



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

MASTER THESIS

Bc. Jakub Háva

**Monitoring Tool for Distributed Java
Applications**

Department of Distributed and Dependable Systems

Supervisor of the master thesis: RNDr. Pavel Parízek, Ph.D

Study programme: Computer Science

Study branch: Software Systems

Prague 2017

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

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

In date

signature of the author

Title: Monitoring Tool for Distributed Java Applications

Author: Bc. Jakub Háva

Department: Department of Distributed and Dependable Systems

Supervisor: RNDr. Pavel Parízek, Ph.D, Department of Distributed and Dependable Systems

Abstract: The main goal of this thesis is to create a monitoring platform and library that can be used to monitor distributed Java-based applications. This work is inspired by Google Dapper and shares a concept called "Span" with it. Spans represent a small specific part of the computation and are used to capture state among multiple communicating hosts. In order to be able to collect spans without recompiling the original application's code, instrumentation techniques are highly used in the thesis. The monitoring tool, which is called Distrace, consists of two parts: the native agent and instrumentation server. Users of this platform are supposed to extend the instrumentation server and specify the points in their application's code where new spans should be created and closed. In order to achieve high performance and affect the running application at least as possible, the instrumentation server is used for instrumenting the code. The tool is aimed to have a small foot-print on the monitored applications, should be easy to deploy and is transparent to target applications from the point of view of the final user.

Keywords: monitoring, cluster, instrumentation, distributed systems, performance

I would like to thank my thesis supervisor Dr. Pavel Parízek for leading me throughout the whole thesis and willing to help with any concerns I've ever had. I would also like to thank H2O.ai for being able to write this thesis under their coordination, particularly to Dr. Michal Malohlava for providing very useful technical advice.

Contents

1	Introduction	3
1.1	Project Goals	4
1.2	Thesis Outline	5
2	Background	6
2.1	Cluster Monitoring Tools	6
2.2	HTrace	10
2.3	Tools for Large-Scale Debugging	10
2.4	Profiling Tools	11
2.5	Byte Code Manipulation Libraries	14
2.6	Communication Middleware	18
2.7	Java Libraries	21
2.8	Logging Libraries	26
2.9	Docker	26
3	Analysis	27
3.1	Limitations of Similar Solutions	27
3.2	Small Footprint and High Performance	27
3.3	Transparency and Universality	30
3.4	Easiness of Use	31
3.5	Easiness of Deployment	32
3.6	Modularity	32
3.7	Summary of Features	33
4	Design	34
4.1	Overview	34
4.2	Example Use Case	35
4.3	Spans and Trace Trees	38
4.4	Native Agent	45
4.5	Instrumentation Server	50
4.6	User Interface	55
5	Implementation Details	59
5.1	Native Agent	59
5.2	Span Injection on Instrumentation Server	61
5.3	Determine the Current Span Saver	61
6	Big Example	63
6.1	H ₂ O In More Details	63
6.2	Building the Core Server and Native Agent	67
6.3	Extending the Core Instrumentation Server	67
6.4	Configuring and Running the Application	70
6.5	Interpreting the Results	70

7	Evaluation	71
7.1	Known Limitations	71
7.2	Platform demonstration	71
7.3	Comparison to Related Work	71
7.4	Future plans	71
8	Conclusion	72
	Bibliography	73
	List of Figures	74
	List of Tables	75
	List of Abbreviations	76
	Attachments	77

1. Introduction

Lately, the volume of data applications need to handle is significantly increasing. In order to support this scaling trend, the applications are becoming distributed for reasons of scalability, stability and availability. Not every task may be solved efficiently by distributed applications, however when it comes to big data, the computational requirements may be higher than single physical node can fulfill. Such distributed applications may run on multiple physical or virtual machines in order to achieve the best performance and the ability to process data significantly large. For this, computation clusters are created where the user interacts with the application as it would be running locally and the cluster should handle the distributed computation internally.

However, with growth of distributed applications there is also increasing demand for monitoring and debugging such applications. Analyzing applications in distributed environment is inherently more complex task comparing to single-node applications where well-known debugging and profiling techniques may be used. Analysis of single-node applications usually focuses on a single standalone applications where most of the information required to reason about it is collected directly from the application. In case of distributed application it is desired to collect the same information as on single-node applications plus, and more importantly, the shared state between the communicating nodes. For instance, an error may occur on one of the computation nodes in the cluster and over time, more and more nodes can become affected. By collecting the relations between the nodes, the analysis tool may be able to use the information to determine where the error initially occurred and how it spread over the time.

Simple solution comes to mind to address this issue. Monitoring or debugging tools used for single-node applications may be attached per each application node and collect the information from the nodes separately on each other. This solution does not require any additional tools however the state between the application nodes would not be preserved unless the monitored application is already designed to send the required information. Most of the applications is not designed to transfer the information used for analysis of the application itself for several reasons. It may be hard or unwanted to design the application in a way that all the information required for analysis are already transferred between communication nodes. New analysis method may require new metrics which in this case would also mean recompiling and new deployment of the application.

For this reason several new monitoring and debugging tools have been developed. These tools are usually built on the code instrumentation technique. This method is used to alter the monitored application's code at run-time in order to collect all relevant information. The significant advantage of this method is that the original application does not have to be changed in order to add additional metrics. Usually, such tools use the instrumentation technique to add special information to the code that is later used to build a so-called distributed stack-trace. A stack-trace in single node application represents the call hierarchy of a method at the given moment. Distributed stack-trace is a very similar concept except that the dependencies between different nodes are preserved and can be seen on the collected stack-trace as well. Therefore, distributed stack-traces give

the possibility to see the desired relations between the applications node. Google Dapper and Zipkin are the most significant available monitoring tools and are discussed later in the thesis.

1.1 Project Goals

This thesis introduces Distrace, a monitoring tool with the similar purpose, sharing some of the concepts mentioned above, however the goals of this work are different and should give the user a new generic and high-performance way how to monitor distributed applications. Main goals of the thesis are to create an open-source generic monitoring library with small footprint on the monitored applications whilst giving the user the possibility to use high-level programming language to define instrumentation points.

Main requirements for the platform are:

- **Small Footprint**

The mentioned cluster monitoring tools can affect the application performance and memory consumption since they perform instrumentation in the same virtual machine as the monitored application. The project is required to have a minimal footprint on the monitored application. The instrumentation is needed in order to inject special information about spans to the application's code, where spans are structures used to collect the shared state between the application nodes. This information is later used for span collection and associating the relationships between spans.

- **Transparency and Universality**

These two requirements are contradictory to each other. The universal tool which could be used for monitoring of majority of applications would either collect just basic information shared about all applications or the user would be required to manually specify the information to be collected specific to the application, which leads to loose of the application-level transparency. This tool tries to find a compromise between these two goals and attempts to support high-level of universality with minimizing the impact on the application itself. The project needs to do some trade-offs between the application-level transparency and the universality of the platform. The goal is to be able to instrument JVM-based applications to minimize the impact on the application's code itself.

- **Easiness of Use**

The application should use high-level programming language for the instrumentation and specifying additional information to be collected. The users of this tool are supposed to work with Java-based language and should not be required to have deeper knowledge about internal Java Virtual Machine structure.

- **Easiness of Deployment**

The complexity of deployment of this tool is also a significant aspect of the tool. In order for developers and testers to use this tool frequently, its deployment and usage has to be relatively easy. This requirement has two

sub-parts. Minimizing the configuration of the monitoring tool to the bare minimum and also minimizing the number of artifacts users of this tool are required to use.

- **Modularity**

The Distrace tool should be designed in a way that some parts of the whole tool may be substituted by user specific modules. For example, the users should be allowed to switch the default user interface to the user interface they prefer without significantly affecting the code of the tool.

The discussion of different approaches for meeting the requirements above are discussed in the Chapter 3.

1.2 Thesis Outline

The thesis starts with the background in the Chapter 2. The purpose of this chapter is to give the reader overview of relevant tools to the thesis such as list of several profiling tools, instrumentation and communication libraries. It also describes the relevant cluster monitoring tools like Google Dapper and Zipkin in more details. The following Chapter 3 contains analysis of the solution and discussion of how specific requirements of the thesis are met. It also mentions the weaknesses of the relevant cluster monitoring tools and describes how the thesis tries to overcome some of the issues. The following Chapter 4 contains design of Distrace. It starts with the Section 4.1 describing overview of the whole tool and depicting the architecture of the whole system. This section is followed by the Section 4.2 demonstrating a simple use case for this tool. Further, the Chapter 4.1 contains sections for each important part of the application such as explanation of spans in the Section 4.3, the instrumentation server in the Section 4.5, the native agent in the Section 4.4 and the user interface in the Section 4.6. The next Chapter 5 describes several interesting implementation details in more depth and is followed by the Chapter 6. The purpose of this chapter is to show a more complex example of how Distrace can be used, in this case on the H2O machine learning platform. The task is to visualize the hierarchy of internal map-reduce calls inside the H2O platform and also to see how long each map and reduce operation last. The last Chapter 7 mentions the current limitations of the tool and also compares different deployment strategies of Distrace and give the comparisons between them.

2. Background

This chapter covers technologies relevant to the thesis. It starts with an overview of similar monitoring tools for cluster-based applications and follows by short overview of tools used for debugging of large scale applications. Different approaches to applications profiling are described in the following part.

In the next several sections the technologies considered to be used in Distrace or used in thesis are introduced. The sections cover libraries for bytecode manipulation, communication, logging and also cover important relevant parts of Java libraries such as JNI and JVMTI. Docker is briefly described at the end of this chapter as it is used as the main distribution package of the whole platform.

2.1 Cluster Monitoring Tools

The most significant and relevant platforms on which this thesis is inspired are Google Dapper, Zipkin and HTrace. All three tools serve the same core purpose, which is to monitor large-scale Java-based distributed applications. Zipkin and HTrace are developed according to Google Dapper design and therefore these three platforms share a few similar concepts. The basic concept shared between the platforms is a concept called *Span*. Spans are structures used to collect the shared state between the application nodes and usually encapsulate a few calls between the neighboring nodes with well-defined start and end of the communication. They are formed in hierarchical structures called trace trees, in which the user can see how distributed calls related on each other. Spans are explained in more details later in the following sections.

The following three sections describe the basics of each of the mentioned platforms. Since Zipkin, HTrace and Google Dapper share some basic concepts, only those parts of each platform that are relevant and interesting to the thesis are described.

2.1.1 Google Dapper

Google Dapper [9] is a proprietary software used at Google. It is mainly used as a tool for monitoring large distributed systems and helps with debugging and reasoning about applications performance running on multiple host at the same time. Different parts of the monitored system does not have to be written in the same programming language. Google Dapper has three main pillars on which is built:

- **Low overhead**

Google Dapper should share the same life-cycle as the monitored application itself to capture also the intermittent events and thus low overhead is one of the main design goals of the tool.

- **Application level transparency**

The developers and users of the application should not need to know about the monitoring tool and are not supposed to change the way how they interact with the system when the monitoring tool is running. It can be

assumed from the paper that achieving application level transparency at Google was easier than it could be in more diverse environments since all the code produced at the Google use the same libraries and share similar control flow practices.

- **Scalability**

Such a system should perform well on data of significantly large scale.

Google Dapper collects the information from distributed applications as distributed traces. The origin of the distributed trace is the communication or task initiator and the trace spans across the nodes in the cluster also participating in the computation or communication.

There were two proposed approaches for obtaining the distributed traces when Google Dapper was developed: black-box and annotation-based monitoring approaches. The black-box approach assumes no additional knowledge about the application whereas the annotation-based approach can make use of additional information via annotations. Google Dapper is mainly using black-box monitoring schema since most of the control flow and RPC (Remote Procedure Call) subsystems are shared among Google, however support for custom annotations is provided via additional libraries build on top of the core system. This gives the developer of an application possibility to attach some additional information to spans, which are very application-specific.

In Google Dapper, distributed traces are represented as so called trace trees, where tree nodes are basic units of work referred to as spans. Spans are related to other spans via dependency edges. These edges represents relationship between parent span and children of this span. Usually the edges represent a set of related RPC calls or similar kind of communication. Each span can be uniquely identified using its id. In order to reconstruct the whole trace tree, the monitoring tool needs to be able to identify the span where the computation started. Spans without parent id are called root spans and serve exactly this purpose. Spans can also contain information from multiple hosts, usually from direct neighborhood of the span. Structure of a span in Google Dapper platform is described in the Figure 2.1. The figure shows a single span encapsulating the client-server communication. It can be seen that events from both nodes are recorded within a single span.

Google Dapper is able to follow distributed control paths thanks to instrumentation of a few common shared libraries among Google developers. This instrumentation is not visible to the final users of the system and therefore Dapper has high-level of application transparency. Instrumentation in Dapper is achieved by tracing three main instrumentation points.

- Dapper attaches a so called trace context as thread-local variable to the current thread when the thread handles any kind of traced control path. Trace context is a small data structure containing mainly reference to a current and parent span via their ids.
- In distributed systems the communication is often done via callbacks. Callback is a method which is usually executed when some execution in different part of the application has finished. Dapper instruments the callback mechanism so when computation is deferred, traced callbacks still carry around



Figure 2.1: Example of a span in Google Dapper, taken from Google Dapper paper [9]

the trace context of the creator and therefore also parent and current span ids.

- Most of the communication in Google is using single RPC framework with language bindings to different languages. This library is instrumented as well to achieve the desired application-level transparency.

Sampling of the captured data has also a positive effect on the low-level overhead of the whole application. As mentioned in the paper, the volume of data processed at Google is significant so only samples are taken at a time.

2.1.2 Zipkin

Zipkin¹ is an open-source distributed tracing system. It is based on Google Dapper technical paper and manages both the collection and lookup of the captured data. Zipkin uses instrumentation and annotations for capturing the data. Some information are recorded automatically, for example time when a span was created, whereas some are optional. Zipkin has also support for custom application-specific annotations.

Zipkin architecture can be seen on Figure 2.2. The instrumented application is responsible for creating valid traces. For that reason, Zipkin has set of pre-instrumented libraries ready to be used in order to work well with the whole Zipkin infrastructure. Spans are stored asynchronously in Zipkin to ensure lower overhead. Once a span is created by the application, it is sent to Zipkin collector. In General, Zipkin consists of four components:

- **Zipkin Collector**
The collector is usually a daemon thread or process which stores, validates and indexes the data for future lookups.
- **Storage**
Data in Zipkin can be stored in multiple ways which makes this is a plug-

¹More information about Zipkin tracing tool is available at <http://zipkin.io>.



Figure 2.2: Zipkin architecture, from Zipkin Architecture [1]

gable component. For example, data can be stored in Apache Cassandra², MySQL³ or can be send to Zipkin user interface right away without any intermediate storage. The last option is good for handling a small amount of data since the user interface is not supposed to handle storing data of big size.

- **Zipkin Query Service**

This component acts as a query daemon allowing the user to query various information about spans using simple JSON (Javascript Object Notation)⁴ API.

- **Web UI**

A basic, but very useful user interface, which allows the user to see whole trace trees and all spans with dependencies between them. The user interface accepts the spans in JSON format. By default, it is available on port 9411.

The Zipkin Web UI is used as a front-end for the monitoring tool developed in the scope of this thesis. More information on how the user interface is used in Distrace tool is described in more details in the Section 4.6.

²Apache Cassandra is a free and open-source distributed NoSQL database. It is designed to handle large amount of data and provide high-availability

³MySQL is an open-source relational database system.

⁴JSON is a lightweight data-interchangeable format based on object notation used in Javascript programming language

2.2 HTrace

HTrace⁵ is a tracing framework created by Cloudera used for monitoring distributed systems written in Java. It is based on Google Dapper as well and shares the same concepts such as spans and trace trees. In order to allow tracing of the application, the users need to manually attach a span identifiers to desired RPC calls (Remote Procedure Calls). These identifiers are then used to create relationships between spans collected on different nodes. HTrace stores the span and trace information in thread-local storage and the user is responsible for making sure this state is transferred to a new thread or node. HTrace has also support for custom spans annotations and thus allows the user to collect application specific information as part of the spans.

The disadvantage of this tool is the need for instrumenting the monitored application in order to allow the tracing and spans collection.

2.3 Tools for Large-Scale Debugging

Standard techniques and tools can be used for debugging distributed applications, however the main purpose of these tool is to debug single node applications and therefore when applying them on nodes in the distributed application the information about dependencies between different nodes in the cluster is not available. Many tools for large-scale debugging exist, but this section just points out basic ideas behind two different approaches - discovering scaling bugs and behavior based debugging.

