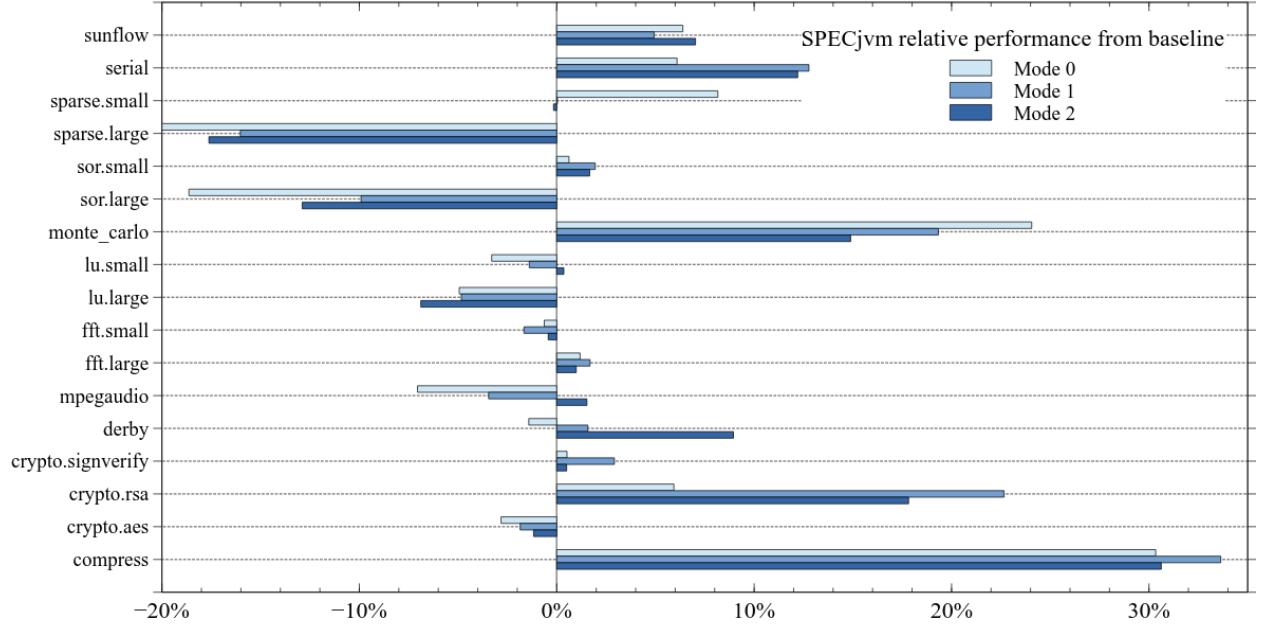


Figure 4.3: Relative performance from baseline for all SPECjvm benchmarks

in Figure 4.9. Again, we also included graphs that show the number of deoptimizations relative to the baseline runs in Figure 4.8 and Figure 4.10.

The measurements show, that when using Mode 1 or Mode 2, we are able to reduce the deoptimizations significantly in all benchmarks except one (gameboy). In Mode 0 there is a clear difference between SPECjvm and Octane. While in SPECjvm the number of deoptimizations is similar to the other modes, in Octane Mode 0 on average increases the number by 30%. Mode 0 also has the worst performance for Octane and we assume the amount of deoptimizations to be one of the reasons for that regression.

And while a low deoptimization number is a good indication of the increased code quality for methods being compiled with cached profiles we could not find a direct correlation between number of deoptimizations and the performance results.

One possible reason is that the amount of deoptimizations does not necessarily describe the performance impact. Especially when considering multi-threaded systems there can be a huge number of deoptimizations in performance uncritical threads that are avoided by using cached profiles and therefore heavily reduce the total counter. But if there is one very important method in a performance critical thread, that has executions that are not reflected in the cached profiles, this method could trigger only very few deoptimizations but still influence performance significantly.

