



# Assessing and optimizing the performance impact of the just-in-time configuration parameters - a case study on PyPy

Yangguang Li<sup>1</sup>  · Zhen Ming (Jack) Jiang<sup>1</sup>

Published online: 08 March 2019

© Springer Science+Business Media, LLC, part of Springer Nature 2019

## Abstract

Many modern programming languages (e.g., Python, Java, and JavaScript) support just-in-time (JIT) compilation to speed up the execution of a software system. During runtime, the JIT compiler translates the frequently executed part of the system into efficient machine code, which can be executed much faster compared to the default interpreted mode. There are many JIT configuration parameters, which vary based on the programming languages and types of the jitting strategies (method vs. tracing-based). Although there are many existing works trying to improve various aspects of the jitting process, there are very few works which study the performance impact of the JIT configuration settings. In this paper, we performed an empirical study on the performance impact of the JIT configuration settings of PyPy. PyPy is a popular implementation of the Python programming language. Due to PyPy's efficient JIT compiler, running Python programs under PyPy is usually much faster than other alternative implementations of Python (e.g., cPython, Jython, and IronPython). To motivate the need for tuning PyPy's JIT configuration settings, we first performed an exploratory study on two microbenchmark suites. Our findings show that systems executed under PyPy's default JIT configuration setting may not yield the best performance. Optimal JIT configuration settings vary from systems to systems. Larger portions of the code being jitted do not necessarily lead to better performance. To cope with these findings, we developed an automated approach, ESM-MOGA, to tuning the JIT configuration settings. ESM-MOGA, which stands for effect-size measure-based multi-objective genetic algorithm, automatically explores the PyPy's JIT configuration settings for optimal solutions. Case studies on three open source systems show that systems running under the resulting configuration settings significantly out-perform (5% - 60% improvement in average peak performance) the default configuration settings.

**Keywords** Just-in-time compilation · Performance optimization · Software configuration · Performance testing · Performance analysis

---

Communicated by: Sven Apel

---

✉ Yangguang Li  
yangguang@cse.yorku.ca

Extended author information available on the last page of the article.

# 1 Introduction

Software performance is one of the crucial factors related to the success and the sustainability of large scale software systems, which serve hundreds or even millions of customers' requests every day. Failure to provide satisfactory performance would result in customers' abandonment and loss of revenue. For example, Amazon reported that one second delay in loading their webpages could result in \$1.6 billion loss in their sales revenue annually (Eaton 2017). BBC has also recently found that 10% of the users will leave their website even if there is merely one additional second of performance delay (Clark 2017). Various strategies (e.g., asynchronous requests (Insights 2017), data compression (Grigorik 2017), just-in-time (JIT) compilation (Oracle Java 8 Advanced JIT Compiler Options 2017), load balancing (What is Load Balancing? 2017), and result caching (Candan et al. 2001)) have been developed to further enhance the performance of these systems.

In general, there are three types of system executions depending on the programming languages: (1) executing natively on top of the operating systems (e.g., C and C++), (2) executing the source code by the interpreters (e.g., Python, PHP, and JavaScript), and (3) executing compiled intermediate artifacts on the virtual machines (e.g., Java and C#). Compared to the native execution mode, systems executed under the interpreted mode (a.k.a., by interpreters or virtual machines) are generally slower due to their additional layers. To cope with this challenge, the JIT compilation is introduced so that frequently executed code snippets are compiled into binaries, which can be executed natively.

Existing works on the JIT compilation focus on the jitting strategies (e.g., method (Cramer et al. 1997) vs. trace-level based code jitting (Bolz et al. 2009)), speeding up the process of the JIT compilations (Jantz and Kulkarni 2013; Gal et al. 2009; Lion et al. 2016), optimizing the performance of the underlying virtual machines (Wimmer and Brunthaler 2013; Würthinger et al. 2017, 2013), and detecting JIT unfriendly code (Gong et al. 2015). Unfortunately, there are very few existing studies which investigate the impact of the JIT configurations on the system performance. Software configuration is one of the main sources of software errors (Xu et al. 2016). The configuration settings of a software system can significantly impact its performance. Most of the existing configuration tuning and debugging studies are focused on the configurations of the studied systems (Duan et al. 2009; Yilmaz et al. 2007; Sopitkamol and Menascé 2005; Jamshidi et al. 2017a). For tuning the configuration settings of interpreters or virtual machines, the focus is mainly on optimizing the performance of the garbage collectors (Singer et al. 2007; Lengauer and Mössenböck 2014; Brecht et al. 2006), except the work by Hoste et al. (2010). In Hoste et al. (2010), Hoste et al. provided an automated approach to tuning the JIT compiler for Java, which is a method-based JIT. Hence, in this paper, we seek to investigate the impact of the tracing-based JIT configurations on the system performance by using PyPy as our case study subject.