2.3.1 Discovering Scaling Bugs

The scalability is one of the most important aspects of distributed systems. It is desirable to know how the platform scales when it process significantly large data and what is the expected scalability trend. It can happen that the platform can run significantly slower on big data than expected after testing on smaller data. This issue is usually called a scaling bug. Tools which can be used to help discovering scaling bugs are for example Krishna and WuKong [6]. Both of the mention tools are based on the same idea. They build a scaling trend based on data batches of smaller size and the observed scaling trend acts as a boundary. The scaling bug becomes observable when the scaling trend is violated. The first tool, Vrisha, is not able to distinguish which part of the program violated the scaling trend. This is however possible in the second tool, WuKong. In comparison to Vrisha, WuKong does not build one scaling trend of the whole application, but creates more smaller models, each per some control flow structure in desired programming language. All these smaller models represent together the whole scaling trend. When the application observes the scaling bug, WuKong is able to locate the developer in the place in the code where the trend is probably violated.

⁵The project is available at <https://github.com/cloudera/htrace>.

2.3.2 Behavior-based Analysis

The second category of tools used for debugging large scale applications are based on behavior analysis. The basic idea behind these tools is creation of classes of equivalence from different runs and processes of the application. Using this approach, the amount of data used for further inspection is reduced down. These tools are especially helpful when discovering anomalies between different observed application runs. For example, STAT - Stack Trace Analysis Tool [6], is a lightweight and scalable debugging tool used for identifying errors on massive high performance computing platforms. It gathers stack traces from all parallel executions, merges together stack traces from different processes that have the same calling sequence and based on that creates equivalence classes which make it easier for debugging highly parallel applications.

The other tool used as an example in this category is AutomaDed [6]. This tool creates several models from an application run and can compare them using clustering algorithm with (dis)-similarity metric to discover anomalous behavior. It can also point to specific code region which may be causing the anomaly.

2.4 Profiling Tools

Profiling is a form of dynamic code analysis. It may be used for example for determining how long execution of each part of the system takes compared to the time of whole application run or for example to determine which part of the application uses the most memory. Profiling tools can be divided in two categories:

- **Sampling Profilers**

Sampling profilers take statistical samples of an application at well-defined points such as method invocations. The points where the application should take samples have to be inserted at the compilation time by the compiler. For example, these profiles are good to collect information about time how long a method run, caller of the method or for example the complete stack trace. However they are not able to collect any application specific information.

- **Instrumentation Profilers**

The instrumentation profilers are based on the instrumentation of the application's source code. They record the same kind of information as the sampling profilers and usually give the developer the ability to specify extra points in the code where the application-specific data are recorded.

Sampling profilers usually have less overhead compared to instrumentation profilers, but on the other hand, instrumentation profilers allow to monitor application-specific parts of the application.

Profilers can be also looked at from different point of view and categorized based on the level on which they operate and are able to record the information.

- **System Profilers**

System profilers operate on operating system level. They can show system

code paths, but are not able to capture method calls for example in Java application.

- **Application Specific Profilers**

Generally, application specific profilers are able to collect method calls within the application. For example, JVM profilers can show Java stack traces but are not able to show the further call sequence on the operating system level.

The ideal solution for monitoring purposes of Java applications would be to have information from both kind of profilers, however combining outputs of these profiler types is not straightforward. The profilers used for collecting traces from both operating system and JVM level are usually called mixed-mode profilers. JDK8u60 comes with the solution in a form of extra JVM argument `-XX:+PreserveFramePointer` [4]. Operating system is usually using this field to point to the most recent call on the stack frame and system profilers make use of this field. In case of Java, compilers and virtual machines don't need to use this field since they are able to calculate the offset of the latest stack frame from the stack pointer. This leaves this register available for various kind of JVM optimizations. The `-XX:+PreserveFramePointer` option ensures that JVM abides the frame pointer register and will not use it as a general purpose register. Therefore, both system and JVM stack frames can appear in single call hierarchy. Using the JVM mixed-mode profilers, it is possible to collect stack traces leading to:

- **Page Faults** - page faults are useful to show which JVM code triggered main memory to grow.
- **Context Switches** - context switches are used to determine code paths that often lead to CPU switches.
- **Disk I/O Requests** - capturing Input/Output information allow to see code paths leading for example to blocking seek operation on the hard-drive.
- **TCP Events** - these traces show code paths going from high-level Java code to low-level system methods such as `connect` or `accept`. They can be used to reason about performance and good design of network communication in much more better detail.
- **CPU Cache Misses** - information about cache misses can be used to optimize Java code to make better use of the existing cache hierarchy.

All the information above can be described on special graphs called Flame graphs.

Flame Graphs

Flame Graphs are special graphs introduced by developer Brendan Gregg. Flame graphs are visualization for sampled stack traces, which allows the hot paths in the code to be identified quickly. The output of sampling or instrumentation profiler can be significantly big and therefore visualizing can help to reason about performance in more comfortable way. The example flame graph is shown on the Figure 2.3.

Flame graph is a graph where:



- Each box represents a function call in the stack.
- The **y-axis** shows stack frame depth. The top function is the function which was at the moment of capturing this flame chart on the CPU. All functions underneath of it are its ancestors.
- The **x-axis** shows the population of traces. It does not represent passage of time. The function calls are usually sorted alphabetically.
- The width of each box represents the time of how long the function spent on CPU.
- The colors are not significant, they are just used to visually separate different function calls.

Flame charts can be created in a few simple steps, but it depends on the type of profiler the user wants to use. The three steps are:

1. Capture stack traces. For this step the profiler of custom choice may be used.
2. Fold stacks. The stacks need to be prepared so Flame graphs can be created out of them. Scripts for most of the major profilers exist and may be used to prepare the folded stack trace. The scripts are available on the official page for the Flame Graphs.
3. Generate the flame graph itself again using the helper script.

The purpose of this really short section is just to introduce the idea of Flame graphs because it's one of the future plans to add support for flame graphs into the Distrace monitoring tool developed in the scope of this thesis. For more information about the flame charts please visit the Brendan Gregg's blog⁶.

⁶The blog is available at <http://www.brendangregg.com/flamegraphs.html>.

2.5 Byte Code Manipulation Libraries

Distrace highly depends on the Java bytecode instrumentation and this section gives an overview of four bytecode manipulation libraries considered to be used at the thesis: Javassist, Byte Buddy, CGLib and ASM. Since it's a core feature of the whole platform and affects both the performance and the usability of the tool, the library was thoroughly reviewed before selected. Byte Buddy library was selected and is therefore described in more detail. However the reasons for its selection can be found in the following Chapter 3.

2.5.1 ASM

ASM⁷ is a low-level high-performance Java bytecode manipulation framework. It can be used to dynamically create new classes or to redefine already existing classes. It works on the bytecode level so the user of this library is expected to understand the JVM bytecode in detail. ASM also operates on event-driven model as it makes use of Visitor design pattern⁸ to walk through complex bytecode structures. ASM defines some default visitors such as `FieldVisitor`, `MethodVisitor` or `ClassVisitor`. The ASM project can be a great fit for project requiring a full control over the bytecode creation or inspection since it's low-level nature.

2.5.2 Javassist

Javassist⁹ is a well-known bytecode manipulation library built on top of ASM. It allows the Java programs to define new classes at run-time and also to modify class files prior the JVM loads them. It works on higher level of abstraction compared to ASM so the user of this library is not required to work with the low-level bytecode. The advantage of Javassist is that the injected code does not depend on the Javassist library at all. The code to be injected to the existing bytecode is expressed as instances of Java `String` class. The disadvantage of this approach is that the code to be injected is not subject to code inspection in most of the current IDEs. The strings representing the code are compiled at run-time by special Javassist compiler. This run-time compilation works well for most of the common programming structures but for example auto-boxing and generics are not supported by the compiler [3]. Also it is important to mention that Javassist does not have support for the code injection itself. Therefore, it can be used for specifying the code which alters the original code but external tool needs to be used to inject the code itself.

2.5.3 CGLib

CGLib¹⁰ as another byte-code manipulation library built on top of ASM. The main concepts are build around `Enhancer` class, which is used to create proxies by dynamically extending classes at run-time. The proxified class is then used to

⁷The library is hosted at <http://asm.ow2.org>.

⁸The visitor pattern is a way to separate data and the operations, which can be performed on the data.

⁹Javassist library is hosted at <http://jboss-javassist.github.io/javassist/>.

¹⁰CGLib library is hosted at <https://github.com/cglib/cglib>.

intercept method calls and field access as is defined by the developer. However CGLib lacks comprehensive documentation making it harder to even understand the basics.

2.5.4 Byte Buddy

Byte Buddy¹¹ is fairly new, light-weight and high-level bytecode manipulation library. The library depends only on visitor API of the ASM library which does not further have any other dependencies. It does not require from the user to understand format of Java bytecode, but despite this, it gives the users full flexibility to redefine the bytecode according to their specific needs. Despite it's high-level approach, it still offers great performance [11] and is used at frameworks such as Mockito¹² or Hibernate¹³. Byte Buddy can be used for both code generation and transformation of existing code.

Code Generation

Code generation is done by specifying from which class a new class should be subclassing. In the most generic case, class can be created based on the `Object` class. The newly created class can introduce new methods or intercept methods from it's super class. In order to intercept existing methods (change their behavior and return value), the method to be intercepted has to be identified using instances of the `ElementMatchers` class. These matchers allow the developer to identify methods using for example their names, number of arguments, return types or associated annotations. The whole list of matchers and also examples how code can be generated is greatly described in the documentation of the Byte Buddy library.

The power behind Byte Buddy is also that it can be used to redefine classes at run-time. This is achieved by several concepts, mainly via Transformers, Interceptors and Advice API.

Code Transformation

In order to tell Byte Buddy what method or field to intercept, the place in code which triggers the interception has to be identified. First, a class containing the desired method for instrumentation needs to be located. It can be done by simply specifying the class name or using more complex structures. For example, the element matchers may be used to only consider all classes `A` extending class `B` whilst implementing interface `C` at the same time.

The next step is to define the `Transformer` class itself. Transformers are used to identify methods in the class, which should be instrumented, and they also specify the class responsible for the instrumentation itself. This class may be

¹¹Byte Buddy library is developed by Rafael Winterhalter and is freely available at <https://github.com/raphw/byte-buddy>. The page contains also a full API documentation and code examples.

¹²Mockito is a mocking framework with a clean API allowing developers to write readable tests. More information is available at <http://mockito.org>.

¹³Hibernate is an open source Java persistence framework project and provides object-relational mapping in the Java programming language. More information is available at <http://hibernate.org>.

either `Interceptor` or `Advice` and their description is given in more detail in the following section.

The methods to be instrumented can be specified in the transformer using the element matchers. In more detail, `AgentBuilder.Transformer` interface has a method `transform` which takes `DynamicType.Builder` as its argument. This builder is used to create a single transformer wrapping all the transformers for all classes in the code so the result of this builder can be thought of as a dispatcher of the instrumentation for complete application.

There are two ways how to instrument a class in Byte Buddy. Either via `Interceptors` or via `Advice API`.

Interceptors

`Interceptor` is a class defining the new or changed desired behavior for the method to be instrumented. The demonstration how Byte Buddy uses interceptors is shown on a small example. Let's assume the class `Foo` is the original unchanged class:

```
class Foo {
    String bar() {
        return "bar";
    }
}
```

Let's also assume that the `Interceptor` is of type `Qux`. The interception of the class `Foo` using the defined interceptor looks like this in schematic code:

```
class Foo {
    // Requires your interceptor class to be known
    static Qux $interceptor;
    String bar() {
        return $interceptor.intercept ();
    }
    static {
        // Requires knowing the framework
        $interceptor = ByteBuddyFramework.defineField(Foo.class);
    }
}
```

It can therefore be seen that in case of interceptors, Byte Buddy does not inline the bytecode to the `Foo` class but requires the interceptor class to be available on the machine where the instrumentation takes place. Also the interceptor field, in this case `$interceptor`, needs to be initialized before it is used for the first time. This is handled automatically by Byte Buddy framework using special helper class called `LoadedTypeInitializer`. When the framework discovers that the instrumented class is about to be used for the first time, it initialize the static field to correct interceptor.

However Byte Buddy also provides ways how to handle setting interceptors field manually. It is a more technical tasks but required for example in cases when the instrumentation is happening on different machine where Byte Buddy framework is not available.

1. In Byte Buddy, the initialization strategy can be modified accordingly to the specific needs. The **no-op** strategy can be used, which has the affect of Byte Buddy not trying to initialize the interceptor fields. In this case, the user needs to handle the initialization. **LoadedTypeInitializer** instance can be recorded right before the class is about to be instrumented and this observed initializer can be used later for the initialization. This is especially useful in case the instrumentation happens in different machine than the application is actually running. The initializer can be serialized together with the **Qux Interceptor** class to a machine with the application and a hook in the native agent can be created to call this initializer when the instrumented class is used for the first time.
2. Instead of referring to **Qux** as an instance, it can be delegated to as a static **Qux** class. In this case no initializers have to be used since no interceptor field is required and the interception is performed via static methods. However the interceptor class still needs to be known at run-time. This is demonstrated in a simple code bellow.

```
class Foo {
    String bar() {
        // no need to have interceptor field
        return Qux.intercept();
    }
}
```

3. Instead of using Interceptors API, advice API which in-lines the code to the bytecode of the class itself may be used.

Advice API

Advices are another approach how code can be instrumented in Byte Buddy. This approach is more limited compared to the interceptors, but in cases where it's possible to use it, the code is in-lined into the byte code of the original class and therefore no other dependencies are required. It is also stated in Byte Buddy documentation that performance of Advice API is better compared to the performance if interceptors. However, the instrumentation using Advice API is only allowed before or after the matched method. This is achieved using the **Advice.onMethodEnter** and **Advice.onMethodExit** annotations.

For example, the following method may be used as an advice:

```
class CustomAdvice {
    @Advice.OnMethodExit
    public static void exit(@Advice.This Callback callback) {
        TraceContext tc = TraceContext.getFromObject(callback);
        System.out.println("Method_finished")
        tc.closeCurrentSpan();
    }
}
```

The method within a class on which the advice should operate needs to be defined. This is achieved by creating an instance of **Transformer** class which specifies the method to be instrumented.

```
class CustomTransformer implements AgentBuilder.Transformer {

    @Override
    public DynamicType.Builder<?> transform(
        DynamicType.Builder<?> builder,
        TypeDescription typeDescription,
        ClassLoader classLoader,
        JavaModule module) {
        return builder.visit(Advice.to(CustomAdvice.class)
            .on(ElementMatchers.named("run")));
    }
}
```

Lastly, the class on which the transformer should operate needs to be defined. This is done on the main Byte Buddy instrumentation builder object in the following way:

```
.type(is(Task.class))
.transform(new CustomTransformer())
```

Therefore in this case, the **CustomAdvice** is applied on **Task** on the **run** method.

2.6 Communication Middleware

This thesis consist of several parts written in different languages that need to be able to communicate together. In order to arrange communication in such a heterogeneous environment, following libraries have been inspected.

2.6.1 Raw Sockets

It this case raw sockets is not a library but it is referred to as using raw sockets on their low-level API. Using raw sockets has several advantages and disadvantages. It gives the user full flexibility and the highest possible performance since there isn't any additional layer between the application data and the socket itself. However, integrating different platforms and different languages can be time consuming. Several frameworks have already been created to hide the implementation details of specific platforms so the user does not need to know about the language or underlying platforms differences.

2.6.2 ZeroMQ

ZeroMQ¹⁴ is a communication library built on top of raw sockets. The core of the library is written in C++, however binding into different languages exist. The library is able transport messages inside a single process, between different

¹⁴More information about the library is available at <http://zeromq.org>.

processes on the same node or transfer messages over the network using TCP or also using multicast communication. The library also allows the user to create typologies using one of the many supported communication patterns. For example, publisher-subscriber or request-reply patterns are supported. The library has several benefits compared to raw sockets:

- Hiding the differences between underlying operating systems.
- Message framing - delivering whole messages instead of stream of bytes.
- Automatic message queuing. The internals take care of ensuring the messages are sent and received in correct order. The user can send the messages without knowing whether there are other messages in the queue or not.
- Mappings to different languages.
- Ability to create different topologies. Example of a topology can be one-to-many communication pattern, where one socket can be connected to using multiple endpoints.
- Automatic TCP re-connection
- Support for zero-copy processing of messages.