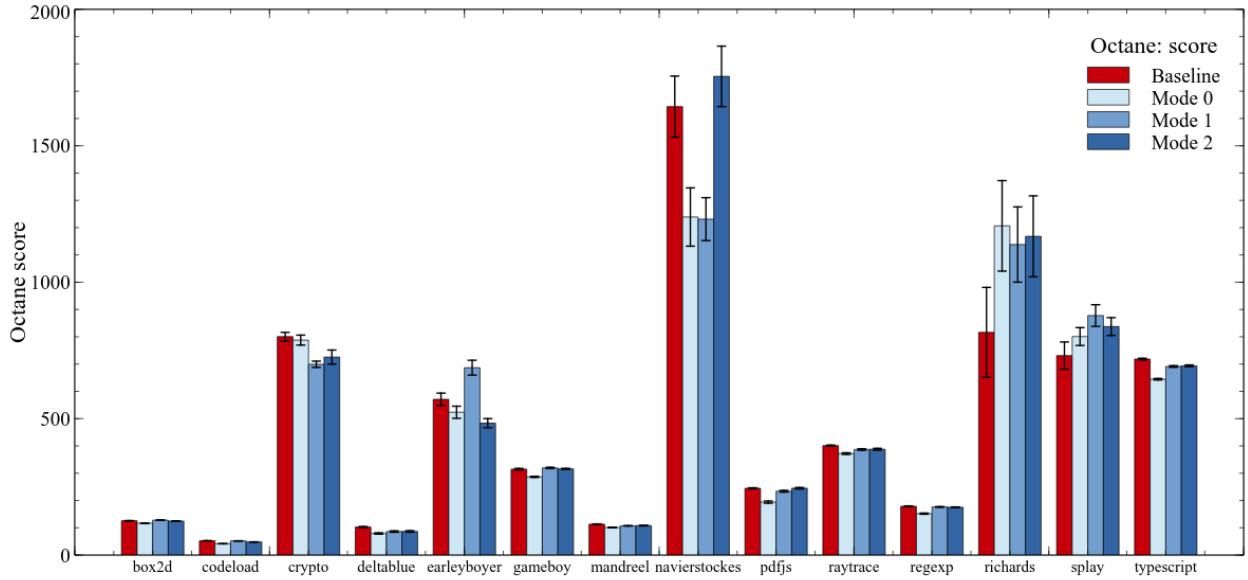


Figure 4.4: Octane benchmarks on all different modes

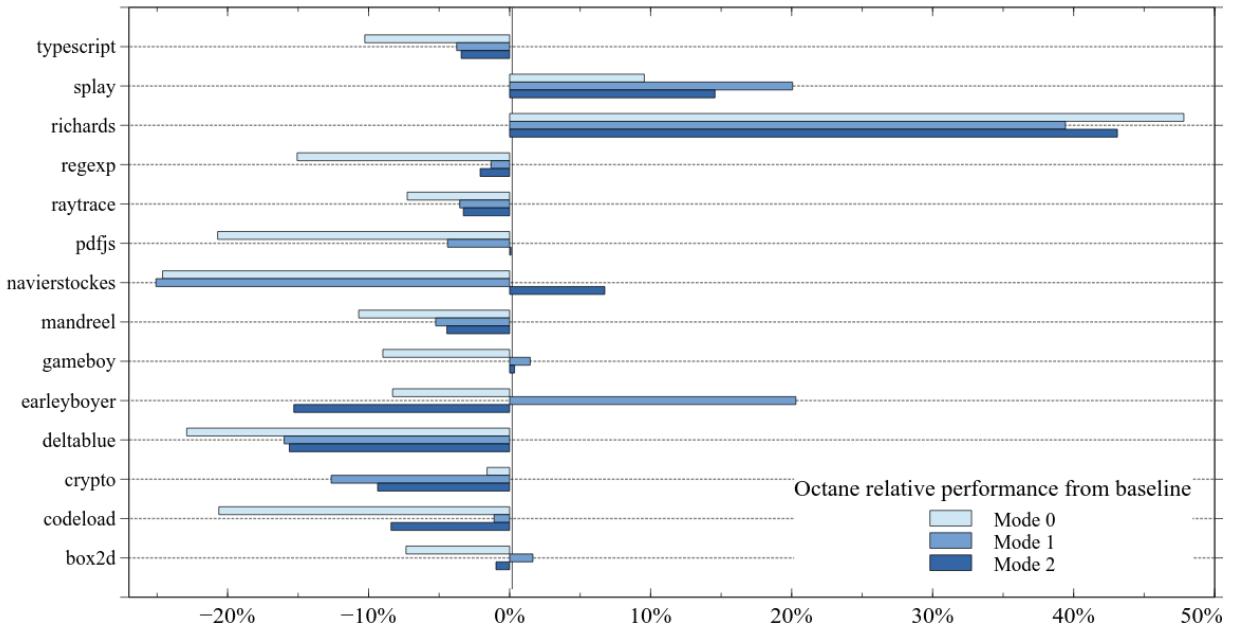


Figure 4.5: Relative performance from baseline for all Octane benchmarks

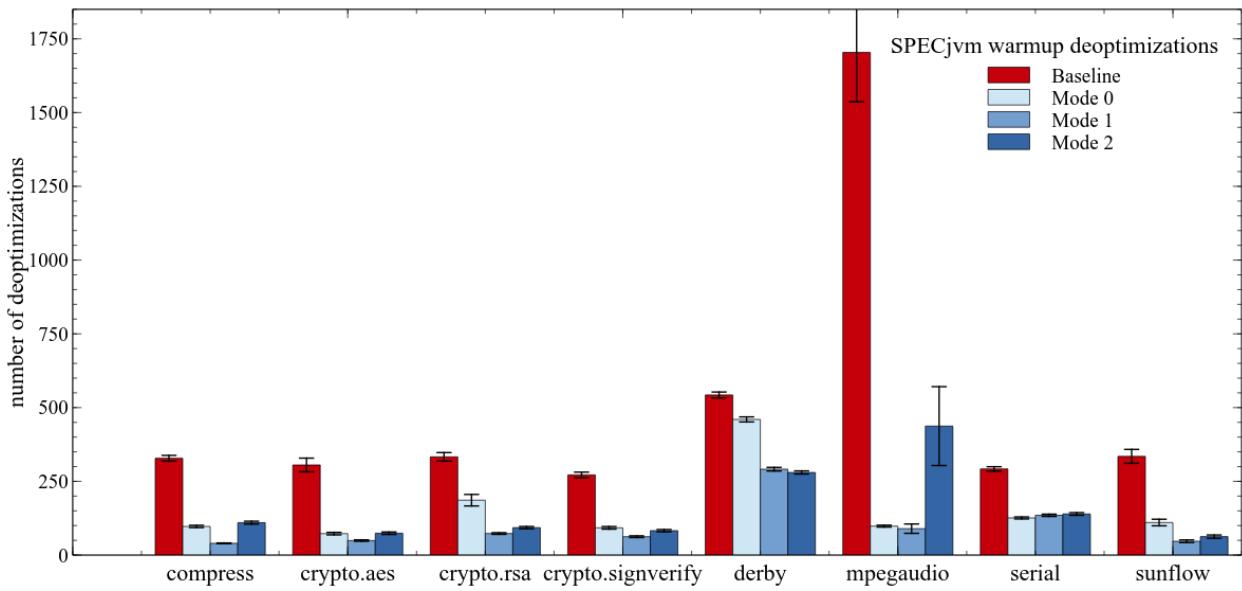


Figure 4.6: SPECjvm deoptimizations of all modes

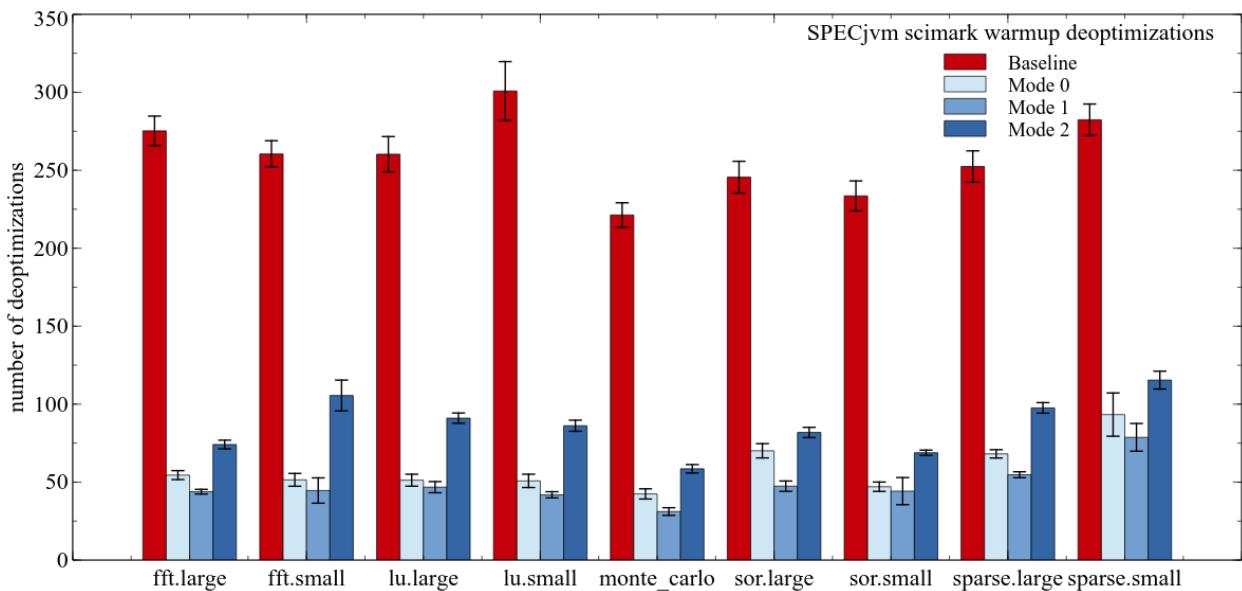


Figure 4.7: SPECjvm scimark deoptimizations of all modes

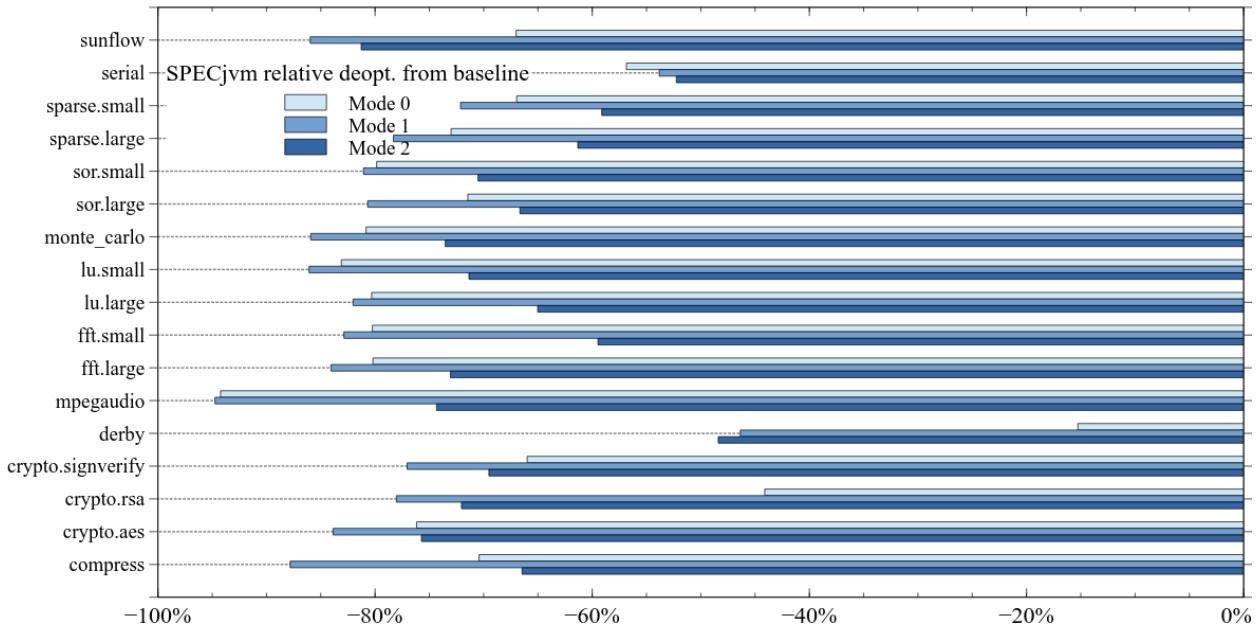


Figure 4.8: Relative deoptimizations from baseline for all SPECjvm benchmarks

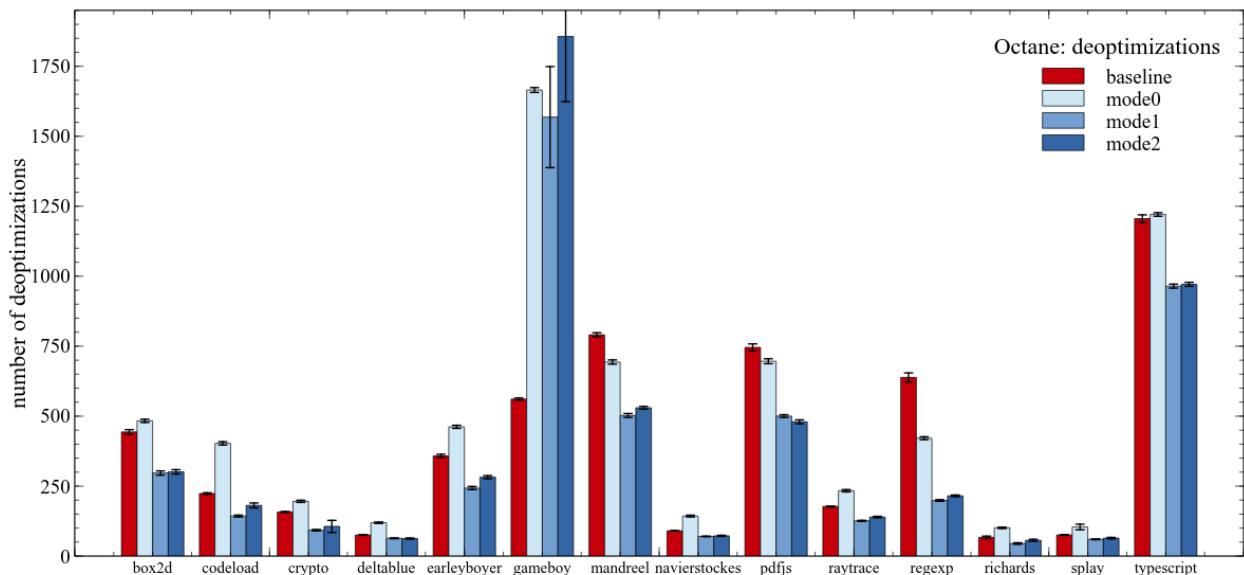


Figure 4.9: Octane deoptimizations of all modes

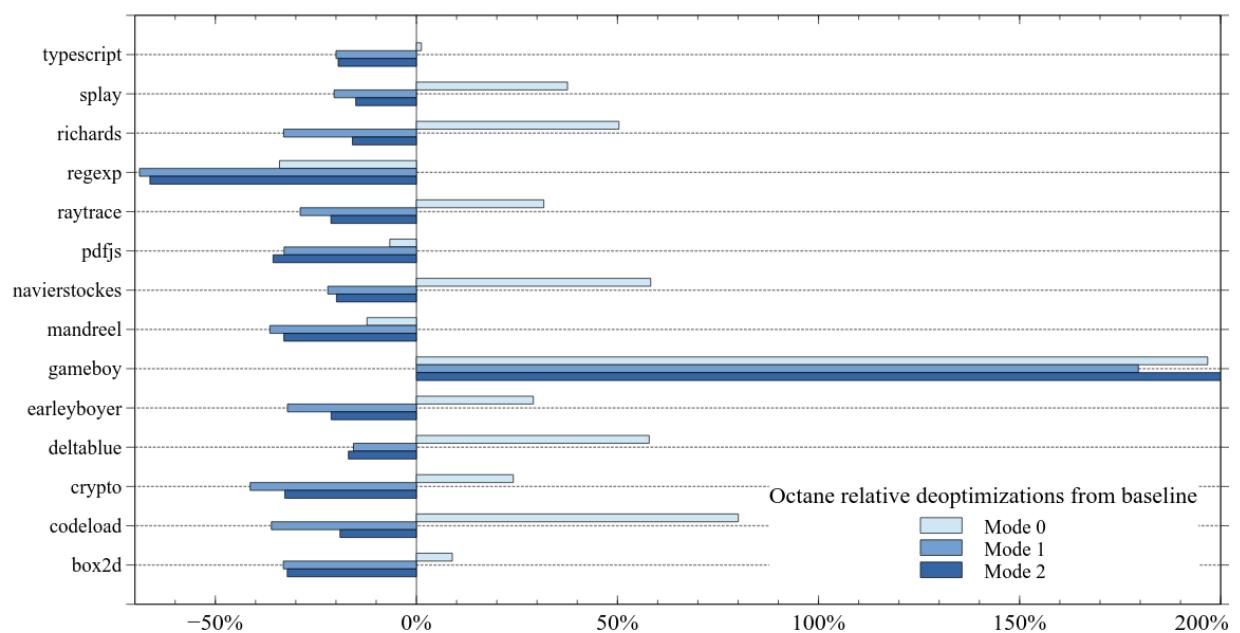


Figure 4.10: Relative deoptimizations from baseline for all Octane benchmarks

4.4 Effect on compile queue

When designing the different CacheProfilesModes, we thought that lowering the compilation thresholds will increase the load on the compiler, especially early in program execution. This results in the compiler not being able to handle all requests immediately and the compile queue fills up. This could delay execution of compiled methods and therefore decrease performance.

We added a new HotSpot flag `-XX:+PrintCompileQueueSize` which allows us to trace the current number of methods that are scheduled for compilation. We selected 6 individual benchmarks and printed their C1 and C2 compile queue for all three CacheProfileModes in Figure 4.11 to Figure 4.18. Keep in mind that the reason for the runs using cached profiles starting their compilations on later time stamps is because they need to parse the cached profile file first. The selected benchmarks and the reason why an individual analysis was performed are listed below:

- **Octane Richards:** This benchmarks achieved the highest performance benefit from using cached profiles in all three modes. We are interested to see if the compile queue load differs from worse performing benchmarks.
- **Octane EarleyBoyer:** This is a benchmark where Mode 1 performs significantly better than the other two modes. Does the compile queue give any indication why this is the case?
- **Octane NavierStokes:** Navierstokes has a 8% performance increase in Mode 2 but a 25% performance decrease in Mode 0 and Mode 1. Motivation is the same as for EarleyBoyer.
- **Octane Deltablue:** This benchmark achieves the highest performance loss from using cached profiles in all three modes. Together with the best benchmark (Richards) does the load on the compile queue give any indication of performance?
- **SPECjvm compress:** Compress is the best performing SPECjvm benchmark when using cached profiles.
- **SPECjvm scimark.sparse.large:** This is the worst performing SPECjvm benchmark when using cached profiles.

We realize that analyzing the compile queue does not really help us understanding the performance variations when using cached profiles. The graphs showing the C1 compile queue size over time do not significantly differ, neither from the baseline, nor is there a difference between the individual modes. The C2 compile queue graphs show different courses but we can not connect this with the benchmark performance.

Figure 4.12 show the C2 compile queue of the Octane Richards benchmark. As expected, due to removing steps from the tiered compilation, we increased the load on C2 in Mode 0 and Mode 1 with compile queue peaks at around 20 scheduled compilations. Nevertheless, these modes achieve close a performance of close to 50% better than the baseline. Mode 2 which was designed to keep the original tiered compilation steps unmodified does not have these peaks but nevertheless achieves

similar performance.

EarleyBoyer's compile queue is displayed in Figure 4.12. Mode 1 performs better than the other two modes and compared to Mode 0 seems to put even more pressure on the compile queue. It is interesting that in this particular benchmark even the baseline version puts a lot of pressure on the compile queue early on.

In Figure 4.16 we see NavierStokes' compile queue. Mode 2 performs best but we could not derive any indications why this is the case from looking at the queue size.

The Deltablue benchmark shown in Figure 4.18 has the worst performance when using cached profiles but the compile queue size looks very similar to the one of the Richards benchmark, where performance is significantly better.

We will omit looking at the SPECjvm benchmarks since they do not offer any new insights. The graphs can be found in the Appendix A.6.

The detailed analysis of the compile queue shows that our thoughts about the effect on the compile queue were not unfounded for most of the selected benchmarks. However, we were not able to relate these influences to actual performance effects. Especially, overloading the compile queue does not necessarily negatively affect performance.

This indicates that the performance differences are even more related to the actual code quality.