Zero-copy in ZeroMQ

The library tries to apply concept called zero-copy whenever possible. When high-performance is expected from a system or network, copying of data is usually considered harmful and should be minimized as possible. The technique of avoiding copies of data is known as zero-copy.

Example of data copying can be transferring data from a memory to a network interface or from an user application to an underlying kernel. It can be seen that zero-copy can't be implemented at all layers. For instance, without copying the data from the kernel to the network interface, no data could be actually exchanged. However, ZeroMQ can achieve zero-copy at least on the application message level so the users can create ZeroMQ messages from their data without any copying which is a significant performance benefit.

2.6.3 Nanomsg

Nanomsg¹⁵ is a socket library shadowing the differences between the underlying operation systems. It offers several communication patterns, is implemented in C and does not have any other dependencies. Generally, it offers very similar features to ZeroMQ since it's heavily based on it.

Unlike ZeroMQ, Nanomsg matches the full POSIX¹⁶ compliance. The author of the library states that since it's implemented in C, the number of memory allocations is drastically reduced compared to C++, where, for example, C++ STL containers are used. Also compared to ZeroMQ, objects are not tightly bound to

¹⁵More information about the library is available at <http://nanomsg.org/>.

¹⁶POSIX is a family of standards used to maintain compatibility between Unix-like operating systems by specifying well-known interface for system methods [5].

particular threads. This gives the user flexibility to create their custom threading models without significant limitations. Nanomsg also implements zero-copy technique at additional layers which again leads to performance benefits compared to ZeroMQ [10].

As in ZeroMQ, Nanomsg supports the following transport mechanisms:

- **INPROC**

In-process communication is used for transporting messages within a single process, for example between different threads. In-process address is an arbitrary case-sensitive string starting with `inproc://`.

- **IPC**

Inter-process communication allows several processes to communicate on the same node. The implementation uses native IPC mechanism available on the target platform.

On Unix-like systems, IPC addresses are references to files, where both absolute and relative path can be used. The application has to have appropriate rights to read and write from the IPC file in order to allow the communication.

On Windows, the named pipes are used. The address can be arbitrary case-sensitive string containing any character except the backslash. On both mentioned platforms, the address has to start with the `ipc://` prefix.

- **TCP**

TCP is used to transport messages in a reliable manner to a single recipient in a reachable network. The address in format `tcp://interface:port` has to be used when connecting to a node. When binding a node to a specific address, the address in format `tcp://*:port` has to be used.

Nanomsg can be used via it's core C library, but several language mappings for different languages exist as well, which makes working with the library easier.

C++11 Mapping

Nanomsgxx¹⁷ a C++11 mapping for Nanomsg library. It is a small layer built on top of the core library making the API more C++11 friendly. Especially, there is no need to explicitly tell when to release resources, since it's handled automatically in the class descriptors. The `nnxx::message` abstraction over Nanomsg `nn::message` automatically manages buffers for zero-copy and also errors are reported using the exceptions which are sub-classes from `std::system_error`.

Java Mapping

Several Java bindings of Nanomsg library exists, but only jnanomsg library¹⁸ is described here. This language binding is built on top of JNA (Java Native

¹⁷More information about Nanomsgxx mapping is available at <https://github.com/achille-rousseau/nanomsgxx>.

¹⁸More information about jnanomsg mapping for Nanomsg library is available at <http://niwinz.github.io/jnanomsg/latest/>.

Access)¹⁹ library. It offers all the functionalities offered by the core library, but also introduces non-blocking sockets exposed via a callback interface.

2.7 Java Libraries

This section describes some fundamental Java-related libraries and technologies on which this thesis heavily depends. Firstly, Java Virtual Machine Tool Interface (JVMTI) is described, followed by the basic introduction to the Java Native Interface.

2.7.1 JVMTI

The JVM Tool Interface is an interface used by development and monitoring tools for communication with JVM. It allows the user to monitor and control the application running in Java virtual machine. An application communicating with the JVM using JVMTI is usually called an agent. Agents are notified via events occurring inside JVM and can react upon them. Agents run in the same process as the application itself, which reduces the delay of the communication between the application and the agent. Since JVMTI is an interface written in C, agents can be written in C or C++.

JVMTI supports two modes of how an agent can be started. It can be started either in **OnLoad** phase or in **Live** phase. In the **OnLoad** phase, the client is started together with the application and agent location can be specified using 2 arguments:

- `-agentlib:<agent-lib-name>=<options>`

In this case, the agent library name to load is specified and it is loaded using platform specific manner.

- `-agentpath:<path-to-agent>=<options>`

In this case, the path to a location of the agent library is specified and the agent is loaded from there.

In the **Live** phase, the agent is dynamically attached to a running application. This approach is more flexible since it is not required to specify the agent library to monitored application in advance. However, it brings several limitations as well.

The goal of this section is not to describe full JVMTI functionality, but just give the reader a brief introduction to the interface. For more details about JVMTI, please visit the official documentation²⁰. The following sections try to briefly describe the important parts of JVMTI relevant to the thesis.

¹⁹Java Native Library is a library allowing to use native shared libraries from Java without using JNI or native code. For more information about the library, please see <https://github.com/java-native-access/jna>.

²⁰The official JVM Tool Interface documentation describes the full API and is available at <https://docs.oracle.com/javase/7/docs/platform/jvmti/jvmti.html>.

JVMTI Agent Initialization

When an agent is started, the following JVMTI method is always called:

`Agent_OnLoad(JavaVM *jvm, char *options, void *reserved)`.

This method should contain the agent initialization specific for the application.

Usually, the agent initialization process consists of several phases:

1. Optionally, parse arguments passed to the JVMTI agent.
2. Initialize JVMTI environment in order to be able to communicate with the observed application. JVMTI does not handle threads switches automatically, so proper locking and thread management fully depends on the user code.
3. Register capabilities of the JVMTI agent. The capabilities specify what are the operations the JVMTI agent can perform. The agent can be, for example, allowed to re-transform classes or react to different class hook events.
4. Register events the agent should react to. JVMTI does not inform the agent about all events by default, these events have to be manually defined.
5. Register callbacks for the events the agent is interested in. Even though the JVMTI supports more events, the interesting events are:
 - `cbClassLoad`
 - `cbClassPrepare`
 - `cbClassFileLoadHook`
 - `callbackVMInit`
 - `callbackVMDeath`
6. Optionally, initializing phase is also good for creating locks which may be later used for synchronization between different JVMTI threads.

The user of JVMTI is also required to manually implement queuing and locking when processing multiple JVMTI events at the same time since the framework is not designed to handle these cases [7].

JVMTI basic callbacks

As mentioned above, there are several events sent from the observed Java application. When instrumenting the applications code, the following callbacks are used to capture the mentioned events.

- `cbClassLoad` - triggered when a class has been loaded by target JVM.
- `cbClassPrepare` - triggered when a class has been prepared by target JVM. All static fields, methods and implemented interfaces are available at this point, but no code has been executed at this phase.

- **cbClassFileLoadHook** - triggered when virtual machine obtains class file data, but before the class is loaded. Usually, class instrumentation is based on this hook, since the callback contains output argument for the instrumented bytecode.
- **callbackVMInit** - triggered when virtual machine is initialized.
- **callbackVMDeath** - triggered when virtual machine has been closed. This event is triggered in both planned and forcible stop.

2.7.2 JNI

Java Native Interface is a framework allowing Java code running in Java Virtual Machine to call native applications (usually written in C or C++). It also allows native applications to access and call Java methods. All JNI operations require instance of class **JNIEnv** to be available. This environment is always bound to a specific thread and manages the connection to the Java virtual machine. When calling a Java method from the native application, the correct method has to be first found. This is achieved by specifying the types and method signature of the method.

Java Types Mapping

For each Java primitive type there is a corresponding native type in JNI. Native types always start with the **j** as the prefix, for example **boolean** is a Java type whereas **jboolean** as a native type. All other JNI reference types are referred to via **jobject** class. This means that Java arrays are accessed as **jobject** as well since at this level they are referred to as Java objects. The most important question is how the types in method signatures can be specified. There is a mapping assigning each type a signature, which is used exactly for this purpose. The following table is based on the JNI documentation²¹. and describes the mapping in detail:

Type Signature	Java Type
Z	boolean
B	byte
C	char
S	short
I	int
J	long
F	float
L fully-qualified-class ;	fully-qualified-class
[type	type[]
(arg-types) ret-type	method type

For example, the method:

`xx.yy.Person foo(int n; boolean[] arr, String s);` has the following signature: `(I[ZLxx/yy/Person;])Lxx/yy/Person;`

²¹The original and more complete documentation for JNI types mapping is available at <http://docs.oracle.com/javase/7/docs/technotes/guides/jni/spec/types.html>.

Note that in JNI, the elements in fully qualified class name are separated by slashes instead of dots.

Example JNI Method Call

The method below demonstrates how JNI can be used to call a Java method `getClassLoader` from the native environment.

```
jobject getClassLoaderForClass(JNIEnv *jni, jclass clazz){
    // Get the class object's class descriptor
    // (jclass inherits from jobject)
    jclass clsClazz = jni->GetObjectClass(clazz);
    // Find the getClassLoader() method in the class object
    jmethodID methodId = jni->GetMethodID(
        clsClazz,
        "getClassLoader", // name of the Java method
        "()Ljava/lang/ClassLoader;"); // method signature
    return (jobject) jni->CallObjectMethod(clazz, methodId);
}
```

It can be seen that the reference to the method needs to be obtained at first. This reference is used later for the method invocation itself. For the performance reasons, it is a good practice to cache the references to Java methods or objects which are accessed from JNI often, since creating the reference has some initial overhead.

2.7.3 Relevant Aspects of the Java Language

This section covers selected areas of the Java programming language and run-time platform relevant to the thesis. It briefly describes the class loading process for dynamically loaded classes. This is followed by explanation of two important class loaders relevant to the thesis and lastly, `ServiceLoader` class is shortly described.

Class Loading Process

Java allows programs to load classes dynamically at run-time. This is achieved by the following process:

1. **Loading** - Load the bytecode from a class file.
2. **Linking** - Linking is the process of incorporating a new class to the run-time state of the JVM. This phase consists of 3 sub-phases:
 - (a) **Verification** - Ensure that type in the binary format is correct and respects JVM restrictions.
 - (b) **Preparation** - This phase consist of allocation the memory for fields inside the loaded type.
 - (c) **Resolution** - This phase is optional and depends on JVM implementation. Resolution is the process of transformation symbolic references in the type's constant pool into direct references. The implementation

may decide to behave in lazy way and delay resolution for the time when the type is accessed for the first time or behave in eager way and resolve all types in advance. Constant pool contains all references to variables and methods found during the compilation time at the class file.

3. **Linking Phase** - class (static) variables are initialized to initial values.

Relevant Class Loaders

There are several class loaders used natively in Java. However, this section describes only two, which are referenced later in the thesis.

- **Bootstrap class loader**

This class loader is used to load system classes. When using native agent, even classes loaded by bootstrap class loader can be instrumented and thus behavior of standard Java classes can be changed.

- **sun.reflect.DelegatingClassLoader**

This class loader is used on the Oracle JVM as the effect of a mechanism called *inflation* [2]. Usually reflective access to a method or a field is initially performed via JNI calls. When Oracle JVM determines that there is a repetition in calling the same method or the same field via JNI (reflection), it creates a synthetic class²², which is used to perform this call without the JNI. This has initial speed overhead, but at the end it speeds up the reflection calls. The classes created for this purpose are loaded and managed by exactly this class loader.

ServiceLoader Class

ServiceLoader class is used to locate and load service providers. Service provider is an implementation of a service that is usually defined as set of methods inside an abstract class or interface.

Service loader is used to load specific service providers at run-time. The group of service providers to be loaded can be specified via the service type, usually type of an interface or abstract class. The available service providers have to be defined in the META-INF folder of the application's JAR (Java Archive) distribution. For example, imagine there is a service A and two implementations, **Impl1** and **Impl2**. In that case META-INF folder should contain text file named A containing lines:

```
Impl1
Impl1
```

Service loaders can therefore be used to extend the application capabilities without changing the source code. When a new implementation of a service should be supported, it just needs to be registered inside the META-INF folder and the application will automatically use the new service provider together with the rest of the service providers defined earlier.

²²Class created dynamically at run-time

2.8 Logging Libraries

Logging can have negative effect on the performance of the application, but sometimes it's necessary to have information from various application runs to be able to locate bugs or discover wrong configuration. Since one of the thesis's requirements is low-overhead, the selection of logging library is important for the performance of the Distrace as well.

Spdlog²³ is a fast, header only logging library written in C++11 on which this project is based on. It allows both synchronous and asynchronous logging and custom message formatting.

2.9 Docker

Docker²⁴ is an open source project used to pack, ship and run any application as a lightweight container. It is used to package applications in prepared environments so the user does not need to worry about configuration and downloading the correct dependencies for the application.

Docker Compose is an extension built on top of Docker allowing the user to specify multi-container startup script. This script can define dependencies between different containers which leads to a simple and automated way how to start a group of related applications in separated environments using one single call.

²³More information about Spdlog is available at <https://github.com/gabime/spdlog>.

²⁴More information about Docker is available at <https://www.docker.com>.

3. Analysis

This chapter contains discussion of different approaches to fulfill the desired requirements. It also provides the arguments for the selection of some specific libraries later used in the Distrace tool. It starts with the section summarizing the limitations of the similar monitoring solutions, and is followed by several sections where each of them is dedicated for a single requirement. There we discuss in more detail what solution is the best for meeting the desired requirement. This chapter ends by the summary of desired and unwanted features.

3.1 Limitations of Similar Solutions

In this thesis, we try to overcome some of the limitations of the relevant monitoring tools and give users the alternative solution. One goal of the Distrace tool is to be an open-source solution, which is in contrast to the proprietary Google Dapper described in the Section 2.1.1. Google Dapper is also targeted to only specific type of applications sharing the same code structure and libraries used at Google. Therefore, it can be said that there is a certain lack of the universality at Google Dapper.

Zipkin, more described in the Section 2.1.2, is stable open-source tracing system. Zipkin creates several pre-instrumented libraries which may be used at the monitored application to communicate with the Zipkin backend. These libraries try to minimize the amount of code changes in application sources, however the user is still required to add custom annotations or change the default tracing mechanism. The similar holds for very similar HTrace tool which is described in more details in the Section 2.2. HTrace also perform the instrumentation in the same JVM (Java Virtual Machine) where the application is running which can have performance impact on the application itself.

The Distrace tool proposed by this thesis attempts to address these problems and find solutions how applications can be instrumented without the need of changing the source code and have the minimal performance impact on the running application. It however shares the same concepts as *Span* and *Distributed traces* with the already mentioned platforms.

3.2 Small Footprint and High Performance

The mentioned cluster monitoring tools, such as Google Dapper and Zipkin, can affect the application performance and memory consumption since they perform instrumentation in the same virtual machine as the monitored application. One of the thesis requirements is to have minimal footprint on the monitored application.

Since the instrumentation is the core functionality of the system, it directly affects the performance of the application. Different instrumentation methods are discussed in the following few paragraphs to give the reader insights into how each method can affect the application's performance. Two standard ways of instrumentation exist - either using the agent or Java agent. The advantages and

disadvantages of these two approaches are discussed here together with the arguments for the final solution which is actually a compromise of both techniques.

3.2.1 Java Agent

Java agents are used for instrumenting the applications on the Java-language level, where the user does not need to worry about the JVM internals. For example, Byte Buddy, Javassist or CGLib may be used for this purpose. Usually, the programmer extends and creates custom class file transformers and the agent internals take care of applying the code when required.

The advantage of this approach is obvious - the ability to write the instrumentation in the high-level language without the knowledge of the underlying bytecode. The distributions of Java Agents is also platform independent since they are packaged inside JAR (Java ARchive) files as the rest of Java classes.

The disadvantages of this approach are usually the performance and the flexibility of the agent. Objects created by Java agents are affected by garbage collection of the monitored application and thus can have negative impact on the application itself. Also the objects created from the agent are put on the application's heap and therefore consume the memory of the application. For the reasons above, it can happen that the observed information via Java Agent can be influenced by the monitoring process itself. Java agents also can not be used to respond to internal JVM events. It is important to note that Java system classes can not be instrumented when using this method.

3.2.2 Native Agent

Native agents are used for monitoring and instrumenting the applications in the low-level programming language (C, C++) using JVMTI and JNI. Native agents are written as native libraries for specific platforms and therefore the packaging is not platform independent.