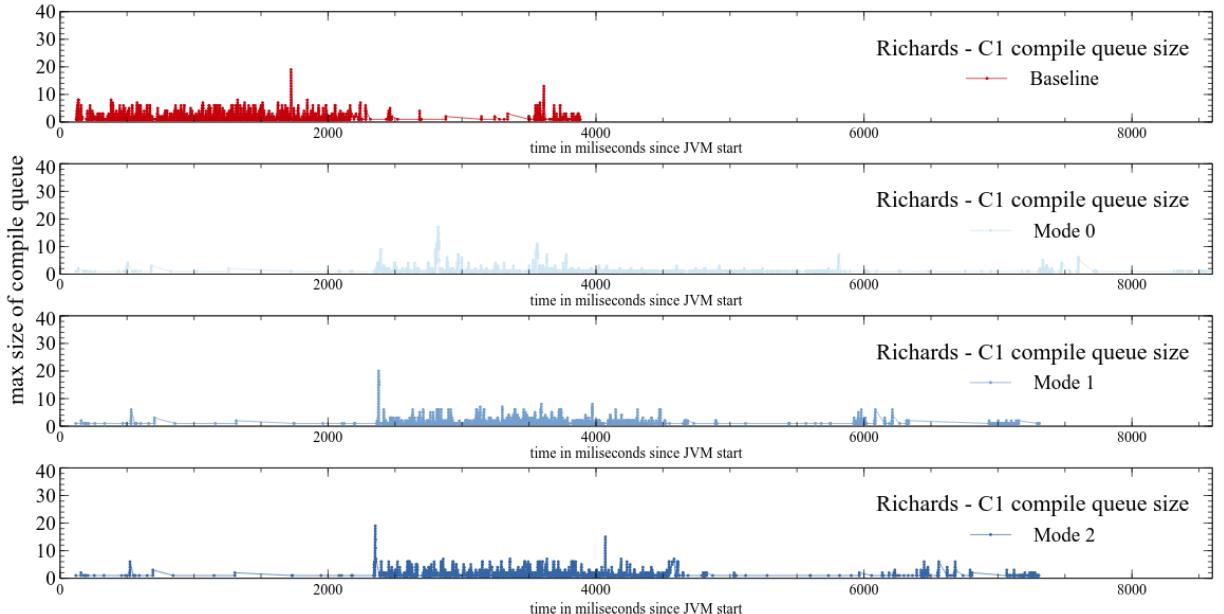


Figure 4.11: C1 Compile queue size over time Octane Richards benchmark

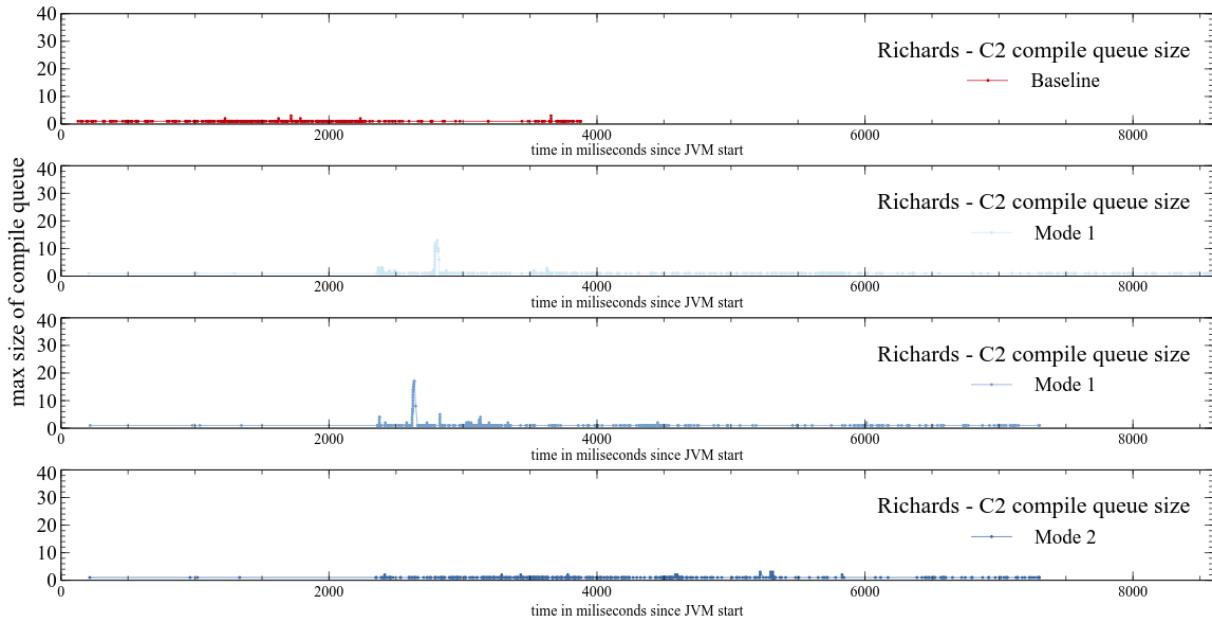


Figure 4.12: C2 Compile queue size over time Octane Richards benchmark

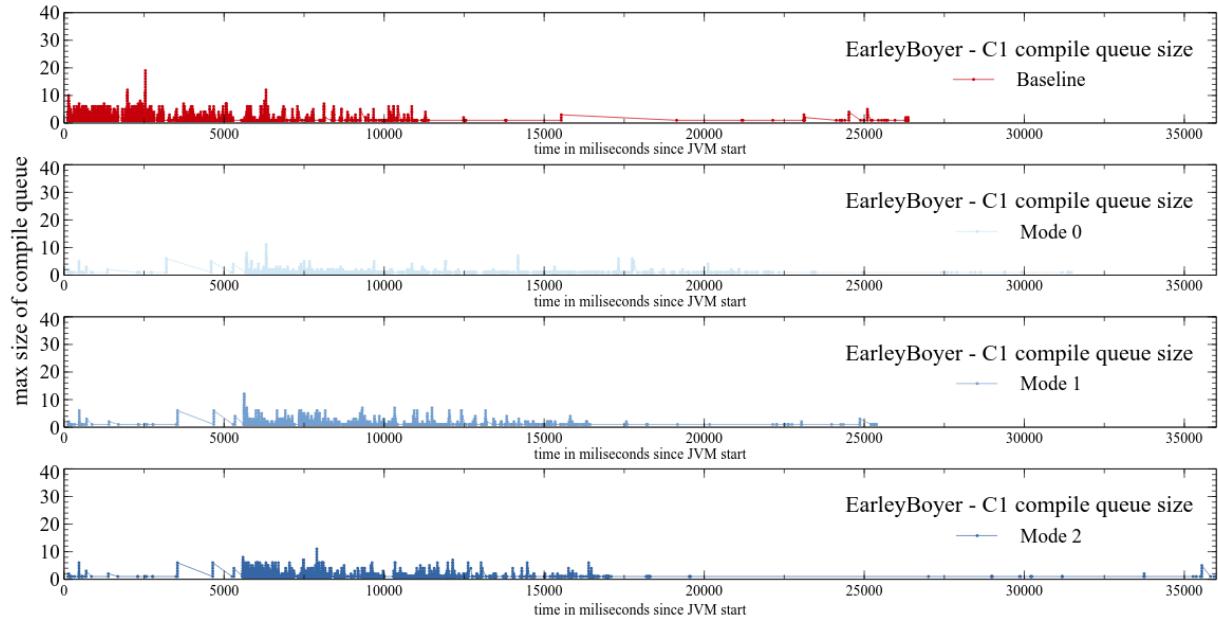


Figure 4.13: C1 Compile queue size over time Octane EarleyBoyer benchmark

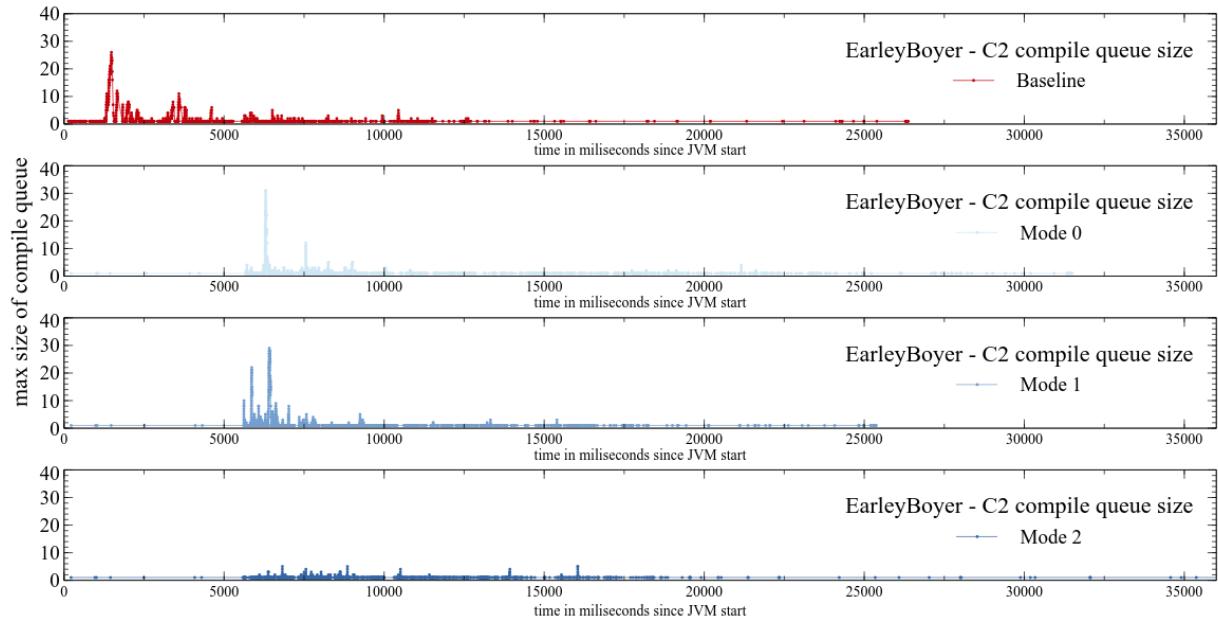


Figure 4.14: C2 Compile queue size over time Octane EarleyBoyer benchmark

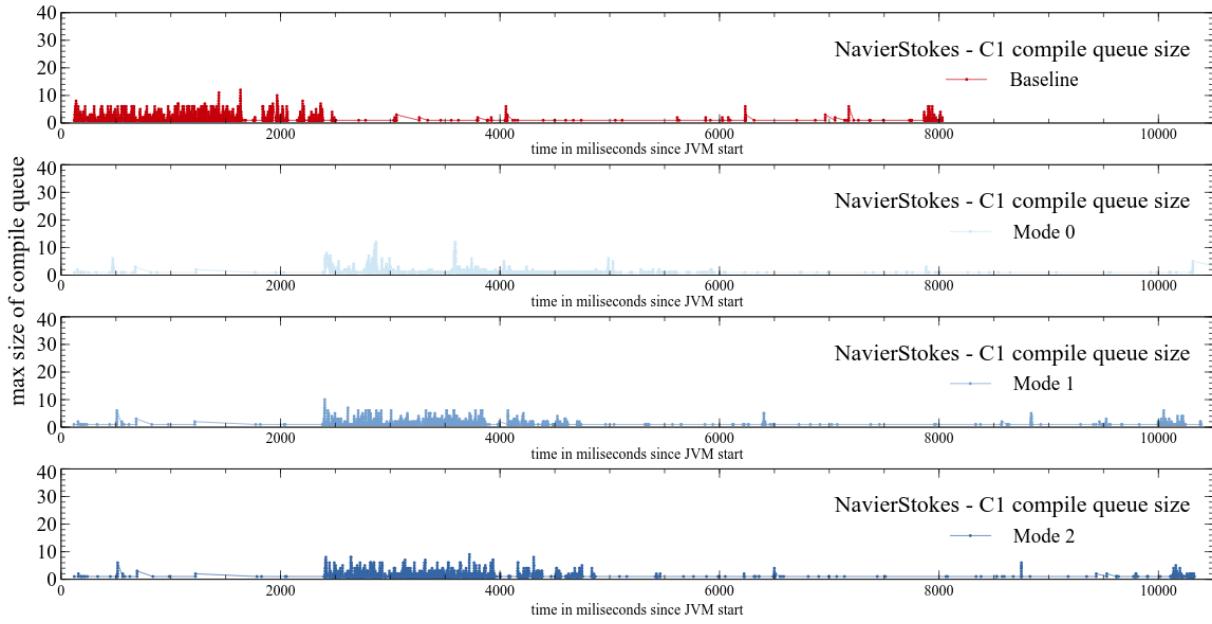


Figure 4.15: C1 Compile queue size over time Octane NavierStokes benchmark

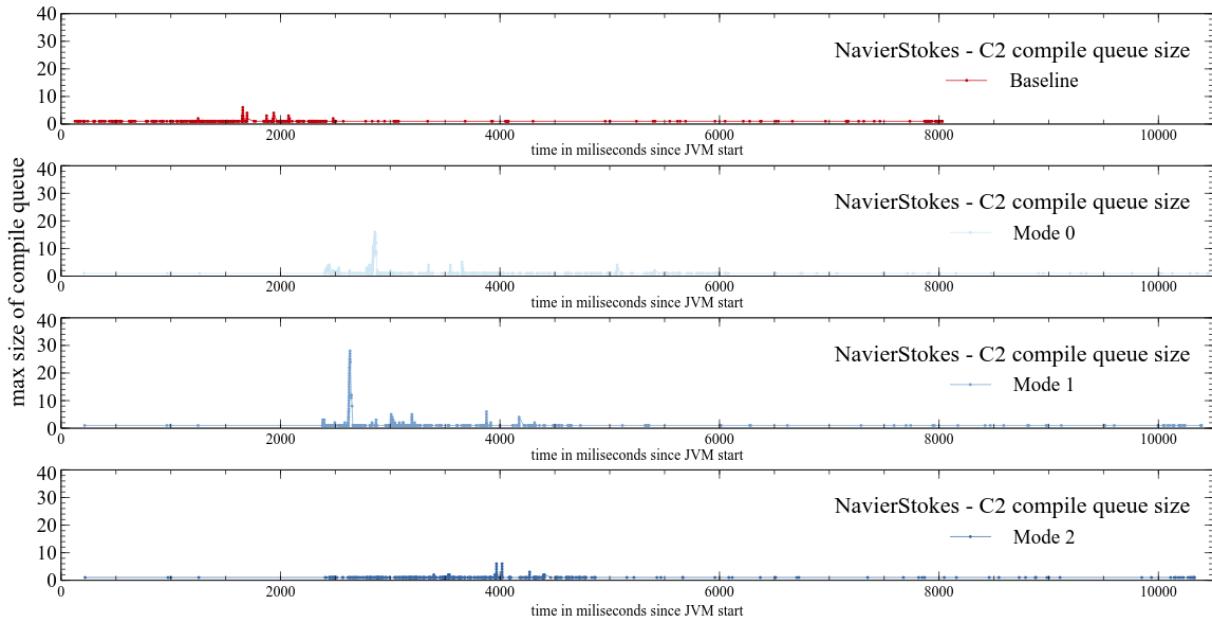


Figure 4.16: C2 Compile queue size over time Octane NavierStokes benchmark

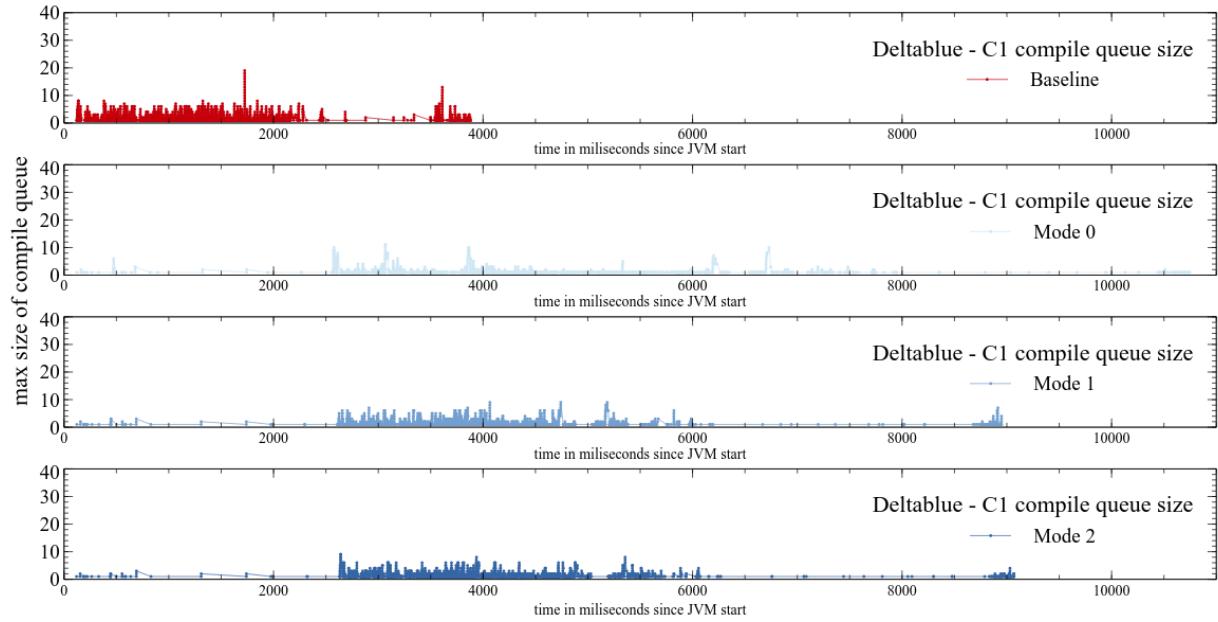


Figure 4.17: C1 Compile queue size over time Octane Deltablue benchmark

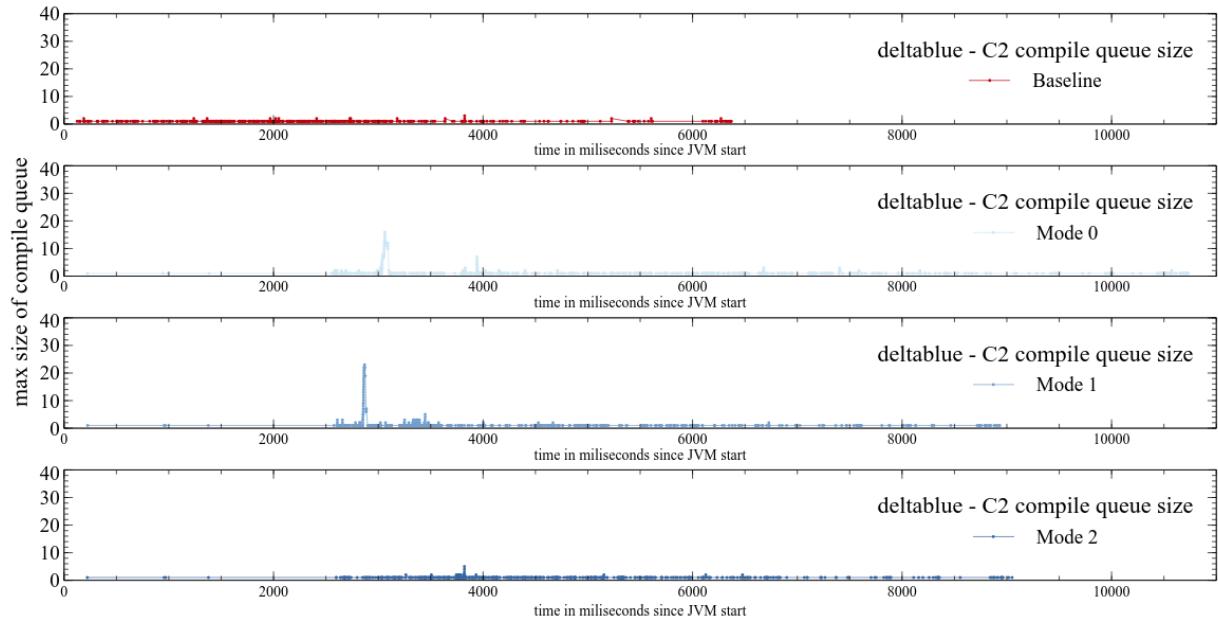


Figure 4.18: C2 Compile queue size over time Octane Deltablue benchmark

4.5 Number and type of compilations

In this section we take a look on how cached profile modify the ratio of C1 and C2 compilations and if there is a correlation between percentage of methods using cached profiles and the resulting performance.