The disadvantage of this method can be that the agent has to be written in non-Java language, but on the other hand, this approach gives the developer the full flexibility in the instrumentation and monitoring of the JVM state. For example, even the System classes can be instrumented using this approach and callbacks may be created to respond to several JVM internal events such as start or end of the garbage collection process, creation of a new instance of specific class or thread switches. The native agent is running as part of the Java process and therefore any resource-demanding computation can have negative performance impact on the application's performance as well as in the previous agent approach. However objects created in the native agent are not subject to garbage collection and are not created on the heap, except when they are created using JNI in the target Java application.

The significant technical disadvantage of this approach is that it does not provide any helper methods to help with the code instrumentation and generally, there is a lack of stable instrumentation libraries written in C++ or C that could be used inside the agent. The developer of the native agent has to write all the required methods for extracting the relevant parts of the bytecode and the instrumentation itself.



Figure 3.1: Sketch of the chosen approach.

3.2.3 Chosen Method

The desired solution should be able to instrument the code without affecting the performance and memory of the monitored application whilst still having access to internal JVM state and allowing the developer to use high-level programming language to write the instrumentation code. For these reasons, the compromise between the proposed approaches has been chosen, together with introduction of special process used specifically for the instrumentation. The solution can be seen on the Figure 3.1.

In more detail, the native agent is used to communicate with the application being monitored. When a class needs to be instrumented, the native agent sends the class's bytecode to a special instrumentation server JVM. This machine handles the instrumentation and sends the instrumented bytecode back to the native agent. Therefore the native agent is only used for collecting important internal JVM information, sending classes for instrumentation and receiving the instrumented classes back. The instrumentation does not happen in the same JVM as the monitored application, which allows the tool to have minimal performance impact on the application. Also, since the instrumentation is not done in the native agent, but in the instrumentation server based on Java, this machine can use any of the available bytecode manipulation tools which operates on Java-language level and therefore, there is no need to implement the bytecode manipulation library completely from scratch in the native agent. Another advantage of this solution is that even the application's Java system classes may be instrumented in Java language on the instrumentation server.

Byte Buddy library was selected for the bytecode manipulation within the instrumentation server as it allows to write the instrumentation in Java without the deep knowledge about the Java bytecode, which is necessary when using ASM library. Also, the library is supposed to have a really good performance results based on the benchmarks conducted by the author of the library [11]. Compared to Javassist, the code is not written inside Java strings, which has the effect that the instrumentation code can be validated in today's IDE (Integrated Development Environment) during compilation time and bugs in the instrumentation code can be found easier. Compared to CGLib, the API of the library is well-documented and the library is under active development. Byte Buddy is also highly configurable library, which was also the significant reason for choosing it as the tool for instrumenting classes inside different JVM then where they are actually used. Achieving the instrumentation in the secondary JVM turned out to be challenging part of the thesis and the technical aspects of the solutions are

described later in the thesis.

The disadvantage of the chosen approach is that the native agent has to send the bytecode to the instrumentation server and wait for the instrumented bytecode. However several optimizations have been implemented to minimize this delay as much as possible. More information about these optimizations can be found in the Section 4.5

3.2.4 Alternative solutions

Two alternative instrumentation solutions were analyzed but rejected at the end. The first alternative solution was to perform the instrumentation right in the agent, even for the price of affecting the application's performance by instrumenting in the same JVM. However this solution required to write the instrumentation from scratch in C++ or C language, since there are no stable libraries for this purpose. Even though this would be possible, it would take significant amount of time and also, it was not the goal of the thesis to create such a library. The performance impact on the application's was also important reason for rejecting this method.

The other alternative solution was based on the idea of running multiple Java Virtual Machines inside one native process. In particular, that would mean running the application and the instrumentation server inside the same process. This would have the same negative performance impacts as the solution above, however it would allow the developer to perform the instrumentation in Java programming language compared to C++ or C. Also all the communication between the machines would be only inside one single process compared to the chosen solution where the communication needs to be handled over the network or between different processes. However, as of JDK/JRE 1.2, creation of multiple virtual machines in a single process is not supported [8].

3.2.5 Chosen Communication Layer

The selection of the tool used for the communication between the native agent and the instrumentation server was also important decision and we choose Nanomsg for the tool purposes. Comparing to raw sockets approach, Nanomsg hides the platform specific aspects of the socket communication. It has also several performance benefits and general improvements over the well-known ZeroMQ library, such as better threading and more sophisticated implementation of the zero-copy technique. The mappings of this library into Java and C++ languages mentioned in the Section 2.6.3 were also perfect fit into the tool.

3.3 Transparency and Universality

One of the most important goals of the thesis is to achieve high level of application transparency while ensuring the tool universality. Each of these two requirements directly affects the second one and therefore we needed to find compromise between these two.

The rejected approach was to create an universal monitoring tool similar to Google Dapper, but allow to monitor mostly shared aspects of distributed ap-



Figure 3.2: .

plications. This solution would give the user great flexibility, universality and would also ensure that users don't need to extend the library. However, implementing such a solution was decided as not a feasible task. Every platform or application is different in its architecture or in the way how it communicates and therefore identifying the shared parts between all distributed applications is extremely challenging task. Also, such instrumentation tool could only instrument very basic information about Java-based programs.

The architecture of the chosen approach can be seen on the Figure 3.2. The chosen approach for the monitoring tool was to design it as an general extendable instrumentation library with two kinds of users - developer and end user. The tool should be available as a library with well-defined methods for defining the instrumentation points and specifying the custom annotations. The developer of the application is responsible for extending this library and based on it, create application-specific monitoring tool. Therefore, the developer has to have understanding of the monitored application and has to know the type of information, which are to be collected. The developer also takes care of defining where a new span starts and when existing span ends. The end-user of this final library is just responsible for starting the application with the monitoring tool attached and does not have to have understanding of the application internals.

The advantage of this approach is that only core instrumentation library needs to be created that is universal and generic to all Java-based applications. On the other hand, the application's developer is required to extend the library in order to provide the desired monitoring functionality. However, for this price, the original application can remain unchanged and from the end-user point of view, the monitoring tool is completely transparent to the monitored application.

In more detail, the core monitoring library acts as the instrumentation server mentioned in the previous section. It can be seen on the Figure 3.2 that by extending the core library, the developer creates an application-specific instrumentation server, which specifies the parts of the original application to be instrumented. The application-specific instrumentation server then communicates with the native agent, which is universal to all Java applications. It is important to mention that the instrumentation server is used only for instrumentation. The instrumented code is responsible for storing data to the user interface.

3.4 Easiness of Use

This requirement is highly connected to the previous one. The easiness of use is also important goal of the application. The usage of the core monitoring system

is separated into two user groups, developers and end users. The goal is to not require from the end user to know the internal structure of the application and the monitoring platform. This is achieved by assuming that the developer, who is responsible for extending the library, can handle this technicalities. The user is supposed to work with the user interface and to read the results of the monitoring run.

However, it is also our goal to ensure that the developer responsible for extending the core monitoring tool needs to write a few lines of code as possible and does not need to have a deep knowledge about the JVM internals or application bytecode. This is also the reason why Java language and Byte Buddy library are used for the instrumentation on the server. The Byte Buddy library provides several very concise ways how to define the application-specific instrumentation. To shield the developer from the internals, all low-level core code is hidden from the developer in the native agent, which is universal to all Java applications and the developer is not supposed to change its implementation.

3.5 Easiness of Deployment

In order to ensure the easiness of deployment, the number of artifacts used by the tool should be as low as possible. Also the application should have understandable and relatively small configuration so the users can effectively set up the application for their desired needs.

Based on the discussion in the previous sections, the tool should have only two final artifacts - the universal native agent and the core monitoring server, which is supposed to be extended by the application developer for specific application needs.

The deployment of this tool should be also simplified at certain use-cases by starting the application-specific instrumentation server automatically from the native agent. In most cases, the user should only specify path to the instrumentation library and attach the native agent to the application. Several deployment strategies exist and are discussed later in the thesis.

3.6 Modularity

It is the also goal of the tool to allow the user to replace some of the default application modules of the whole platform by specific custom modules without the need of rewriting the complete application.

The extendable modules should be the interface presenting the observed data and the collectors bringing the data from the application nodes to the user interface. Whilst default implementation are available in the tool, the users have the possibility to plug-in custom user interface or more advanced data collectors. The modules have to meet some specific criteria in order to replace the default implementations, however the core implementation is left up to the user.

The main reason for this solution is that a lot of monitoring frameworks already exist. For example, an user interface provided by some different tool may be already deployed or some platform with already defined data collection may be used. The thesis tries to support these use-cases and tries to minimize the

changes of the environment where this monitoring tool runs. This also leads to easier deployment of the platform.

3.6.1 Selection of the User Interface

Zipkin UI was selected as the default user interface for the Distrace tool. The main reasons for its selection were the simplicity of the interface and ease of use. It also fulfills our visualization requirements as it allows the user to see dependencies between spans and also whole trace tree as well. However, as mentioned above, the monitoring platform is not tightly-coupled to this user interface. It is described later in the thesis how the user can create a plug-in allowing to use custom user interface.

3.7 Summary of Features

This section just summarize the list of desired and not-desired features based on the previous background and analysis.

Desired features and properties of the thesis based on the analysis are:

- Ability to collect distributed traces and spans via the instrumentation.
- Native agent implemented in C++ to be attached to a Java application. The native agent is universal to all Java applications and is not supposed to be changed.
- Instrumentation server implemented in Java used for the instrumentation. The instrumentation server acts as the core library, which can be extended by the application's developer for the specific application. Byte Buddy library is to be used within the server for the code transformations and definitions of the instrumentation points.
- Nanomsg library is to be used as the communication layer between the native agent and the instrumentation server.
- Support for the Zipkin user interface. The user interface is used to visualize the collected spans.
- Ability to replace the default user interface with the custom user interface.

Selected features and properties which are not desired by this thesis:

- Support for creating a custom user interface.
- Support for creating an universal instrumentation tool, which would not require the developer to specify the instrumentation points and specify start and end of spans.
- Implementing bytecode instrumentation in C++.
- Instrumenting the Java bytecode in the same machine as the running application.

4. Design

This chapter describes design of the whole platform in details, however implementation specifics of some parts of the tool are described in the following chapter. This chapter starts with an high-level overview of the complete platform and interactions between the parts and follows by a simple use-case to give the reader idea how this tool is mean to be used.

Spans and their format are described next followed by design of the native agent and instrumentation server. This chapter ends by describing the used Zipkin user interface and also JSON format in which the user interface accepts the data from the instrumentation server.

4.1 Overview

Main purpose of this tool is to collect distributed traces. In order to achieve that, the thesis is based on the concept of Spans. Spans are used to denote some specific part of the communication between the communicating nodes and are the important elements for building the whole trace trees. Trace trees consists of several spans and represent the complete task or communication where a span inside the trace represents usually a few remote procedure calls between two neighboring nodes. The node initiating the trace creates so called parent span and new calls started within this span create new nested spans. Created spans can be processed using different span savers and can be send to the user interface using various data collectors.

Spans are processed using Span Savers. Span savers are used to save spans in desired format on disk or on the network for further data collection and collected data are used for Spans visualizations in the user interface. The user interface receives the spans from the span savers or data collectors and present them to the user in a form of trace trees.

Definition of when a new span is to be created and when an existing span needs to be closed is done by a developer by extending the core instrumentation server library. The created instrumentation server is then used for instrumenting the classes of the original application at which information about spans needs to be preserved. In order to obtain the class files, the native agent runs as part of the monitored application and sends the instrumentation server the desired classes. The native agent is the core part of the whole platform. It is attached to the monitored application and additionally to the instrumentation, it is used to obtain various low-level information from the application.

The tool therefore consists of three main parts:

- Native Agent
 - Is used to obtain byte-code for the instrumentation
 - Is used to actually apply the instrumented byte-code
- Instrumentation Server
 - Instruments the classes obtained from the native agent



Figure 4.1: Basic relationship between the major components. Instrumentation server communicates with the native agent mainly in order to instrument classes. The Application's communicates with the user interface by sending the spans. Spans can be send either via data collection agent or one of the default span savers explained later in this chapter.

- Is also a base library for custom application instrumentation server
- Can contain implementations of custom span savers
- User Interface

The graph 4.1 denotes the basic relationships between the major parts. Each part is explained in more detail later in this chapter.

4.2 Example Use Case

In order to demonstrate the use case for which this architecture is a good fit, a simple architecture is shown in this section. The example consists of three modules. The client module used for submitting a task, the execution module and the module used for exporting the data. These modules can be represented as a different threads in a single applications or as separated application nodes. In this example, the user always passes a task to the client module. The client module performs some pre-processing and sends the task to the computation module. This module performs the computation and sends the data to the exporter module which saves the data on disk and informs the client of the task completion. The architecture of the example can be seen on the figure 4.2.

The goal is to be able to discover how long the transfers between different modules last and how long the processing on each module last. It is also assumed that the platform does not collect this information already otherwise the cluster monitoring tool would not be required. For simplicity, let's also assume that each module performs the functionality in a single method. The following code sections give the schematic code of each method.

- **The client:**

```

public acceptTask(Task task){
    preprocessTask(task)
    ...
    sendTaskForComputation(task)
    ...
    waitForExporterToFinish()
}

```

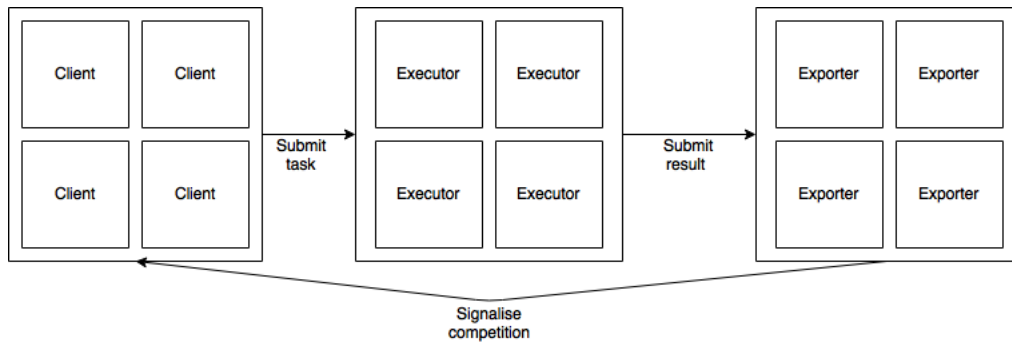


Figure 4.2: Example architecture.

```
...
}
```

- **The executor:**

```
public execute(Task task){
    TaskResult result = executeTask(task)
    ...
    sendResultForVisuzalization(result)
}
```

- **The exporter:**

```
public export(TaskResult result){
    saveToDatabase(result)
    ...
    notifyClient ()
}
```

In order to be able to collect this type of information and be able to reason about the relationship between the nodes, these methods needs to be instrumented. The instrumented code should look as in the schematic code bellow. For achieving this, the developer is supposed to extend the instrumentation server which acts as the base library and provides the developer with several helper methods used to specify the instrumentation points for their applications.

- **The instrumented client:**

```
public acceptTask(Task task){
    TraceContext tc = TraceContext.create()
    tc.attachOnObject(task)
    Span s = tc.openSpan("Main_Client_Span")
    s.addAnnotation("tskReceived", timestamp)
    ...
    preprocessTask(task)
    s.addAnnotation("tskPreprocessed", timestamp)
    ...
    sendTaskForComputation(task)
```



```

...
Result res = waitForExporterToFinish() // blocking method
...
tc.closeCurrentSpan()
}

```

The `TraceContext.create` creates a new trace and the `TraceContext.attachOnObject` attaches the trace context on the `task` object, which is passed around the network. The method `openSpan` opens a new span encapsulating the client computation. The `addAnnotation` is used to add application specific information to the current span. The `closeCurrentSpan` method is used to close the current span and export the content using the provided span saver. In the default case, the data are sent to the user interface directly.

- **The instrumented executor:**

```

public execute(Task task){
    TraceContext tc = TraceContext.getFromObject(task)
    tc.openSpan("Executor_span")
    TaskResult result = executeTask(task)
    tc.attachOnObject(result)
    ...
    sendResultForExport(result) // non-blocking method
    tc.closeCurrentSpan()
}

```

In this case, we don't create a new trace context but get the existing one from the `task` object on the input. A new nested span is created. This span has the original span created within the client application as its parent.