We continue focusing on the 6 individual benchmarks selected in Section 4.4. We use the new HotSpot flag `-XX:+PrintCacheProfiles` that prints out the Level of each compilation and whether or not it uses cached profiles.

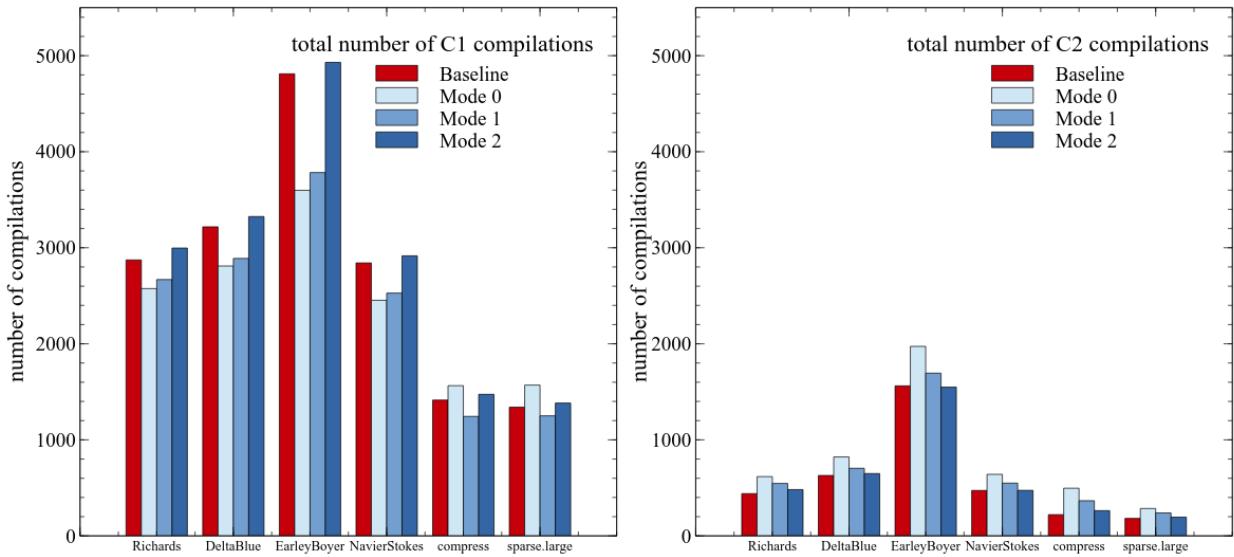


Figure 4.19: Number of compilations for some specJVM and octane benchmarks

Figure 4.19 shows the total amount of compilations split in C1 and C2. We see that the Octane benchmarks and the SPECjvm benchmarks behave differently. While the 4 Octane ones achieve a lower amount of C1 compilations in `Mode 0` and `Mode 1`, `Mode 2` is similar to the baseline. The two SPECjvm benchmarks have more C1 compilations in `Mode 0`, less in `Mode 1` and the same amount in `Mode 2` compared to the baseline.

The changes in the amount of C2 compilations are very similar in all benchmarks. Using `Mode 0` and `Mode 1` results in more C2 compilations than the baseline and `Mode 2` again achieves around the same.

If we recall the differences between the modes these results make sense. `Mode 0` lowers the thresholds of C1 compilations in case the method has a cached profile and compiles with C2 instead. This reduces the number of C1 compilations in favor of more C2 compilations. `Mode 1` does not lower the thresholds but it still promotes some C1 compilations to C2 compilations which is also seen in

the numbers. Mode 2 leaves the tiered compilation completely untouched, we see that it also does not influence the total number of compilations significantly.

Furthermore, we are interested how many of the compilations make use of cached profiles. Ideally, this should be 100% but there could always be methods that have no cached profiles available because they have not been compiled when generating the profile. Additionally, we intentionally do not use any profiles for compilation Level 1 and Level 2 as described in Section 3.3. The ReplayCompilation functionality does not support certain methods, e.g. lambda expressions. Since this CacheProfiles implementation is based on ReplayCompilation it will also not compile these methods using cached profiles.

In Figure 4.20 to Figure 4.25 we show cake diagrams that visualize the portion of specific compilation types. The results show, that the proportions are very similar if we compare different benchmarks of the same CacheProfilesMode. Mode 0 and Mode 1 look similar with Mode 1 invoking less compilations using cached profiles. In Mode 2 we see Level 2 compilations appearing due to the changed tiered compilation transitions. The number of Level 3 compilations is nearly unchanged compared to Mode 1 because these are compilations of methods where no profiles from C2 compilations are available. The Level 2 compilations only happen if such a profile is available and usually result in another Level 4 compilation.

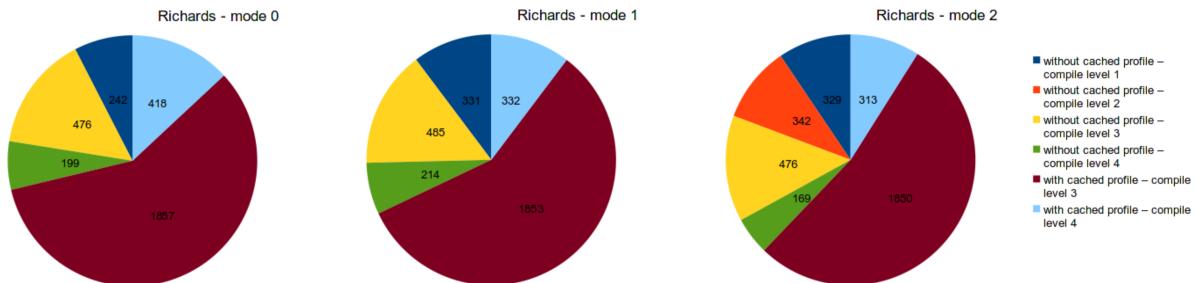


Figure 4.20: Ratio of compilations Octane Richards benchmark

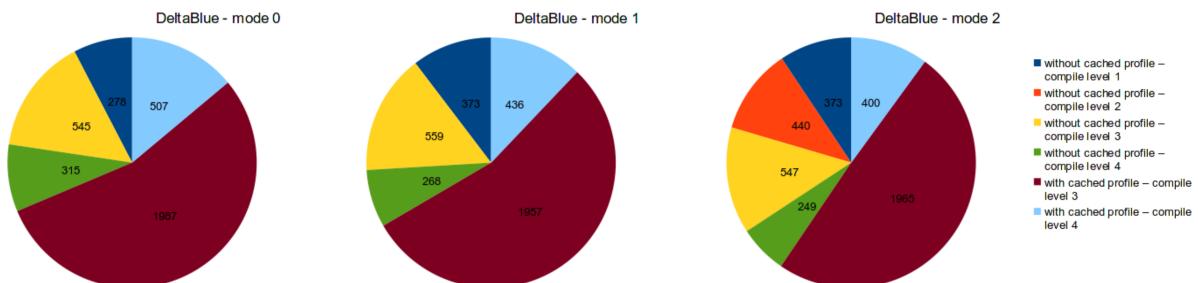


Figure 4.21: Ratio of compilations Octane Deltablue benchmark

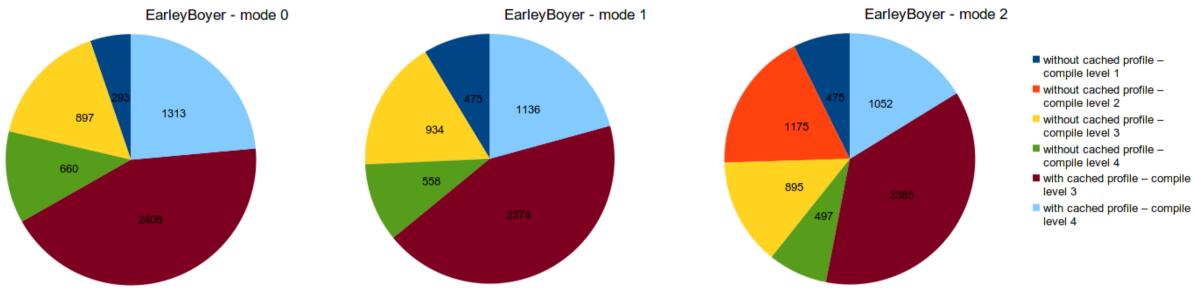


Figure 4.22: Ratio of compilations Octane EarleyBoyer benchmark

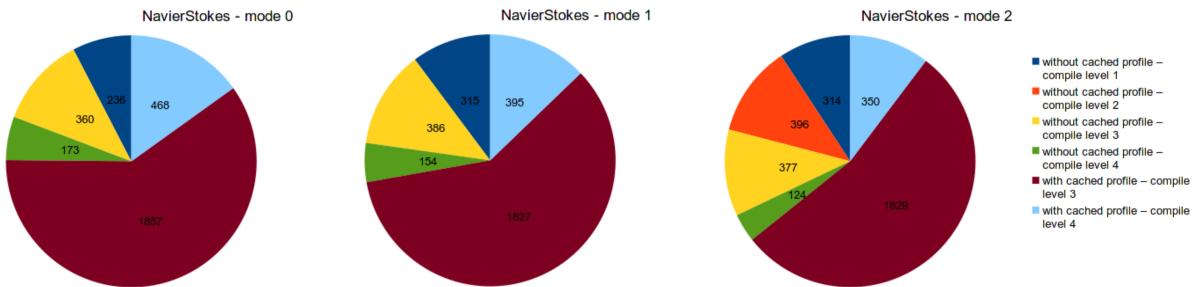


Figure 4.23: Ratio of compilations Octane NavierStokes benchmark

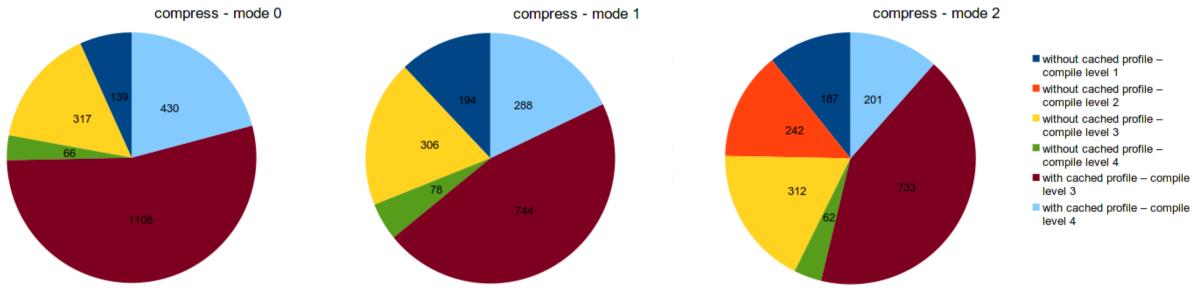


Figure 4.24: Ratio of compilations SPECjvm compress benchmark

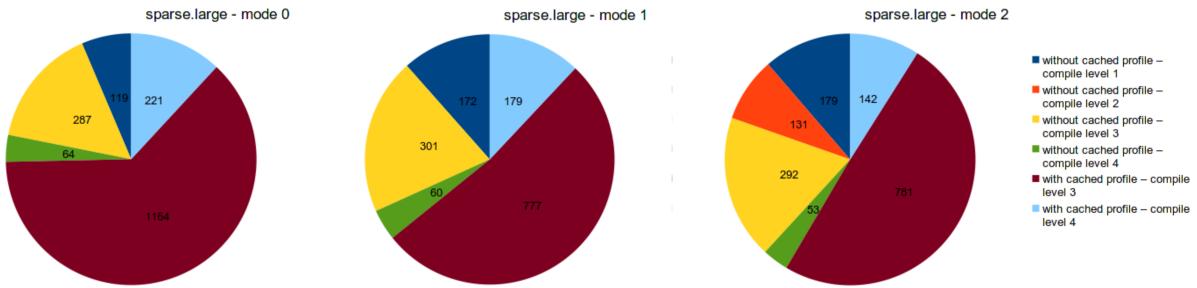


Figure 4.25: Ratio of compilations SPECjvm sparse.large benchmark

4.6 Effect of interpreter profiles

Our system makes use of two types of cached profiles. Profiles that are gathered by the interpreter and used by the C1 compiler and profiles that are gathered by a C1 compiled method and used when compiling with C2.

We added a HotSpot flag that allows us to specify the minimum level of a compilation that dumps profiles (`-XX:DumpProfilesMinTier=level`). Previous measurements were done using Level 3 which dumps profiles during C1 and C2 compilations.

However, we are also interested how the system performance changes when only C2 compiler profiles are used. The system will then only use cached profiles where a C2 compilation took place in the previous profile generation run. We use the same setup as before and run the individual SPECjvm (see Figure 4.26 and Octane (see Figure 4.27 benchmarks.

Most of the benchmarks do not show significantly different results compared to Section 4.2, where both types of cached profiles where used. There are a few benchmarks where individual modes now improve the performance while having a performance drop when using C1 and C2 profiles (e.g. NavierStokes Mode 0). But we also experience the other way around for example in benchmark Splay Mode 2. In these individual cases, we believe that for example a benchmarks C1 compilation does not profit from having cached profiles and therefore using them will even decrease performance. This happens because methods are first compiled by using cached profiles and after ten deoptimizations are compiled using generated profiles.

The results let us conclude that the performance differences to the baseline are mostly due to the code quality of C2 compilations. Even though the number of C1 compilations is usually a lot higher than the number of C2 compilations, C2 compilations seem more crucial to the methods performance. Since C2 is the maximum compilation level, a program might spend most of it's time in C2 compiled code.

This again shows that it is very hard to predict JVM performance.

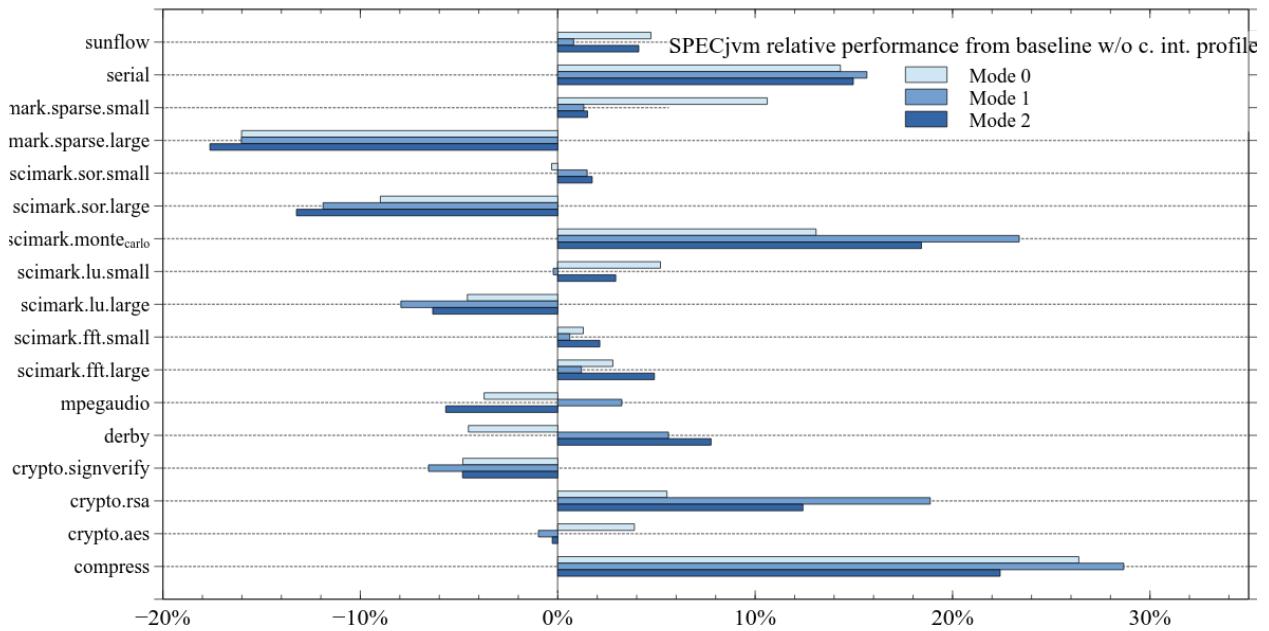


Figure 4.26: Relative performance from baseline for all SPECjvm benchmarks without using cached interpreter profiles

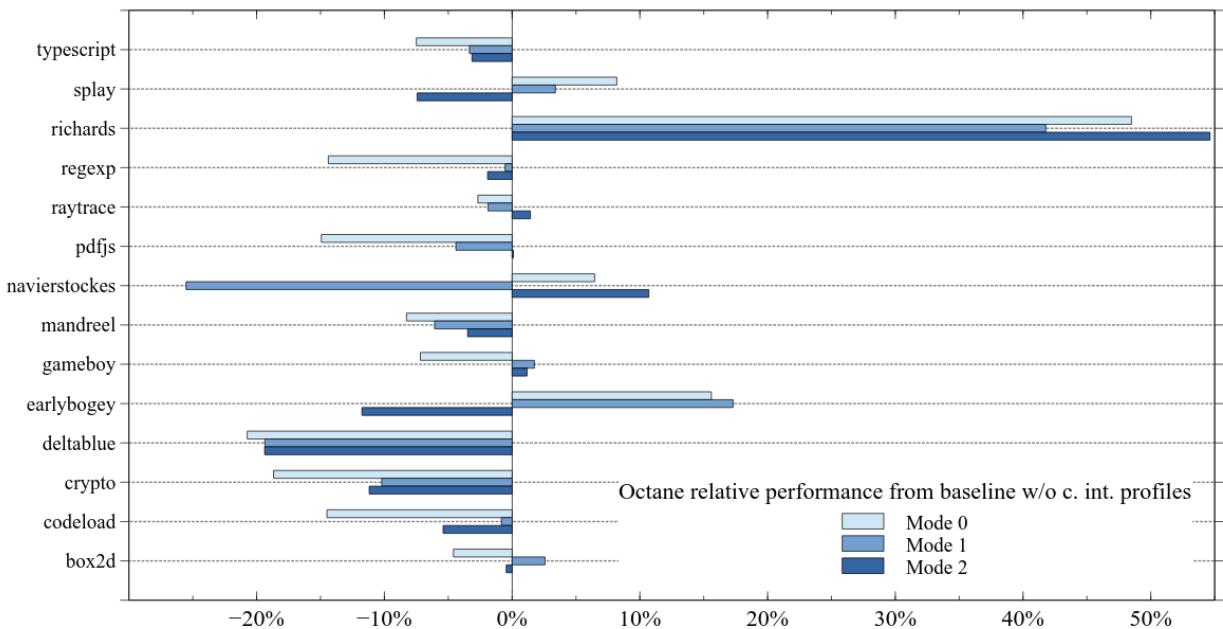


Figure 4.27: Relative performance from baseline for all Octane benchmarks without using cached interpreter profiles

4.7 Effect of intrinsified methods

Most modern JVMs implement use *method intrinsics* to further optimize commonly used Java core library methods [7]. This means that the JIT compiler will not compile the method based on the Java bytecode, but instead replace it with a predefined, hand-written, and manually optimized assembly code snippet. The current list of methods where intrinsics are available can be found in the code reference [11].

Intrinsics are mostly used in C1 and C2 compilations and the emitted code is independent of the current available profiling information. This means, if many methods of a benchmark are intrinsified, the influence of profiles and therefore cached profiles as well decreases. We want to know whether this has a large influence on the performance of our cached profiles. An intrinsified method has no advantage of having a rich profiles available but will still be influenced by modified compilation thresholds. E.g. lowering the threshold will intrinsify methods earlier and therefore speed up execution, independent of the cached profile.

The results of both benchmark suites with disabled method intrinsics can be found in Figures X and Figure Y. We see that there are small performance differences in individual benchmarks but we can not conclude a major influence to the behavior of cached profiles and their influence on performance.

5 Possible improvements

While this thesis provides a first look on implementing a system to reuse profiles from previous JVM runs it leaves room for several improvements.

- There is a time overhead for parsing the cached profiles file when the JVM boots up. This could be lowered by finding more compact ways of storing profiles to disk or storing the cached profiles in main memory.
- The data structure of the cached profiles is a simple C heap array. The lookup is $O(n)$ with n being the number of cached method compilations. This could be improved to $O(\log(n))$ by using a more complex tree-like data structure.
- Additionally, all method compilations get dumped and stored. The system could be improved in a way, so only the last compilation record of a method is kept in the file. This would decrease the size of the cached profile file and decrease the parsing time.
- Currently, only the profiles from a single run are used. A possible improvement is to use multiple executions for gathering the cached profiles and merge the profiling information to achieve even more complete profiles.
- We think it could also be beneficial to be able to modify the cached profile. That would allow the JVM user to manually improve profiling information by using his knowledge of the method execution which might not be available to the compiler.

6 Conclusion

Current Java Virtual Machines (JVM) like HotSpot gather profiling information about executed methods to improve the quality of the compiled code. This thesis presents a way to cache these profiling information, so they can be reused in future runs of the JVM.

The expected benefits are a faster warmup of the JVM because the JVM does not need to spend time profiling the code and can use cached profiles directly. Furthermore, since the cached profiles originate from previous compilations where extensive profiling already happened, compilations using these profiles produce more optimized code which decreases the amount of deoptimizations.

We show, using two benchmark suites, that cached profiles can indeed improve warmup performance and significantly lower the amount of deoptimizations.