- **The instrumented exporter:**

```

public export(TaskResult result){
    TraceContext tc = TraceContext.getFromObject(result)
    tc.openSpan("Exporter_span")
    saveToDatabase(result)
    ...
    tc.closeCurrentSpan()
    notifyClient ()
}

```

In this case, the meaning of the methods is the same as above.

The developer should extend the base instrumentation server by steps how to instrument the classes in order to have a similar format as above. The extended instrumentation server is described on the figure 4.3.

The extended instrumentation server is then run on each node or on the network and is used to perform the instrumentation requested from the application. The native agent has to be attached to all nodes of distributed application prior it's start and the path to the extended instrumentation server JAR needs to be set as an mandatory argument. The default span sever is used in this case and



Figure 4.3: Structure of extended instrumentation server JAR.



Figure 4.4: Example trace in case of the example application.

the collected spans are send right to the Zipkin UI. The default endpoint for the user interface is used when not defined as an argument to the native agent.

A single collected trace from this application should look as shown on the figure 4.4.

Therefore, it can be seen that the only part the developer needs to work with is the extension of the instrumentation server to specify the custom instrumentation points, otherwise the rest of the work is done automatically. The end user is only responsible for starting the application with the agent attached.

4.3 Spans and Trace Trees

As mentioned briefly in the previous section, spans are used to gather the information about the distributed calls or so called, distributed stack traces. Points in the code where spans are created and closed are defined as part of the Instrumentation server but since it's the most important concept in the thesis, we explain them in the separated section.

Spans are the main concept behind capturing the distributed traces. They are special classes which instances are injected to instrumented application's classes to keep track of the communication and the state between the nodes in the distributed application. Usually, the initiator creates so called parent span and new calls started within the span create new nested spans. Collected spans can be processed using different span savers and can be sent to the user interface using various data collectors.

Spans has several mandatory and optional fields. The mandatory fields are

trace id, span id and parent span id. Trace id represents one complete distributed call among all interacting nodes on the cluster. This field is attached automatically when a new root span is created. Root span is a first span created inside a trace. The root span does not have parent id field set up and therefore the user interface back-end can distinguish between regular spans and root spans and therefore can identify the start of the whole trace. Parent id of a span is always id of span from which the span received a request to perform some task. The span and its parent span can be located on the same node or on different nodes as well. The first variant can be useful in cases where the developer requires to trace several threads as separated spans within a single node.

Span have several additional fields, which are later used in the user interface. The fields are:

- **Timestamp** - when the span started.
- **Duration** - how long the span lasted.
- **Annotations** - annotations which are used to carry additional timing information about spans. For example time when span has been received on the receiver side or the time the span has been processed at the receiver side can be set using the annotations.
- **Binary annotations** - annotations which can be used to carry around application specific details. We can use these annotations to transfer information between communication nodes inside of spans. For example, one node can store number of bytes sent during the request and the receiver can use this information to calculate overall number of bytes received from this particular node.

Each span has also an internal field called **flags**. The developer may store flags important to the instrumentation and these flags are transferred as part of the spans. Flags are not sent to the user interface and are only used for instrumentation purposes.

Each annotation, both binary and regular, has also endpoint information attached. This element consist of:

- **ip** - ip of the node on which this event was recorded
- **port** - port on which the service which recorded the span is running.
- **service name** - service name which is used to group different traces by names and can be later used to filter traces using this field in the user interface.

4.3.1 Span IDs

It is also important to mention that span id and parent span id are created randomly. This is done in order to allow parallel spans in the same control flow without overlapping as can be seen on the figure 4.5.

If the ids were not random, the parallel spans would be creating spans with ids in the same linear sequence and therefore these spans would be overlapping as can be seen on the figure 4.6

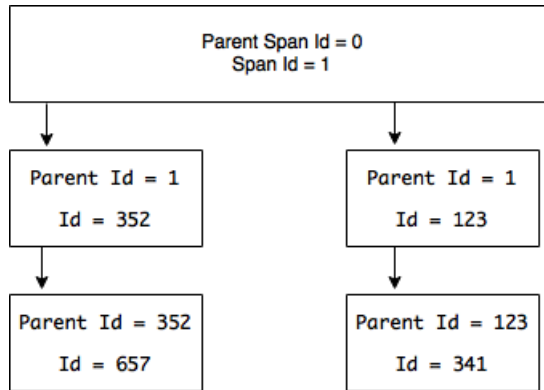


Figure 4.5: Generating span ids randomly ensures that they don't overlap when they are created in parallel.

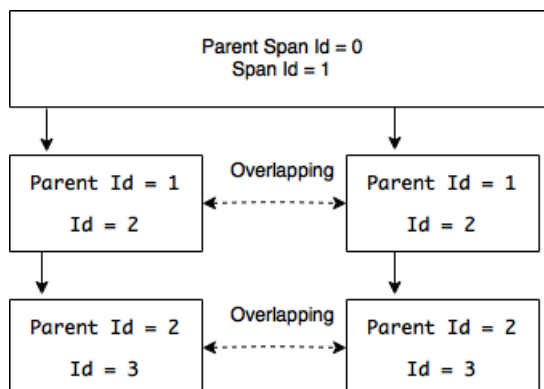


Figure 4.6: Generating span with the same linear sequence leads to span overlapping.

The following sections contain information about how spans are exported for external communication with the user interface and also how spans are created using `TraceContext` and `TraceContextManager` classes.

4.3.2 Span Savers & Data Collectors

Spans are exported from the application using the span savers. Span saver for the application may be selected via one of the native agent configuration properties. It is important to mention that all spans in the application are using the same, because span saver is static field of the `Span` class. Span savers need to be available to the monitored application and therefore are brought to the application from the instrumentation server via the native agent during the agent initialization phase.

The thesis provides two default implementation of span savers but also allows the developer to create new span savers. New span savers may be useful in cases when the developer wants to export the span data in a format used by different user interface or to use custom data collector. Data collector is a service which collects the data from the specified place and store them in central data storage available by all application's node. This thesis does not implement data collection service as many services exists for this purpose.

Data collected within spans are internally represented in JSON format understandable by the Zipkin user interface. This is also the reason why the thesis contains support for working with JSON data and it is explained in more detail later in the Implementation chapter. This format can be again changed by the custom span saver.

Span saver implementations need to extend from the abstract ancestor defining common methods for each span saver. Also, in order to be able to use the saver automatically in the code, it has to have a constructor with single `String` argument accepting saver arguments. The arguments format is defined in the case of default span savers however the developer may use any format in case of custom span savers. The common ancestor, `SpanSaver` abstract class has two abstract methods:

- **saveSpan**. This method is used for saving the span. Custom span saver implementation may save the data on local disk or send over network. The destination is not limited by the code. Internally, the **saveSpan** method is called asynchronously in a separated thread to allow asynchronous span processing, which has a performance benefit.
- **parseAndSetArgs**. The instrumentation agent accepts also special argument which contains arguments for the defined span saver. Each span saver is responsible for parsing the span saver arguments.

As mentioned above, the tool provides two default simple span savers:

- **DirectZipkinSaver** - The default span saver sends the collected span asynchronously to the user interface right away without storing the data on disk to be collected by an external data collector. In this case the functionality of the span saver and the data collector are handled by this single saver. This span saver should be used for demonstration purposes only since it could overload the user interface or network when processing and sending high



Figure 4.7: Using the Zipkin span saver to save spans directly to Zipkin user interface without the data collector agent.

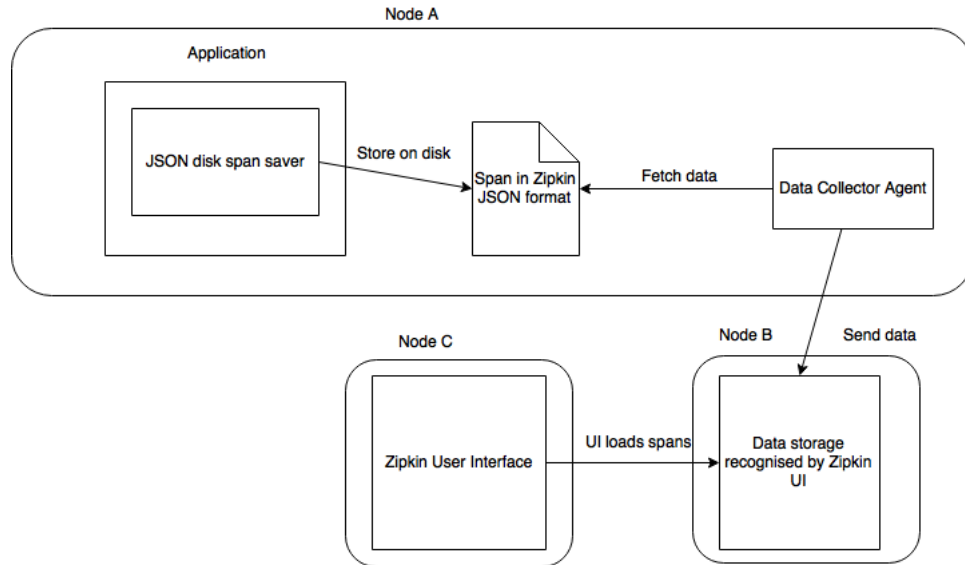


Figure 4.8: Using the JSON disk saver together with the data collection service but Zipkin user interface .

amount of collected spans in a short amount of time to the user interface as the UI is not prepared to store large amounts of data in the memory.

This saver accepts a single argument which is ip and port of the Zipkin UI service. The diagram 4.7 shows how the Zipkin span saver is used.

- **JSONDiskSaver** - The second available span saver saves the collected data asynchronously on disk in the format known to Zipkin UI for future collection to the Zipkin user interface by custom data collector. Together with some well-known data collection agent, this is a preferred way of transferring spans from the application to the Zipkin user interface in the production. This saver accepts single argument which is a directory where collected spans are saved. The diagram 4.8 shows how JSON disk span saver is used.

Additionally, the diagram 4.9 shows how a custom span saver may be used.

In order to give the developer the flexibility to add new savers without changing the internals, the span savers have to be registered in the META-INF directory of the extended instrumentation server JAR file. This ensures that the service loader can find all implementations of the `SpanSaver` abstract class. The reason



Figure 4.9: Using custom span saver together with the data collection service and custom user interface.

why the classes needs to be discovered is explained in the following Native Agent section.

To make the developer life easier, the **AutoService** library¹ is used. Instead of manually registering the implemented span savers into META-INF directory, they can be annotated in the code using the **AutoService** annotation with a single argument specifying the abstract parent, in this case **SpanSaver**. The library takes care of registering the classes automatically in the desired folder in correct format so the human error is minimized.

4.3.3 Trace Context

Trace context is a class used for storing the information about the current span and also for creating new spans and closing current spans. Trace context is always attached to a specific thread. This is done in order to allow multiple threads to have different computation state and therefore the platform is able to capture multiple distributed traces at the same time on the same node. Singleton instance of class **TraceContextManager** is used for attaching the threads to the trace contexts and vice-versa. It has a few methods allowing the developer to attach trace context to a specific thread and also to get trace context which is attached to a current thread.

Each trace represented by a trace context is uniquely identified by **Universally unique identifier (UUID)** of type one is created. This UUID version combines 48-bit MAC address of the current device with the current timestamp. This way it is ensured that two traces created at the same time on different nodes can't have the same identifier. The identifiers are created using a C++ library called sole (<https://github.com/r-lyeh/sole>) in the native agent and are made available to the Java code via a published native method.

The trace contact has method **sopenNestedSpan** and **closeCurrentSpan**. The

¹The AutoService library is available at <https://github.com/google>.

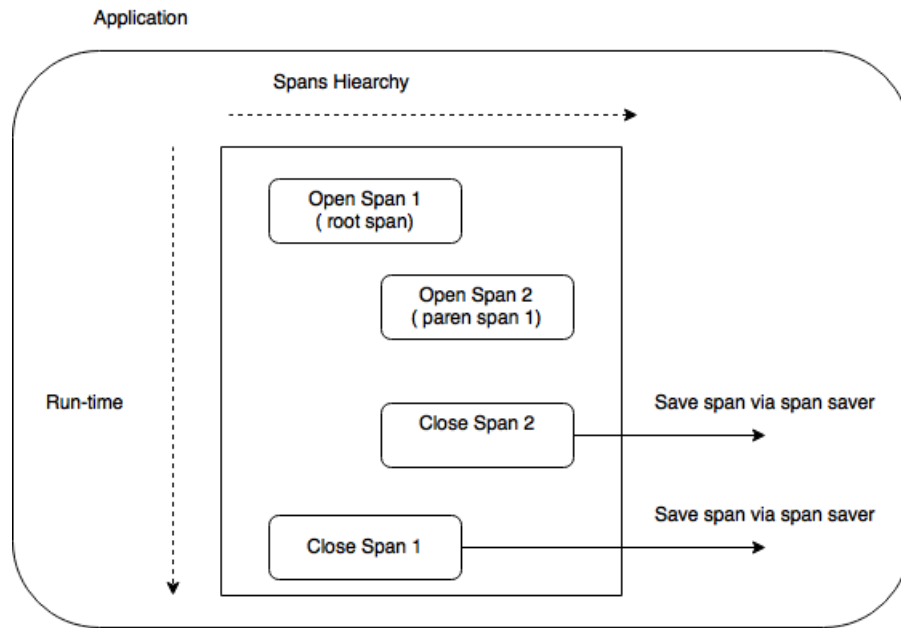


Figure 4.10: Creating and closing spans.

first method is used to create a new nested span and set the newly created span as the current one. Nested span is a span which sets its parent id to the current span. A root span is created in case when no current span exists. The second method is used to close the current span. Closing the span triggers the span saving operation on the used span saver and the parent span becomes the current span.

The diagram 4.10 shows how spans are created, closed and saved. in this chapter.

4.3.4 Transferring Span Information

In order to capture the shared state between the nodes in the application or even between the threads on the same application node, the span details and trace context have to be transferred between the threads or between the nodes. Distrace prepares several methods for attaching trace context to either current thread or an object which acts as a carrier of the span information.

Usually when transferring the trace context between the threads on the same node the copies of the trace context have to be used. When transferring the trace context between the different nodes this is not an issue.

The methods the developer can use for obtaining and attaching the trace context are:

- `static create()` - creates a new trace context.
- `static getFromObject(holder)` - gets the existing trace context from the holder object.
- `static getFromThread(thread)` - gets the existing trace context from the specified thread.

- `static getFromCurrentThread()` - gets the existing trace context from the current thread.
- `attachOnObject(holder)` - attaches the trace context to the holder object.
- `attachOnThread(thread)` - attaches the trace context to the specified thread.
- `attachOnCurrentThread()` - attaches the trace context to the current thread.
- `deepCopy` - creates a deep copy of the trace context. It is usually used in cases where child spans are processes in parallel by multiple threads. In this case, the copy of trace context with the same id is shared among all these threads, but they operate on very own object. This is done in order to be able to monitor parallel spans within a single trace without having to face race conditions on a single trace context object.

The methods above can also be chained and for example a trace context can be obtained from the holder object, deep copy created and the newly created copy attached to a new holder object.

There are also different variants how spans can be closed. Span can either closed by the same node or thread who created it. In this case attaching the trace context information to the current thread is usually sufficient. It is also however possible that span can be closed by different thread or node then where it was initially created. In this case the copy of a trace context is usually created and attached to an object which acts as an carrier of the trace information.

4.4 Native Agent

The native agent is used for accessing the internal state of the monitored application and also to instrument classes so they can carry the span and trace identifiers between the application nodes. The main agent task is to check whether a class is required to be instrumented and if yes, send the class for the instrumentation to the instrumentation server and wait for the instrumented code.

The native agent consist of several parts. The most important parts are:

- **Bytecode parsing module.**
The classes in this module are used to parse the JVM bytecode in order to discover the classes dependencies for further instrumentation. Byte code parsing is a technical task described in the following Implementation chapter.
- **InstrumentorAPI.**
The `InstrumentorApi` class provides several methods which are used to communicate with the instrumentation server JVM. All the queries to the server are done via instance of this class.
- **AgentCallbacks.**
All callbacks used in the native agent are defined in this namespace.

- **AgentArgs.**

The `AgentArgs` class contains all the logic required for argument parsing.

- **NativeMethodsHelper.**

The `NativeMethodsHelper` class is used for registering native methods defined in C++. These methods can be later used from the Java code without worrying of the low-level implementation.

- **Utilities module.**

This module contains several utility namespaces. The most important utility namespaces are `AgentUtils` and `JavaUtils`. The first one contains methods for managing the JVM TI connection and for registering the JVM TI callbacks and events. The second one is used to simplify work with Java objects in the native code via JNI.