The ideas and functionality is implemented in the HotSpot JVM (openJDK9). It provides the user of the JVM several choices to use the system and allows fine-grained selection of the cached methods. We believe that this can be a valuable asset in scenarios where a fast JVM warmup is appreciated and greatly help to decrease performance fluctuations of JIT compiled code.

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

A.2 Cached profile example

Listing A.1: Example of cached profiling information

A.3 Code changes

Table A.1: Code lines changed, inserted, deleted, modified and unchanged

class	changed	inserted	deleted	modified	unchanged
src/share/vm/ci/ciClassList.hpp	2	2	0	0	121
src/share/vm/ci/ciEnv.cpp	58	56	0	2	1283
src/share/vm/ci/ciEnv.hpp	5	5	0	0	469

TODO

A.4 SPECjvm benchmark

This list gives a short description of the benchmarks that are part of the SPECjvm 2008 Benchmark Suite. The list is directly taken from <https://www.spec.org/jvm2008/docs/benchmarks/index.html> and put in as a reference.

- **Compress:** This benchmark compresses data, using a modified Lempel-Ziv method (LZW). Basically finds common substrings and replaces them with a variable size code. This is

deterministic, and can be done on the fly. Thus, the decompression procedure needs no input table, but tracks the way the table was built. Algorithm from "A Technique for High Performance Data Compression", Terry A. Welch, IEEE Computer Vol 17, No 6 (June 1984), pp 8-19.

This is a Java port of the 129.compress benchmark from CPU95, but improves upon that benchmark in that it compresses real data from files instead of synthetically generated data as in 129.compress.

- **Crypto:** This benchmark focuses on different areas of crypto and are split in three different sub-benchmarks. The different benchmarks use the implementation inside the product and will therefore focus on both the vendor implementation of the protocol as well as how it is executed.
 - **aes** encrypt and decrypt using the AES and DES protocols, using CBC/PKCS5Padding and CBC/NoPadding. Input data size is 100 bytes and 713 kB.
 - **rsa** encrypt and decrypt using the RSA protocol, using input data of size 100 bytes and 16 kB.
 - **signverify** sign and verify using MD5withRSA, SHA1withRSA, SHA1withDSA and SHA256withRSA protocols. Input data size of 1 kB, 65 kB and 1 MB.
- **Derby:** This benchmark uses an open-source database written in pure Java. It is synthesized with business logic to stress the BigDecimal library. It is a direct replacement to the SPECjvm98 db benchmark but is more capable and represents as close to a "real" application. The focus of this benchmark is on BigDecimal computations (based on telco benchmark) and database logic, especially, on locks behavior. BigDecimal computations are trying to be outside 64-bit to examine not only 'simple' BigDecimal, where 'long' is used often for internal representation.
- **MPEGaudio:** This benchmark is very similar to the SPECjvm98 mpegaudio. The mp3 library has been replaced with JLayer, an LGPL mp3 library. Its floating-point heavy and a good test of mp3 decoding. Input data were taken from SPECjvm98.
- **Scimark:** This benchmark was developed by NIST and is widely used by the industry as a floating point benchmark. Each of the subtests (**fft**, **lu**, **monte_carlo**, **sor**, **sparse**) were incorporated into SPECjvm2008. There are two versions of this test, one with a **large**dataset (32Mbytes) which stresses the memory subsystem and a **small**dataset which stresses the JVMs (512Kbytes).
- **Serial:** This benchmark serializes and deserializes primitives and objects, using data from the JBoss benchmark. The benchmark has a producer-consumer scenario where serialized objects are sent via sockets and deserialized by a consumer on the same system. The benchmark heavily stress the Object.equals() test.

- **Sunflow:** This benchmark tests graphics visualization using an open source, internally multi-threaded global illumination rendering system. The sunflow library is threaded internally, i.e. it's possible to run several bundles of dependent threads to render an image. The number of internal sunflow threads is required to be 4 for a compliant run. It is however possible to configure in property specjvm.benchmark.sunflow.threads.per.instance, but no more than 16, per sunflow design. Per default, the benchmark harness will use half the number of benchmark threads, i.e. will run as many sunflow benchmark instances in parallel as half the number of hardware threads. This can be configured in specjvm.benchmark.threads.sunflow.
- **XML:** This benchmark has two sub-benchmarks: XML.transform and XML.validation. XML.transform exercises the JRE's implementation of javax.xml.transform (and associated APIs) by applying style sheets (.xsl files) to XML documents. The style sheets and XML documents are several real life examples that vary in size (3KB to 156KB) and in the style sheet features that are used most heavily. One "operation" of XML.transform consists of processing each style sheet / document pair, accessing the XML document as a DOM source, a SAX source, and a Stream source. In order that each style sheet / document pair contribute about equally to the time taken for a single operation, some of the input pairs are processed multiple times during one operation.

Result verification for XML.transform is somewhat more complex than for other of the benchmarks because different XML style sheet processors can produce results that are slightly different from each other, but all still correct. In brief, the process used is this. First, before the measurement interval begins the workload is run once and the output is collected, canonicalized (per the specification of canonical XML form) and compared with the expected canonicalized output. Output from transforms that produce HTML is converted to XML before canonicalization. Also, a checksum is generated from this output. Inside the measurement interval the output from each operation is only checked using the checksum.

XML.validation exercises the JRE's implementation of javax.xml.validation (and associated APIs) by validating XML instance documents against XML schemata (.xsd files). The schemata and XML documents are several real life examples that vary in size (1KB to 607KB) and in the XML schema features that are used most heavily. One "operation" of XML.validation consists of processing each style sheet / document pair, accessing the XML document as a DOM source and a SAX source. As in XML.transform, some of the input pairs are processed multiple times during one operation so that each input pair contributes about equally to the time taken for a single operation.

A.5 Octane benchmark

What follows is an overview of the benchmarks Octane consists of. The original list can be found on <https://developers.google.com/octane/benchmark>.

- **Richards:** OS kernel simulation benchmark, originally written in BCPL by Martin Richards (539 lines).

- Main focus: *property load/store, function/method calls*
- Secondary focus: *code optimization, elimination of redundant code*
- **Deltablue:** One-way constraint solver, originally written in Smalltalk by John Maloney and Mario Wolczko (880 lines).
 - Main focus: *polymorphism*
 - Secondary focus: *OO-style programming*
- **Raytrace:** Ray tracer benchmark based on code by Adam Burmister (904 lines).
 - Main focus: *argument object, apply*
 - Secondary focus: *prototype library object, creation pattern*
- **Regexp:** Regular expression benchmark generated by extracting regular expression operations from 50 of the most popular web pages (1761 lines).
 - Main focus: *Regular expressions*
- **NavierStokes:** 2D NavierStokes equations solver, heavily manipulates double precision arrays. Based on Oliver Hunt’s code (387 lines).
 - Main focus: *reading and writing numeric arrays.*
 - Secondary focus: *floating point math.*
- **Crypto:** Encryption and decryption benchmark based on code by Tom Wu (1698 lines).
 - Main focus: *bit operations*
- **Splay:** Data manipulation benchmark that deals with splay trees and exercises the automatic memory management subsystem (394 lines).
 - Main focus: *Fast object creation, destruction*
- **SplayLatency:** The Splay test stresses the Garbage Collection subsystem of a VM. SplayLatency instruments the existing Splay code with frequent measurement checkpoints. A long pause between checkpoints is an indication of high latency in the GC. This test measures the frequency of latency pauses, classifies them into buckets and penalizes frequent long pauses with a low score.
 - Main focus: *Garbage Collection latency*
- **EarleyBoyer:** Classic Scheme benchmarks, translated to JavaScript by Florian Loitsch’s Scheme2Js compiler (4684 lines).
 - Main focus: *Fast object creation, destruction*
 - Secondary focus: *closures, arguments object*

- **pdf.js:** Mozilla’s PDF Reader implemented in JavaScript. It measures decoding and interpretation time (33,056 lines).
 - Main focus: *array and typed arrays manipulations.*
 - Secondary focus: *math and bit operations, support for future language features (e.g. promises)*
- **Mandreel:** Runs the 3D Bullet Physics Engine ported from C++ to JavaScript via Mandreel (277,377 lines).
 - Main focus: *emulation*
- **MandreelLatency:** Similar to the SplayLatency test, this test instruments the Mandreel benchmark with frequent time measurement checkpoints. Since Mandreel stresses the VM’s compiler, this test provides an indication of the latency introduced by the compiler. Long pauses between measurement checkpoints lower the final score.
 - Main focus: *Compiler latency*
- **GB Emulator:** Emulates the portable console’s architecture and runs a demanding 3D simulation, all in JavaScript (11,097 lines).
 - Main focus: *emulation*
- **Code loading:** Measures how quickly a JavaScript engine can start executing code after loading a large JavaScript program, social widget being a common example. The source for this test is derived from open source libraries (Closure, jQuery) (1,530 lines).
 - Main focus: *JavaScript parsing and compilation*
- **Box2DWeb:** Based on Box2DWeb, the popular 2D physics engine originally written by Erin Catto, ported to JavaScript. (560 lines, 9000+ de-minified)
 - Main focus: *floating point math.*
 - Secondary focus: *properties containing doubles, accessor properties.*
- **TypeScript:** Microsoft’s TypeScript compiler is a complex application. This test measures the time TypeScript takes to compile itself and is a proxy of how well a VM handles complex and sizable Javascript applications (25,918 lines).
 - Main focus: *run complex, heavy applications*