4.4.1 Agent Initialization

The agent is initialized via the same phases as described in the JVM TI Agent Initialization section of the Background chapter. The following JVM TI events are especially important to the thesis: `VM Init`, `VM Start`, `VM Death`, `Class File load Hook`, `Class Prepare` and `Class Load`. Callbacks are registered for all the mentioned events so the native agent can react to them accordingly in the code.

As part of the initialization process, the agent is responsible for either connecting to or starting a new instrumentation server. In case the native agent was started in the shared mode of the instrumentation server, the agent tries to connect to already existing server and the server is shared between all application's nodes. In the local instrumentation mode, the instrumentation server is started as a separated process automatically and the connection is established with the server using the inter process communication. In this case, each application node has dedicated instrumentation server.

The callback registered for the `VM Init` event is responsible for loading all additional classes from the instrumentation server as part of the initialization as well. The additional classes are for example `Span`, `TraceContext` or custom implementations of `SpanSaver` abstract class. These classes are used in the instrumented code and therefore have to be available to the monitored application. The native agent is designed in a way that developers are not supposed to change the code of it. All the extension are supposed to be done within the instrumentation server. Therefore, the instrumentation server is asked at the initialization phase for the list of all additional classes and they are sent to the native agent. The agent puts all the received classes on the application's class-path so they are available to the instrumented code.

4.4.2 Instrumentation

Code for handling the instrumentation is part of the callback for the `Class File load Hook` event. The callback has the bytecode for the class being loaded as its input parameter and allow the developer to pass a new instrumented bytecode as the output parameter. The process of instrumentation is described here, however the technical details are described in the following chapter.

The process consist of several stages:

1. Enter the critical section. It can happen that the class file load hook is triggered multiple times and in order to not confuse the instrumentation server, the lock has to be acquired before the instrumentation of a class starts.
2. Firstly, the check whether the virtual machine is started is done. If the virtual machine is started and initialized, the instrumentation continues, otherwise the instrumentation for currently loaded class is skipped without asking the instrumentation server since it the class is a system class and it is not desired to instrument Java system classes at this moment.
3. Attach JNI environment to the current thread. Since the JVMTI and JNI does not have automatic thread management, it's up to the developer to take care of correct threading management.
4. Discover the class-loader used for loading the class
5. Parse the name of the class being loaded. Even though the callback provides input parameter which should contain the name of loaded class, at some circumstances it can be set to NULL even though the class name is available in the bytecode. Instead of relying on this parameter, the bytecode is parsed and the class name is found manually.
6. Decide whether the instrumentation should continue. This check is based on the used class loader and name of the class being loaded. Classes loaded by the **Bootstrap** class loader and in case of Sun JVM, from `sun.reflect.DelegatingClassL` are not supposed to be instrumented. The **Bootstrap** class loader is used to load system class and the second mentioned class loader is used to load synthetic classes and in both cases, it's not desired to instrument classes loaded by these class loaders. There are also some ignored classes for which the instrumentation is not desired. Example of these classes are the classes loaded during initialization phase from the instrumentation server and the auxiliary classes generated by the Byte Buddy framework. Auxiliary classes are small helper classes Byte Buddy is using for instance for accessing the super class of the currently instrumented class. Therefore the instrumentation continues only If the class is not ignored and not loaded by ignored class loader.
7. The instrumentation server is asked whether it already contains the loaded class or not and also if the class should be instrumented. The agent does not know which classes are to be instrumented and it therefore needs to query the server. The classes for instrumentation are marked by developer when extending the instrumentation server library using simple Byte Buddy API. If the server does not contain the class, the native agent sends the class data to the instrumentation server, parse the class file for all the dependent classes and send all dependent classes to the instrumentation. This step is repeated throughout the dependency scan recurrently until the loaded class does not have any other dependencies or until all dependencies is already

available on the server. All dependencies for the currently instrumented class have to be available on the server in order to perform the instrumentation.

8. At this stage, the class is already on the instrumentation server and all dependencies for this class as well. The native agent waits for the instrumented bytecode to be send from the server.
9. Exit the critical section.

Even though the class is fully instrumented and the instrumented bytecode is available to the agent, the process is not completely done. The instrumentation library used at the instrumentation sever (Byte Buddy) is using so called **Initializer** class to set up special interceptor field in the instrumented classes. It is a static field which references the instance of the class interceptor - class defining the instrumentation code. This field is automatically set by Byte Buddy framework in most of the cases, but since in case of this thesis the instrumentation is done on different JVM then where the code is actually running, it needs to be handled explicitly. In order to set this field by corresponding **Initializer** class, both the initializer class and interceptor class need to be available on the agent. The instrumentation server sends the initializer class together with the instance of interceptor during the instrumentation of the class and the agent registers the interceptor and initializers with the instrumented class for later use since the static interceptor field can be set up when the class is used for the first time. The initializers are loaded during **Class Prepare** event. This event is triggered when the class is prepared but no code has been execute so far.

The callback for **Prepare** event is also used to register the native methods for the class being loaded. Registering the native method to the class makes it available from Java programming language.

Several technical difficulties had to be dealt with during the development. For example, cyclic dependencies when instrumenting the class had to be properly handled. Also ensuring that the dependencies for the instrumented class are also instrumented in the correct order has been a significant challenge. The different attempts for the solution and the final solution is described in the following chapter since it's highly implementation specific.

4.4.3 Instrumentation API

The Instrumentation API is used to communicate with the instrumentation server. It provides low-level methods for sending data in form of byte arrays or strings and the corresponding methods for receiving the data. On top of these method several methods are built to make the communication easier. The most important methods are:

- **sendClassData** method sends bytecode to the instrumentation server.
- **isClassOnInstrumentor** method checks whether the bytecode for the given class is already on the instrumentation server or not.
- **instrument** method triggers the instrumentation and returns the instrumented bytecode.

- **loadInitializersFor** method is for loading the initializers for specific class.
- **loadDependencies** method is used to load all dependent classes and upload them on the instrumentation server. The dependency is uploaded only in case it's not already available on the instrumentation server.
- **shouldContinue** method checks if the class on its input is allowed to be instrumented.
- **loadPrepClasses** method loads all dependent classes in the agent initialization phase.

4.4.4 Native Agent Arguments

The native agent accepts several arguments which can be used to affect the agent behavior. In local instrumentation server mode, several arguments affect also the sever started from the agent. Available arguments are:

- **instrumentor_server_jar** - specifies the path to the instrumentation server JAR. It is a mandatory argument in case the instrumentation server is supposed to run per each node of monitored application.
- **instrumentor_server_cp** - specifies the classpath for the instrumentation server. It can be used to add application specific classes on the server classpath which has the effect that the monitored application does not have to send to the server these classes if they need to be instrumented or if some class to be instrumented depends on them.
- **instrumentor_main_class** - specifies the main entry point for the instrumentation server. It is a required argument in case of local instrumentation server mode.
- **connection_str** - specifies the type of connection between native agent and the instrumentation server. It is a mandatory argument in shared instrumentation server mode in which case the value is in format **tcp://ip:port** where ip:port is address of the instrumentation server. Otherwise, the agent and server communicates via inter-process communication and the argument can be set in format **ipc://identifier** where identifier specifies the name of pipe in case of Windows and name of the file used for IPC in case of Unix. However this value is set automatically at run-time if not explicitly specified as an argument.
- **log_dir** - specifies the log directory for the agent and when running in local server mode, specifies the log directory for the server as well.
- **log_level** - specifies the log level for the agent and when running in local server mode, specifies the log level for the server as well.
- **saver** - specifies the span saver type. The value can be either **directZipkin(ip:port)** where ip:port is address of the Zipkin user interface or **disk(destination)**, where destination represents the output directory for the captured spans. Custom span savers are supported as well. In

that case, the format of the value is a fully qualified name of the span saver with arguments in parenthesis, for example as `com.span.saver(arguments)`]

- **class_output_dir** - specifies the output directory for several helper classes received from the instrumentation server. This value is automatically set if not configured explicitly.
- **config_file** - specifies path to a configuration file containing the agent configuration. It can contain all arguments mentioned above, except the configuration file argument. Argument entries in the configuration file are in the format `arg=value` and each entry is on a new line of the configuration file.

4.5 Instrumentation Server

The instrumentation server is responsible for instrumenting the bytecode received from the native agent in separated JVM and also acts as the base library for the instrumentation for specific applications. The developer extending the instrumentation server can use prepared method to define custom instrumentation points without touching the internals of the native agent.

This section covers several design aspects of the instrumentation server, leaving the implementation details on the following sections. The core instrumentation on the server is handled by the Byte Buddy code manipulation framework. The native agent asks the server if the class being loaded is required to be instrumented. If yes, the server receives the bytecode, performs the instrumentation and sends the data back to the agent. The server does not contain any application state, in particular it does not take track about the distributed traces. The information about traces is contained in the application's instrumented classes.

The platform was designed to be configurable and deployment of instrumentation server is supported via two approaches. The instrumentation server can be either on the network available to all the application nodes and can be shared by all applications. This has the advantage of caching the instrumented classes. So when any class is instrumented for the first time, it is saved and the instrumentation is not performed for other nodes but the class is immediately sent. The disadvantage of this solution is higher latency between the agent and the instrumentation server since they are usually not on the same node. In this case the instrumentation server has to be manually started in advance. Architecture of this scenario is depicted on the diagram 4.11.

The other deployment method is that the instrumentation server runs on each application node. This has the advantage of faster communication since inter-process communication is used to communicate between monitored JVM and the instrumentation server. The disadvantage of this solution is that all classes have to be instrumented on each node since there is no communication between the instrumentation servers. In this solution, the server is started automatically during the native agent initialization. Architecture of this scenario is depicted on the diagram 4.12.

Except from the cached classes, the server does not contain any application state and it just reacts to the agent requests. It can accept four type of requests:



Figure 4.11: Architecture with shared instrumentation server. The dotted lines represents the communication between instrumentation server and the agent whilst the regular lines represents data collection from the agent to the UI



Figure 4.12: Architecture with separated instrumentation server. The meaning of the lines is the same as on the diagram above.

- Request for code instrumentation.
- Request for storing bytecode for a class on the server.
- Request for sending all helper classes needed by the agent such as the `Span` class or `TraceContext` class.
- Request to check whether the server contains specific class or not.

The server interacts in more ways with the agent, however they are just sub-parts of the communication initiated by one of these 4 request types.

The instrumentation server needs to deal with several technical problems. The main issue is that the classes which are about to be instrumented require all other dependent classes to be available. The other issue is instrumenting the classes with circular dependencies. The server also performs several optimizations to provide faster response to the agent such as caching the instrumented classes or minimizing the communication when possible. The technical aspects of these issues and the optimizations mentioned above are described in the following sections.

4.5.1 Instrumentation

The instrumentation of the class is triggered by the agent and it's done in two stages. The first stage informs the client whether the class is already on the instrumentation server or not. The second stage is the instrumentation itself. The first stage is initiated by the agent. The server performs the check for class availability in 3 phases:

1. Check whether the instrumented bytecode for this class is available.
2. If not, check whether the original bytecode for this class is available.
3. If not, check if the class can be loaded using the server's context class loader. This handles the cases where the user builds the instrumentation server together with the application classes or adds the application classes on the instrumentation server classpath for optimization reasons.

The server informs the agent if it does not have the bytecode for the class available and in that case the agent sends the class to the server. The server registers the received bytecode under the class name. The agent therefore does not have to send the class next time since it's already cached on the instrumentation server. The second stage follows the first stage immediately. If the server already contains the instrumented class in the cache, the instrumented class is sent right away without instrumenting the class again. If the cache is empty, the class is instrumented and put into the cache.

The code instrumentation is handled by `MainAgentBuilder` and `BaseAgentBuilder` classes. The instrumentation server expects instance of `MainAgentBuilder` on the input of its `start` method. This is an abstract class containing single abstract method `createAgent(BaseAgentBuilder builder, String pathToHelperClasses)`, where the builder is a wrapper around the Byte Buddy `AgentBuilder` class, which is used to define the class transformers.

The developer needs to implement this method and specify on which classes and on which methods the instrumentation should happen. Since Byte Buddy is used for writing transformers and interceptors, please read more about Byte Buddy in the Byte Buddy section. The server provides several helper methods for creating the transformers and interceptors, which are less verbose than the standard Byte Buddy approaches.

Each transformer has to have associated either an interceptor or advice defining the code to be injected to the original code. Each interceptor implementation has to implement `Interceptor` interface. This is required for the server to be able to discover all interceptors at run-time without the need for changing the internals of the server. Each implementation of the interceptor needs to register itself in the META-INF directory of the generated JAR in the same way as the span savers mentioned in the previous section. Custom service loader is then used to locate all classes implementing the `Interceptor` interface.

The advices may be used without any special annotations since Byte Buddy in-lines the code defined by the advices into the original code and therefore there is no need to transfer the Advice classes to monitored application.

Even though Byte Buddy takes care about the instrumentation, the `BaseAgentBuilder` class is internally properly configured so the instrumentation is defined exactly as desired. The class implements four Byte Buddy listeners used for informing us about the instrumentation progress and allow us to react on the process of the instrumentation. The listeners are:

- `onTransformation` listener is called immediately before the class is instrumented. Implementation of the listener in the thesis also sends the agent all auxiliary classes required by the instrumented class and the initializers used for setting the static interceptor field on the instrumented class.
- `onIgnored` listener is called when the class is not instrumented. The class is not instrumented when the user does not define any transformer for the specified class.
- `onError` listener is called when some exception occurred during the instrumentation.
- `onComplete` listener is called when instrumentation process completed. It is called after both of `onTransformation` and `onIgnored` listeners.

Byte buddy requires dependencies for the instrumented class to be available. They are needed because the instrumentation framework needs to know signature of all methods in several cases, for example when the method is overridden in the child class. The dependencies are all the classes specified in the class file such as type of the methods return value or arguments, super class or implemented interfaces. By default, Byte Buddy tries to find these dependencies using two classes - `LocationStrategy` and `PoolStrategy`. The first class is used to tell Byte Buddy where to look for the raw bytecode of dependent classes. The classes are loaded by context class loader by default, but since the classes are received over the network, custom `InstrumentorClassLoader` class loader is used to handle the class loading. It is a simple class loader which keeps the cache of the classes

received from the agent and when a request for instrumentation comes, instead of looking into the class files, it loads the data from the cache in the memory.

However, Byte Buddy internal API does not work with raw bytecode for scanning the further dependencies and obtaining the metadata for the classes. It uses classes `TypeDescription` and `PoolStrategy` for this purpose. The first class has a constructor accepting the `Class` class and created instance contains metadata for the class such as the signature of all methods and fields, list of all interfaces or for example list of constructors. The second class is used for caching the type descriptions so they are not created every time the class is accessed.

So in overall, class lookup is done in the following two steps:

1. Check whether type description for the class is available. If yes, load the type description from the cache.
2. If the type description is not available, load the class using the `InstrumentorClassLoader`, create type description for the class and put it in the cache.

4.5.2 Custom Service Loader

In order to allow the developer to extend the base instrumentation library service loaders for loading the extensions are used. The service loader is used for two types:

- Custom span savers. Each span saver inherits from the abstract class `SpanSaver`.
- Custom Interceptors. Each interceptor implements the interface `Interceptor`.

The user can create custom span savers and interceptors by either inheriting the desired class or implementing the required interface and put the name of the class inside the text file in the META-INF directory in the JAR file. The text file has to have the same name as the abstract class or the interface the implementation is for. For example, when user creates a new `Interceptor` called `x.y.InterceptorA`, the file `Interceptor` in the META-INF folder has to contain line `x.y.InterceptorA`.

Java provides service loader for this purpose. However the standard Java implementation looks up the classes defined as above and automatically creates new instances using the well-known constructors. For the thesis purposes this was unwanted as it's only required to obtain `Class` object representing the available implementation. Therefore a custom service loader was created for this purpose. This loader works in very similar way as the standard Java one, but instead of returning the instances of loaded services it just returns classes of available services.

4.5.3 JSON Generation

The data inside spans are internally stored as instances of `JSONValue` class since in order to support the communication with the default Zipkin UI they need to be exported as JSON. JSON is a lightweight format for exchanging data where the syntax is based on Javascript object notation.

The JSON handling is based on the <https://github.com/ralfstx/minimal-json> library, however custom simplified implementation was created which fits the theses requirements. Also the number of dependencies is lowered by this decision.

This JSON support is designed via several classes:

1. **JSONValue**. The abstract ancestor of all JSON types. This type defines common methods to all implementation.
2. **JSONString**. Class representing the string type.
3. **JSONNumber**. Class representing the numeric types.
4. **JSONLiteral**. Class representing the literals **null**, **true** and **false**.
5. **JSONArray**. Class representing the JSON arrays. It has support for adding new elements into the array.
6. **JSONObject**. Class representing the JSON objects. It has support for adding a new items in the object.

Each **JSONValue** can be printed as valid JSON string where the printing is driven by **JSONStringBuilder** class. This class is also responsible for escaping the characters according to JSON standards. The default printer prints the data without any formatting as one line, however **JSONPrettyStringBuilder** prints the data in more human-readable format. The second printer is usually used for debugging purposes and the first one for real usage as the size of the data is smaller in this case.

4.6 User Interface

. The user interface receives spans and presents them in a hierarchical way so the relationships between different nodes can be seen easily. The important feature of the user interface is that the data for a single span can be sent incrementally. This means that several JSONs representing the same span can be sent with different annotations and the user interface merges these spans into single one and presents all annotations under the given span. This allow the tool to send part of data from the sender side and part of data from the receiver side directly to the user interface instead of sending the data back on forth to send them as one single complete span.

The thesis is using Zipkin as default user interface. The default data format for exporting spans is designed in order to be understandable by this user interface. The user is however still able to change the data format to support custom user interface via custom span saver class. This section gives an overview of Zipkin user interface and describes the Zipkin data model

Each span in the UI is clickable and all the additional information can be seen at that level. In this thesis the stack trace are also collected at each span for monitoring purposes. Example of such information screen can be seen on the figure 4.14.

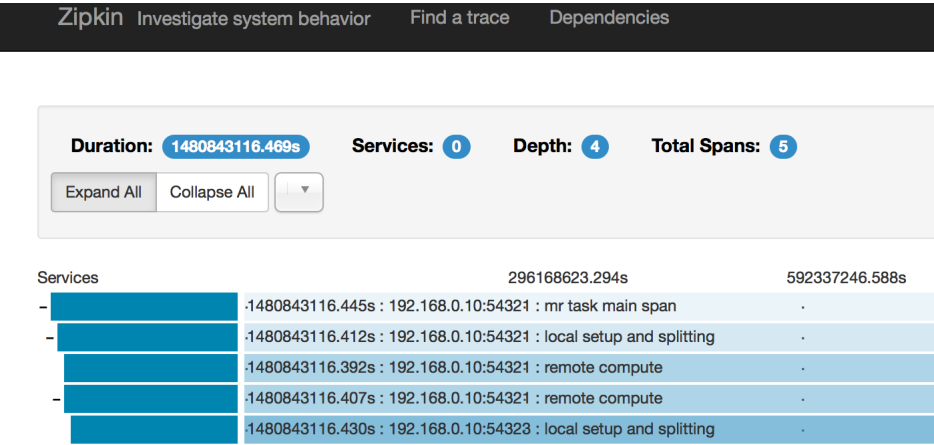


Figure 4.13: Example of Zipkin UI

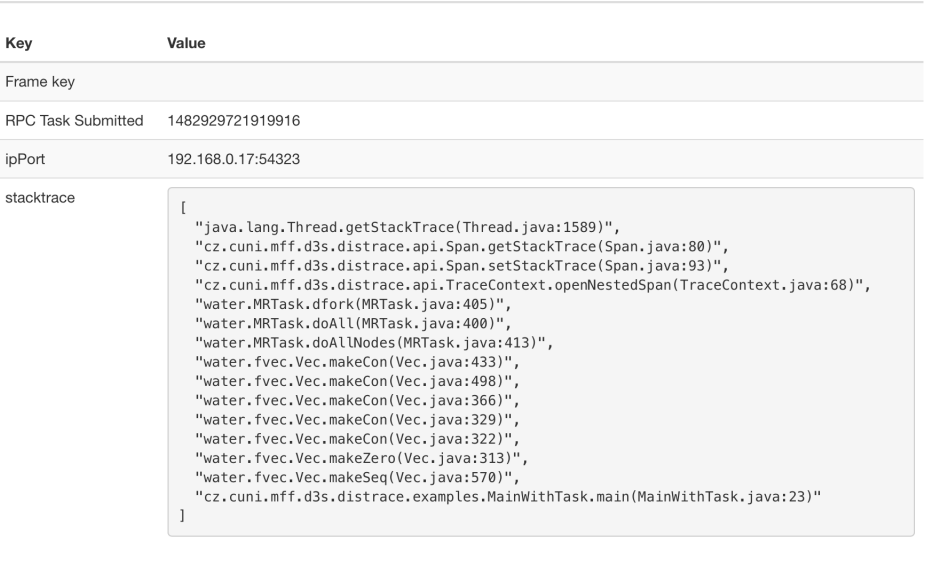


Figure 4.14: Example of the detail span information.

4.6.1 Zipkin Data Model

Zipkin requires data to be sent in JSON format. Requests to UI are sent as JSON arrays where the array elements are the spans. Zipkin understands the following fields of Span object:

- **traceId** - unique id representing the complete trace. It can be either 128 or 64 bit long.
- **name** - human readable span name
- **id** - id of this span. At the current implementation, Zipkin UI supports span ids only to be 64-bit long.
- **parentId** - parent id of the current span.
- **timestamp** - the time when the span was created.
- **duration** - the duration of the span. It is the duration between the span creation and span closing.
- **annotations** - array containing standard Zipkin annotations. These annotations can be handled by user interface in specific way since the user interface understands the meaning of the content. The documentation specifies the following annotations:
 - **cr** : timestamp of client receiving the span
 - **cs** : timestamp of client sending the span
 - **sr** : timestamp of server receiving the span
 - **ss** : timestamp of server receiving the span
 - **ca** : client address
 - **sa** : server address
- **binaryAnnotations** - array of custom annotations. For example collected stack traces are sent as a binary annotation.

Except the *annotations* and *binaryAnnotations* fields, the fields are of simple string or number type. Annotations are objects with the three fields - annotation value, annotation name and the endpoint. Endpoint is another object specifying the address and port at the code where the span or particular annotation was recorded. Endpoints can also specify service name which may be used to search for particular spans.

Full example of data sent to Zipkin can be:

```
[
{
  "traceId": "123456789abcdef",
  "name": "query",
  "id": "abcd1",
  "timestamp": 1458702548467000,
  "duration": 100743,
```

```

    "annotations": [
      {
        "timestamp": 1458702548467000,
        "value": "sr"
        "endpoint": {
          "serviceName": "example",
          "ipv4": "192.168.1.2",
          "port": 9411
        },
      }
    ],
    "binaryAnnotations": [
      {
        "key": "bytes_sent",
        "value": "1783"
        "endpoint": {
          "serviceName": "example",
          "ipv4": "192.168.1.2",
          "port": 9411
        },
      }
    ]
  }
]

```

5. Implementation Details

This chapter explains several technical implementation details. It starts with explanation of the bytecode parsing and instrumentation at the native agent part of the complete tool. The following section covers relevant parts of the instrumentation server which are the instrumentation together with creation of transformers and estimators. This chapter ends with a brief explanation of how the spans are exported to the Zipkin user interface and also, how the spans can be exported to a custom data format.

5.1 Native Agent

The native agent consists of several interesting technical parts. This section covers the instrumentation itself and also explains considered approaches during the development. The problem of instrumentation server requiring the dependencies for each instrumented class is explained together with the problem of instrumenting the classes with cyclic dependencies. The final solution is explain as well.

In the following part, the internals of how JVM bytecode is parsed is explained.

5.1.1 Instrumentation

In general, the native agent does not perform the instrumentation but gets bytecode for the the required class, sends the bytecode to the instrumentation server and applies the instrumented bytecode after receiving it from the client.

The instrumentation server required all dependencies to be available for the currently instrumented class. This means that all other classes mentioned as part of method and field signatures, super classes or interfaces has to be available on the instrumentation server. To achieve this, two solutions have been tried but only the second solution shown to be feasible.

The first and unsuccessful solution was based on the fact that several on class file load hooks may be executed at several time in different threads. When the application loads a class, the class load hook event is triggered for it and its bytecode is made available. In this method, the new class file load hook event was artificially enforced via the `RetransformClasses` method. This method accepts array of classes for which the hook should be re-thrown. In order to continue with the instrumentation of the original class, all dependent classes have to be instrumented in this solution first. This also means that in this approach, the classes with cyclic dependencies are not supported. In order to instrumented such a class, all dependencies have to be instrumented first which is also the class itself.

This solution had also different problem. Since a number of dependencies can be significant, the problem of too many threads being opened at a single time has also appeared.

The second and currently used solution is based on the fact that the Java class files may be accessed as a resource using the class loader of the class currently being load. Disadvantage of this solution that the developer may override the `getResourceAsStream` method on the custom class loader and not provide

access to the class files. This is a limitation of the thesis. However, when a such event happens, the instrumentation does not end with the exception but first, the attempt to load the class using a different class loader is done.

In this solution, the instrumentation server is first asked whether the current class should be instrumented based on the server's extensions. If the class is designed to be instrumented, its bytecode is sent to the instrumentation server (only in case if the bytecode for the class is not already available). Then, dependencies are scanned via parsing the raw JVM bytecode which is explained in detail in the following section. Dependency loading is recursively called for each new dependent class until the class does not have any other dependencies or if all the dependencies are already uploaded to the instrumentation server. Once all dependencies for a class have been sent to the server, the instrumentation is invoked and the agent waits for the new bytecode.

Also several helper classes are sent back to the native agent at the stage where the class is checked whether it should be instrumented or not. The classes are auxiliary Byte Buddy classes and also instances of `LoadedTypeInitializer` class. The initializers are sent as serialized instances and therefore their defining class has to be available in the application. This is achieved in the agent initialization phase where several required classes are sent to the application from the server. The instances are saved to a map which maps the initializer name to its serialized representation. This initializer is later used during the class preparation phase to set the static interceptor field of the instrumented class as mentioned in the previous sections.

The auxiliary classes are classes created at run-time during the instrumentation on the server and have to be available on the applications machine as well. This is achieved by loading the bytecode for the auxiliary class, saving the the class as a java class file on the disk and making it available by adding the class on the application's class path.

5.1.2 Byte Code Parsing

Byte code parsing is necessary feature of the thesis and is required in order to be able to get the list of dependent classes on a class currently being loaded. No sufficient C++ implementation has not been found and therefore a custom parsing module has been implemented. The parsing module is based on the Apache Commons BCEL Java package and the simplified version but still with the same logic has been rewritten to C++.

The main class used for parsing is `ClassParser` which contains `parse` method accepting the bytecode of a class and defines also several accessors for the parsed data such as the super class name and complete reference, list of all implemented interfaces, list of all methods or list of all defined fields and their types.

The bytecode starts with the several important parts:

- **Magic id.** Magic id is a first integer stored in each bytecode and contains always 0xCAFEBAFE number.
- **Version.** This part contains actually two shorts, where the first short represents minor Java version and the later one represents major Java version.

- **Constant pool.** Constant pool is a table representing class and interface names, field names and also other important constants. It contains mapping from id representing the type to the fully qualified type name.
- **Class Info.** Class information contains the information whether the bytecode represents a class or an interface, denotes class name and also super class name.
- **Interfaces.** This part contains the number of interfaces this class implements followed by id of type short of each interface. The interface can be looked up using the class pool.
- **Fields.** This part of the bytecode contains the number of fields this class defines together with some additional information for each defined field.
- **Methods.** This part contains the number of defined methods in the bytecode together with additional information per each method such as the number of arguments.

More information about class file structure can be found at <https://docs.oracle.com/javase/spec4.html>. Each part of the class file mentioned above is parsed separately. For accessing the raw bytecode a **ByteReader** class is used. It contains methods for reading different types of data from the bytecode array.

Parsing the magic id and both minor and major versions is straightforward as they are just numbers and can be read using the byte reader class directly. Parsing of the constant pool is more complex. For each entry in the constant pool a constant representing the entry is read. The constant can be of several types such as constant representing the Class symbol, String symbol, Method types or Field types for example. Once the Constant pool is parsed, it can be queried for the specific symbol by its id. Class name and super class name, interfaces, fields and methods are read from the constant pool by using their ids.

5.2 Span Injection on Instrumentation Server

This short section explains how additional fields such as the span information are internally attached to the instrumented classes. The trace information is attached to the instrumented class by adding a new synthetic field with name `___traceId`. This trace id represents the current trace and is used in the code to obtain reference to a current trace context and also current span. A new field is created using the Byte Buddy instrumentation builder using the `defineField` method.

5.3 Determine the Current Span Saver

This section explains how span saver type and also other arguments passed to the native agent are made accessible to the instrumentation saver. The span saver type is defined as part of the native agent and needs to be used at the instrumentation server as well. This is achieved by creating and registering the native method to a class which is defined as part of the instrumentation server. In span

saver case, the abstract class `cz.cuni.mff.d3s.distrace.storage.SpanSaver` is defined at the instrumentation server and contains also native agent for getting the span saver type. The abstract `SpanSaver` class is sent to the native agent during the initialization phase. When this class is used for the first time, the native method implementation is bound to the method `getSpanSaverType` defined in the `SpanSaver` class. Therefore, even the the classes defined as part of the instrumentation server can use method defined on the native agent with a really small overhead. Actually there may be a performance gain since these methods are written as native methods.

6. Big Example

This chapter demonstrates how Distrace can be used on bigger example and also serves as the user manual for creating custom monitoring applications. We will show how all steps from creating the application up to running the application and seeing the observed results.

This example is based on the H₂O¹ open-source fast scalable machine learning platform. This platform supports various methods for building machine learning models, methods such as deep learning, gradient boosting or random forests. The core of the tool is written in Java and clients for different languages exist as well. Internally, H₂O is using map-reduce computation paradigm² to perform various tasks across the cluster.

The goal of this example is to monitor subset of map-reduce tasks and see visualization of the computation process. This can help reasoning about performance of the platform and can discover unwanted delays in computations. This chapter first describes relevant parts of the H2O platform in more details. Then in the following sections we describe in steps how to extend the core instrumentation library for H2O purposes. Lastly, we show how this example can be started and how visual output can be interpreted. This full example is also available in the attached source code of the thesis.

6.1 H₂O In More Details

H₂O is in-memory machine learning platform. The computations are performed in cluster, where the cluster consists of several H2O nodes. All nodes in the cluster are equal and each of them can initialize the computation. Each computation is performed as a map-reduce. This section first describes the format of data used in H2O and how data are stored and then how the computations is performed in the cluster.

6.1.1 Data in H₂O

Data are stored in H₂O in so called **Frames**. Frame is an in-memory distributed data table with columns and rows. The **Frame** is designed in a way that it can handle data which are not possible to fit into a memory of a single machine. Each column is represented by the **Vec** class. This class represents vector of data that is again distributed across nodes in the cluster. Further, each vector is split into multiple **Chunks**, where chunk is part of the vector actually stored on the single node.

It is possible for one node to contain multiple chunks from a single frame and therefore, the number of chunks in a vector does not represent the number

¹More information about H₂O can be found on <https://www.h2o.ai> and <https://github.com/h2oai>

²Map-reduce is a programming model in distributed systems. The basic idea is to split tasks into smaller parts and perform the map operations. The intermediate results are then combined together using the reduce calls until the complete result has been assembled from all the sub-tasks.



Figure 6.1: Structure of H₂O frame and its distribution in the cluster.

of nodes on which the data are distributed. Also, data imported to H₂O are distributed via chunks equally among all the nodes in the cluster, but algorithm may also decide to distribute the data on just several nodes in the cluster. It is also possible to create the frame manually and specify on which nodes the chunks should be stored. Therefore, the frame may be distributed on only a portion of the cluster. Also, usually when chunks are being created on some specific node, chunks of the same size for each column are created on that node. This means that each node storing the data usually have corresponding chunks for all the columns. This can be thought of that each node storing some data has a subset of rows from the full table with all columns.

The Figure 6.1 shows the structure of frame with three columns, where each column is split into two chunks. It also shows how chunks may be distributed in the cluster of size three. It can also be seen that corresponding chunks for each vector have the same size and are stored on the same nodes.