A.6 Additional graphs

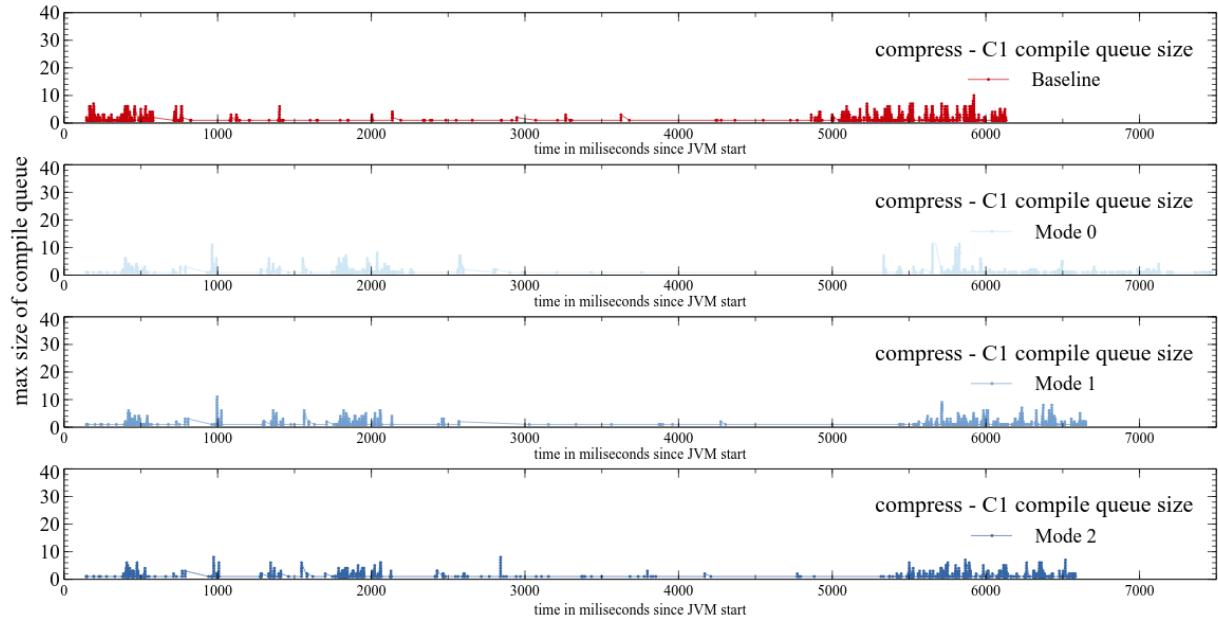


Figure A.1: C1 Compile queue size over time SPECjvm compress benchmark

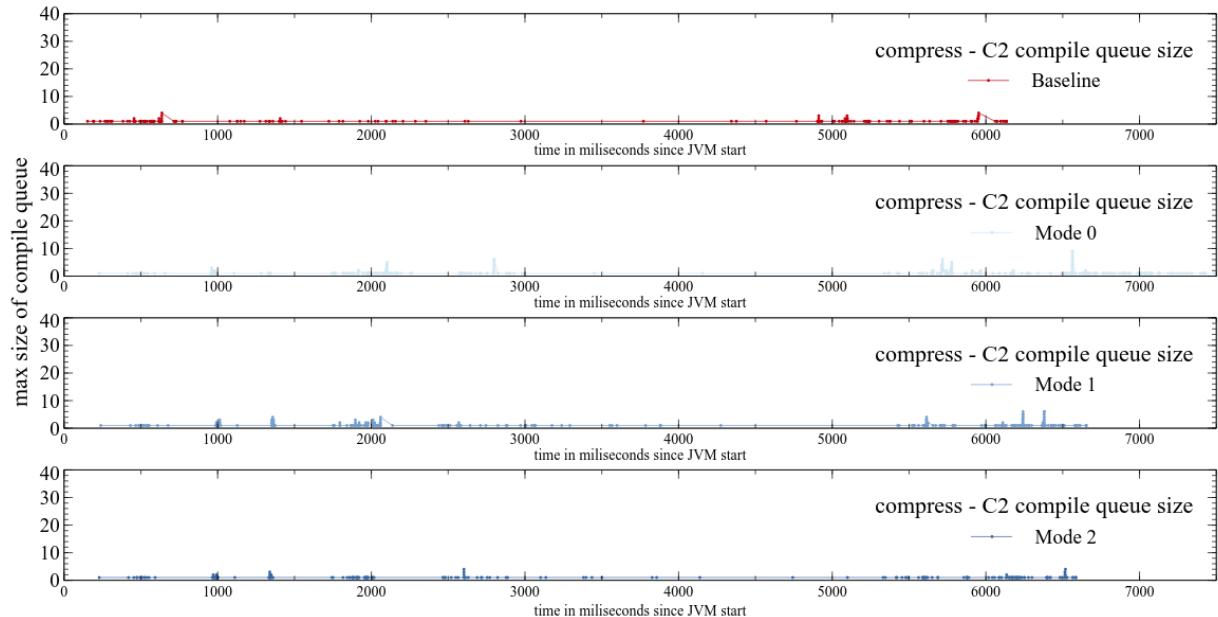


Figure A.2: C2 Compile queue size over time SPECjvm compress benchmark

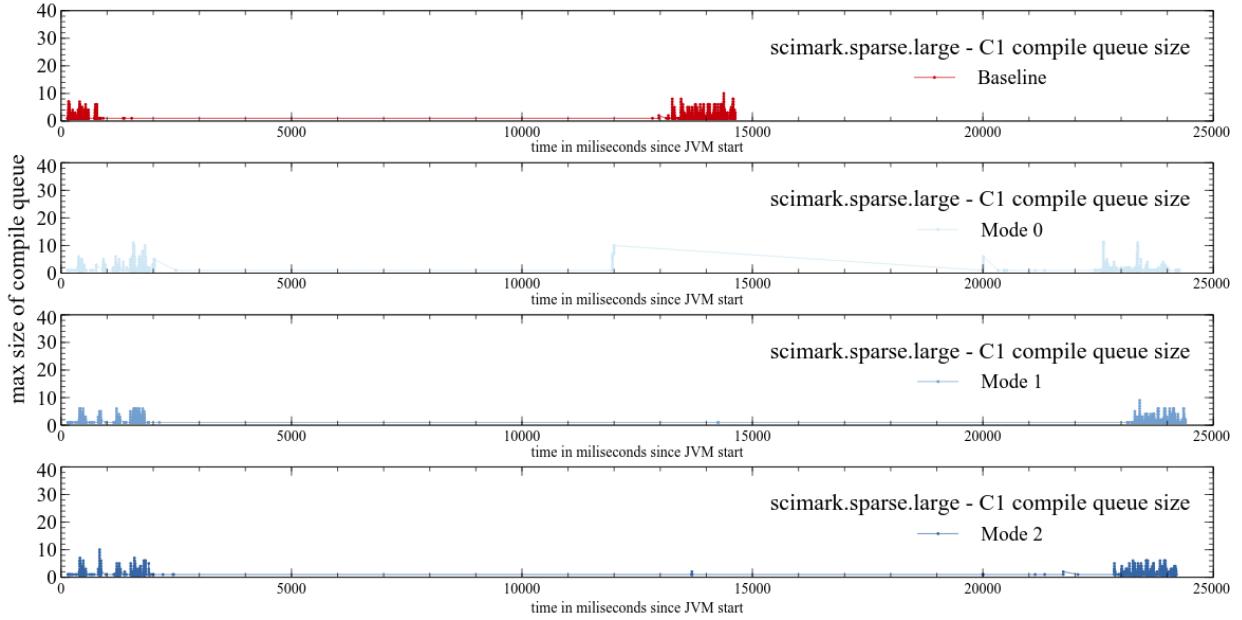


Figure A.3: C1 Compile queue size over time SPECjvm scimark.sparse.large benchmark

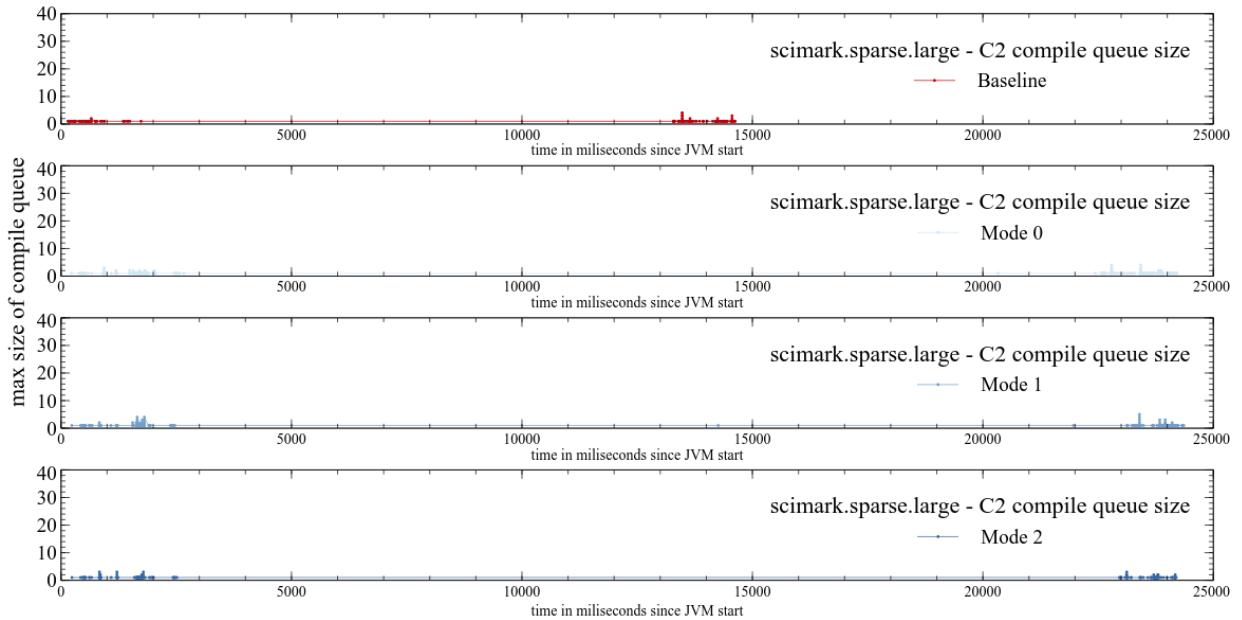


Figure A.4: C2 Compile queue size over time SPECjvm scimark.sparse.large benchmark

Bibliography

- [1] Azul Systems. ReadyNow! Technology for Zing. <http://www.azulsystems.com/solutions/zing/readynow>, 2015.
- [2] Azul Systems. Zing: Java for the Real Time Business. <http://www.azulsystems.com/products/zing/whatisit>, 2015.
- [3] ETH. ETH Data Center Observatory. <https://wiki.dco.ethz.ch/bin/viewauth/Main/Cluster>, 2015.
- [4] Google. Octane 2.0 Benchmark. <https://developers.google.com/octane/>, 2015.
- [5] 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.
- [6] V. Ivanov. JIT-compiler in JVM seen by a Java developer. <http://www.stanford.edu/class/cs343/resources/java-hotspot.pdf>, 2013.
- [7] K. Mo. Intrinsic Methods in HotSpot. <http://www.oracle.com/technetwork/java/javase/tech/3198-d1-150056.pdf>, 2013. Talk from GreenTeaJUG event.
- [8] Oracle Corporation. Code for advancedThresholdPolicy.hpp. <http://hg.openjdk.java.net/jdk9/hs-comp/hotspot/file/63337cc98898/src/share/vm/runtime/advancedThresholdPolicy.hpp>, 2013.
- [9] Oracle Corporation. hprof - A heap/CPU profiling tool. <https://docs.oracle.com/javase/7/docs/technotes/samples/hprof.html>, 2014.
- [10] Oracle Corporation. javac - Java programming language compiler. <https://docs.oracle.com/javase/7/docs/technotes/tools/windows/javac.html>, 2014.
- [11] Oracle Corporation. Intrinsified methods in openJDK 9. <http://hg.openjdk.java.net/jdk9/jdk9/hotspot/file/1aef080fd28d/src/share/vm/classfile/vmSymbols.hpp#l1604>, 2015.
- [12] M. Paleczny, C. Vick, and C. Click. The Java HotSpot™Server Compiler. https://www.usenix.org/legacy/events/jvm01/full_papers/paleczny/paleczny.pdf, 2001. Paper from JVM '01.

- [13] T. Rodriguez and K. Russell. Client Compiler for the Java HotSpot™Virtual Machine: Technology and Application. <http://www.oracle.com/technetwork/java/javase/tech/3198-d1-150056.pdf>, 2002. Talk from JavaOne 2002.
- [14] Standard Performance Evaluation Corporation. SPECjvm2008 Benchmark. <https://www.spec.org/jvm2008/>, 2008.

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

First name(s):

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.