6.1.2 Computation in H₂O

For example, when the user tells H₂O to create a deep learning model based in the data on the input, H₂O sees this as a **Job**. Jobs are used to track long-lifetime user interactions and encapsulate the whole computation from the user point of view. The job can consist of several map-reduce tasks. In H₂O, the class `MRTask` is used as the core implementation of the map-reduce tasks. The map-reduce task is always bound to some H₂O frames on which the computation needs to be performed. This class is used to encapsulate the task, partition it to a smaller tasks and run remote computations among the whole cluster. The `map` operations are called on the leaf tasks to compute the result based on the locally available data and `reduce` calls are used to reduce the result from two sub-tasks into a new task with combined result from the children.

In more detail, the class `MRTask` extends from `DTask`. This class is a general class used in H₂O to represent task remote executed. Further, `DTask` extends from `H2OCountedCompleter`. The last class is a simple wrapper around the Fork/Join execution framework³ allowing the platform to prioritize tasks. Fork/Join (F/J) framework is an implementation of Java `ExecutorService`, which helps with job

³For more information, please read Java documentation for Fork/Join framework available at <https://docs.oracle.com/javase/tutorial/essential/concurrency/forkjoin.html>.

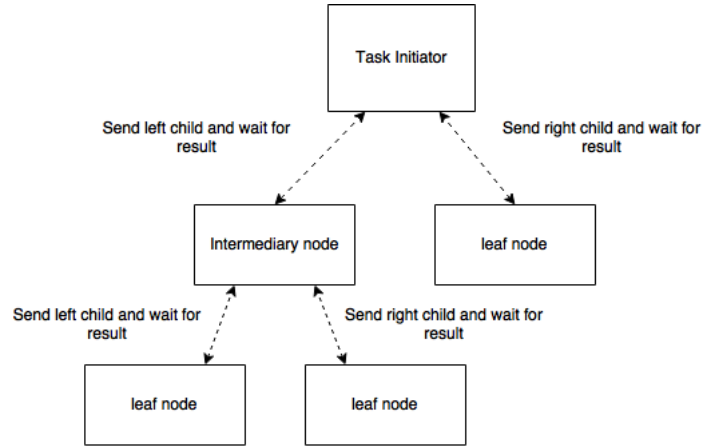


Figure 6.2: High-level overview of execution hierarchy.

parallelization on multiple processors. The Fork/Join thread framework execute tasks in separated threads and can move tasks between threads to ensure the highest possible performance. Each H₂O MRTask is executed as `ForkJoinTask` inside this execution framework. A `ForkJoinTask` is a task wrapper which can run inside a single thread. It is a light-weight wrapper and big number of tasks may be served by a smaller number of actual threads.

The way how H₂O perform computation from high-level point of view can be seen on the Figure 6.2. The task initiator receives the `MRTask` from the parent `Job` or from the user. It splits the task into two new sub-tasks and send these tasks to new nodes in the cluster. The intermediary nodes does the the splitting again and sends the task again to the another two selected nodes in the cluster. The leaf nodes does not split the task anymore.

It is important to say that the `MRTask` is always distributed to all nodes in the cluster. This figure shows how it is ensured that each node will be participating in the computation, however we still miss the computation step itself. Each node of the cluster who receives a task also submit this task for computation into the Fork/Join execution framework. This computation performs the mapping operation on all the chunks, which are available locally on the node executing this task, for the frame associated to the task. The `operation` follows the `map` operation, however we need to first ensure that the child tasks, from which we want to combine results together, are already finished. This is ensured by child tasks signaling the parent tasks when the work has been finished so parent tasks can know when they can start reducing the results.

The computation on the single node is shown in the following pseudo-code

```

MRTask task = ... // task received from the parent or from the user
MRTask left = split(task, start1, end1)
MRTask right = split(task, start2, end2)
remoteCompute(left)
remoteCompute(right)
H2O.submitTask(task) // submit this task into F/J
..
// task is taken from F/J for execution
task.map(..)
  
```

```
task.waitForComplete() // wait for the child tasks to finish  
task.reduce(left , right)
```

```
notifyComplete(task) // notify parent of completion
```

The split method accept also indices representing nodes in the cluster on which this sub-task and further sub-sub-tasks may be executed.

The last missing piece of information is what is done when the node has more chunks available for the frame associated with the task, since each `map` operation is executed only per single chunk. In case the node has multiple chunks available for the frame, the node always locally submits two new tasks into the F/J framework, each having half of the chunks. This is recursively repeated until we have tasks of size one chunk which are processed normally. These locally-split tasks signalize to their parents that they are done. This signalization goes up the tree until the original task is marked finished.

6.1.3 Methods for Instrumentation

For purpose of this example, a special `MRTask` called `SumMRTask` has been created. This task just perform distributed sum of the range of the numbers. In order to be able to visualize the computation process of this task, the following `MRTask` methods are important for the instrumentation

- `dfork` - this method is called at the initiator node and starts the computation of the whole task.
- `getResult` - this method is also called at the initiator and blocks until the distributed computation finishes.
- `setupLocal0` - this method splits the task, creates sub-tasks for child nodes and finally submits the sub-tasks to the target child nodes.
- `map` - the map operation.
- `reduce2` - the reduce operation.
- `compute2` - this method handles the computation itself by calling the `map` implementation.
- `onComplete` - this method waits for the sub-tasks to finish and then calls `reduce2` on them.

We are also interested how long the remote computation lasted on child nodes. For this reason we need to instrument the following two methods on the `RPC` class:

- `call` - this method is called when the remote computation has been submitted.
- `response` - this method is called when the remote computation has been finished and the child node is signaling that its work is done.

6.2 Building the Core Server and Native Agent

In order to be able to extend the core instrumentation library, we need to build it first. This won't be necessary once the core instrumentation server is published to some online repository of JAR packages⁴. The native agent also needs to be built on the platform where H₂O will be running.

To build the both parts, the developer needs to obtain the Distrace tool sources⁵ and call `./gradlew build`⁶ command, which builds the tool and produces artifacts for the core instrumentation server and the native agent as well. The build process requires several dependencies to be available on the machine such as Java virtual machine and Nanomsg and Nanomsgxx libraries.

In order to run this example without the need to set up these dependencies, the Docker container with this example is prepared. This container creates cluster of three H₂O nodes with the agents attached and also starts the Zipkin user interface on the default port. To run the example in Docker, Docker needs to be available, the projects source directory needs to be available as well⁷ and the following script should be called: `./docker/run-test.sh`⁸. For more information how to start the examples in docker, please read the readme file in the source code directory.

6.3 Extending the Core Instrumentation Server

In order to be able to capture the relationships between tasks and their computation, we need to instrument the method mentioned in the previous section. It is important to always find a good pair of methods which open and close a single span. In case of H₂O the following calls have been identified to form a good spans. The pairs are ordered from the top level pairs encapsulation the whole computation up until the local spans encapsulation single mapping or reduce operations.

1. `dfork` - `getResult`. This pair is used to open a main span which encapsulates the complete computation of a single `MRTask`, because `dfork` is only called on the node where the computation starts and the same holds for `getResult` on the same node after computation finished. The span is open when the first method is entered and closed when the later method is leaved.
2. `setupLocal0` - `onCompletion`. This pair forms a span representing the complete communication on a single node. It encapsulates the local work as well as waiting for the remote work to complete in order to be able to call `reduce`. This span is opened when the first method is entered and leaved when the later method is leaved. This span is well-defined since for each task the `onCompletion` method is called and represents that the work has been finished at all children tasks.

⁴For example, Maven Central Repository.

⁵Distrace is available at <https://github.com/jakubhava/Distrace>.

⁶On Windows, `./gradlew.bat build`

⁷The sources are not actually needed for running the examples in docker, but the script which is used to start the docker machines is available there as well.

⁸On Windows, `./docker/run-test.cmd`

3. **call - response**. This pair forms a span used to track only remote computation. It encapsulates the remote computation and all following sub-tasks created by the remote task. This span is opened when the first method is entered and closed when the second method is leaved.
4. **compute2 - onCompletion**. This pair forms a span used to encapsulate computation on the local node. It encapsulates all the cases of local work. In case of multiple chunks exist for the given tasks on the same node, it is recursively called until we have tasks representing a single chunk. The calls are also recursively confirmed by the **onCompletion** call. In multi-chunk case, only one child task is submitted for execution into F/J thread. The second task is executed immediately in the same thread. This way it is ensured that we can reuse the existing threads as much as possible. This span is opened when the first method is entered and closed when the first method is leaved.
5. **map - map**. This pair represents a single mapping operation. This span is entered when the **map** method is entered and closed when the **map** method is leaved. Therefore this span lasts only for duration of the **map** method call.
6. **reduce2 - reduce2**. This pair represents a single reduce operation. This span is entered when the **reduce2** method is entered and closed when the **reduce2** method is leaved. Therefore this span lasts only for duration of the **reduce 2** method call.

For all of the pairs above, the Advice API is used for instrumentation because it's sufficient to be able to capture just method entered and method exit events. In case we would like to instrument more complex cases, we could use Interceptors API which is described in the Byte Buddy documentation and briefly also in the Section 2.5.4.

Instrumenting the methods is technical tasks and the code can be seen in the example. The trace context information is always attached to the tasks transferred. The trace context is initially created during the **dfork** method call, since it's the main entry point of the computation. When instrumenting **compute2** and **setupLocal0** methods, the deep copy of trace context is attached to the transferred task as the **onComplete** method is usually called from different thread and this ensures there are no collisions and multiple objects accessing the same trace context.

Also, when instrumenting a few of the mentioned methods we need to ensure that correct pairs of spans are created. For this purpose we use support for flags on the **Span** class. Flags allow us to attach additional information to trace context which may be used at the closing side of the span. This is for example useful in cases when a method is used to close several spans at the same time and we need to properly distinguish between the spans.

Once all the advices or onterceptors are created we need to create corresponding transformers which define the association between the method in the monitored application and the method used for the instrumentation. For this purpose, Distrace provides several helper methods allowing us to create transformers in a

very concise way. For example, the transformer defining instrumentation of `call` and `response` method looks like:

```
new BaseTransformer() {
@Override
public DynamicType.Builder<?> defineTransformation(
    DynamicType.Builder<?> builder) {

    Method call = ReflectionUtils.getMethod(RPC.class, "call");

    return builder.visit (Advice.to(RPCAdvices.call.class).on(is(call))).
        visit (Advice.to(RPCAdvices.response.class).on(named("response")));
}
});
```

Once all transformers are defined, we need to associate the transformers with the classes from the application on which they should operate and also create a main entry point of the extended instrumentation server we just created. This is demonstrated on the following code snippet:

```
public static void main(String args[]) {
    new Instrumentor().start(args, new MainAgentBuilder() {
        @Override
        public AgentBuilder createAgent(
            BaseAgentBuilder builder,
            String pathToHelperClasses) {

            return builder.type(isSubTypeOf(MRTask.class))
                .transform(mrTaskTransformer)
                .type(is(RPC.class))
                .transform(rpcTransformer)
            }
        });
}
```

It can be seen in this code that we are starting the instrumentation server using the API provided by the core instrumentation server with the `MainAgentBuilder` instance. This instance is later used as a dispatcher of instrumentation of the whole application.

Later when starting the application with the attached native agent and configuring the agent, we need to explicitly specify the class, which is used as the main entry point. This is exactly the class containing the `main` method with the content above.

Now we need to build the extended instrumentation server. This server has two dependencies: the core instrumentation server and also the H₂O application sources. The first dependency is obvious as we are using the API defined in the library. The second library is required since we used several application classes when defining the instrumentation points. This may not be necessary when instrumenting different applications since we can identify the class to be instrumented for example by its name as `named("fully.qualified.class.name)` instead of `is(Example.class)`. The second option however gives us more free-

dom when defining the instrumentation points and has also a performance benefit. In this case, the application classes are already located in the instrumentor and when they are requested to be instrumented, their original bytecode don't have to be sent from the native agent.

After this step, we should have two artifacts - the native agent library built for our platform from the previous steps and also the extended instrumentation server JAR file from this step.

6.4 Configuring and Running the Application

The application can be started simply by

6.5 Interpreting the Results

7. Evaluation

7.1 Known Limitations

7.1.1 Required Java version

Distrace requires Java 8 to be available. The platform has been tested on several Java 7 implementations and several internal Java bugs occurred. These problems are already fixed in the Java 8.

7.2 Platform demonstration

7.2.1 Deployment Strategies

Instrumentor per Application Node

Instrumentor per Whole Cluster

Optimizing the Deployment

7.2.2 Basic Building Blocks

7.2.3 Basic Demonstration

7.2.4 Optimizing the Solution

7.3 Comparison to Related Work

7.4 Future plans

7.4.1 Integration with well-known data collectors

7.4.2 Add support for Flame charts

charts of performance when having shared server and local instrumentation servers

8. Conclusion

Bibliography

- [1] Architecture, 2017. URL <http://zipkin.io/pages/architecture.html>.
- [2] David Buck. Inflation system properties, feb 2014. URL <https://blogs.oracle.com/buck/inflation-system-properties>.
- [3] Shigeru Chiba. Getting started with javassist, 2015. URL <https://jboss-javassist.github.io/javassist/tutorial/tutorial3.html#boxing>.
- [4] Brendan Gregg. Java Mixed-Mode Flame Graphs at Netflix, JavaOne 2015, nov 2015. URL <http://www.brendangregg.com/blog/2015-11-06/java-mixed-mode-flame-graphs.html>.
- [5] Andrew Josey. Posix® 1003.1 frequently asked questions (faq version 1.15), jun 2015. URL http://www.opengroup.org/austin/papers/posix_faq.html.
- [6] Ignacio Laguna, Zhezhe Chen, Feng Qin, Dong H. Ahn, Bronis R. de Supinski, Todd Gamblin, Gregory L. Lee, Martin Schulz, Saurabh Bagchi, Milind Kulkarni, and Bowen Zhou. Debugging high-performance computing applications at massive scales. *Communications of the ACM*, 58(9):72–81, aug 2015. doi: 10.1145/2667219. URL <https://doi.org/10.1145%2F2667219>.
- [7] Kelly O’Hair. The jvmpi transition to jvmti, jul 2004. URL <http://www.oracle.com/technetwork/articles/java/jvmpitransition-138768.html>.
- [8] Oracle. The invocation api, 2014. URL <http://docs.oracle.com/javase/7/docs/technotes/guides/jni/spec/invocation.html>.
- [9] Benjamin H. Sigelman, Luiz André Barroso, Mike Burrows, Pat Stephenson, Manoj Plakal, Donald Beaver, Saul Jaspán, and Chandan Shanbhag. Dapper, a large-scale distributed systems tracing infrastructure. Technical report, Google, Inc., 2010. URL <https://research.google.com/archive/papers/dapper-2010-1.pdf>.
- [10] Martin Sustrik. Differences between nanomsg and zeromq. URL <http://nanomsg.org/documentation-zeromq.html>.
- [11] Rafael Winterhalter. Testing the performance of 4 java runtime code generators: cglib, javassist, jdk proxy & byte buddy, jul 2014. URL <https://zeroturnaround.com/rebellabs/testing-the-performance-of-4-java-runtime-code-generators-cglib-javassist->

List of Figures

2.1	Example of a span in Google Dapper, taken from Google Dapper paper [9]	8
2.2	Zipkin architecture, from Zipkin Architecture [1]	9
2.3	Flame Graph example	13
3.1	Sketch of the chosen approach.	29
3.2	31
4.1	Basic relationship between the major components. Instrumentation server communicates with the native again mainly in order to instrument classes. The Application's communicates with the user interface by sending the spans. Spans can be send either via data collection agent or one of the default span savers explained later in this chapter.	35
4.2	Example architecture.	36
4.3	Structure of extended instrumentation server JAR.	38
4.4	Example trace in case of the example application.	38
4.5	Generating span ids randomly ensures that they don't overlap when they are created in parallel.	40
4.6	Generating span with the same linear sequence leads to span overlapping.	40
4.7	Using the Zipkin span saver to save spans directly to Zipkin user interface without the data collector agent.	42
4.8	Using the JSON disk saver together with the data collection service but Zipkin user interface	42
4.9	Using custom span saver together with the data collection service and custom user interface.	43
4.10	Creating and closing spans.	44
4.11	Architecture with shared instrumentation server. The dotted lines represents the communication between instrumentation server and the agent whilst the regular lines represents data collection from the agent to the UI	51
4.12	Architecture with separated instrumentation server. The meaning of the lines is the same as on the diagram above.	51
4.13	Example of Zipkin UI	56
4.14	Example of the detail span information.	56
6.1	Structure of H ₂ O frame and its distribution in the cluster.	64
6.2	High-level overview of execution hierarchy.	65

List of Tables

List of Abbreviations

Attachments