Python is nowadays one of the most popular programming languages (Gewirtz 2017). Python has been used extensively to develop real-world business systems, including many large scale and mission-critical systems inside companies like Facebook (Komorn 2016), Google (What is python used for at Google? 2017), and PayPal (Hashemi 2014). Among various implementations of the Python programming language (e.g., CPython, IronPython, Jython, and PyPy), PyPy is generally the fastest (PyPy speed center 2017) mainly due to PyPy's efficient tracing-based JIT compiler (Bolz et al. 2009). Hence, in this paper, we focus on assessing and optimizing the performance impact of PyPy's JIT configuration settings. The contributions of this paper are:

1. This is the first empirical study on assessing and optimizing the impact of the tracing-based JIT configuration settings on system performance.
2. Our experiments are carried out on both the synthetic benchmarks as well as real systems. The empirical findings in this paper can be useful for both software engineering and programming language researchers as well as practitioners.
3. Compared to Hoste et al. (2010) which also used a search-based approach to automatically tuning the JIT configuration settings, many of the details (e.g., the initial setup, the configuration settings, and the evaluation details) are not clear. It is not easy to reapply the approach to other Java-based systems or other JIT compilers. In this paper, we have detailed our search-based configuration tuning approach, (ESM-MOGA) to ease replication.
4. To enable replication and further research on this topic, we have provided a Replication package (2018) which includes the implementation for our configuration tuning framework, PyPyJITTuner, as well as the experimental data.

## 1.1 Paper Organization

The rest of the paper is organized as follows: Section 2 provides some background information regarding the JIT compilation process and explains PyPy's JIT configuration parameters. Section 3 presents an exploratory study on the performance impact of PyPy's JIT configuration settings on two microbenchmark suites. Section 4 describes our automated approach to tuning PyPy's JIT configuration settings. Section 5 demonstrates the effectiveness of our automated approach by applying it to three open source systems. Section 6 provides some discussions based on the results of the case study and presents some future work. Section 7 presents the threats to validity and Section 8 explains the related research works. Section 9 concludes this paper and presents some future work.

## 2 Background

In this section, we will first give an overview of the JIT compilation process in Section 2.1. Then we will explain PyPy's JIT configuration setting in Section 2.2.

### 2.1 An Overview of the JIT Compilation Process

JIT compilers are introduced for systems executed by interpreters (e.g., Python, PHP, and Ruby) or virtual machines (e.g., Java and C#) to further speed up the system performance during runtime. By default, there are no JIT compilations upon the initial system startup and these systems are executed under the interpreted mode by their interpreters or virtual machines. Hence, their performance is usually slower than natively executed systems (e.g., systems programmed in C or C++), whose binaries are executed directly on top of the operating systems. To cope with this limitation, the JIT compiler is introduced so that, during runtime, various parts of the systems are converted into machine executable code (a.k.a., code jitting). However, the code jitting process is usually slow, as it takes time to load and compile the corresponding code snippets. Hence, only the commonly used (a.k.a., "hot") code snippets are usually jitted (Bolz et al. 2009; Oaks 2014). For such systems, there is usually a warmup period after the initial system startup before these systems reach the peak performance (Barrett et al. 2017). During the warmup period, the frequently executed code will be profiled to locate the "hot" spots and various code snippets are being jitted (Lion

et al. 2016). In general, there are two approaches for code jitting depending on their granularity:

- **Method-based JIT Compiling:** if one method has been used many times (a.k.a., “hot method”), the method-based JIT will compile this entire method into the binary executable format. Hotspot (Oracle’s implementation of the Java Virtual Machine) and Chakra (Microsoft’s JavaScript engine) use the method-based JIT Compiling.
- **Trace-based JIT Compiling** is more fine-grained, in which only the commonly executed code path (a.k.a., “hot path”) inside a method is compiled into the binary executable format. PyPy and TracingMonkey (Mozilla’s JavaScript engine) use the trace-based JIT Compiling.

For the method-based JIT compiler, there will be a threshold value (e.g., 1500 as the default value for the configuration parameter `-XX:CompileThreshold` in Oracle’s HotSpot JVM), which defines the number of invocations for a particular method before this method is considered to be hot. As soon as a method has been called 1500 times, the whole method will be compiled into the binary executable format.

For the trace-based JIT compiler, the process is a bit more complicated. We will use the sample code snippet shown in Fig. 1 to explain. There can be various configuration parameters which define a particular code path to be hot. For example, in PyPy, there is a configuration parameter, called *threshold*, which defines the number of times a loop has to be run before it can be considered hot. During the system execution, the PyPy JIT compiler counts the number of iterations for each loop and all code paths in the loop will be potential candidates for code jitting. For example, the code lines marked with star (\*) in Fig. 1 are the resulting jitted lines, if the method *test* is executed in PyPy under the default configuration setting. After the loop reaches 1039 (the default value for *threshold*) iterations, the JIT compiler starts to trace the code path in the next iteration. And the code path which contains the *if* branch, and the second *elif* branch will be recorded and compiled into efficient machine code. The first *elif* and the *else* branches are not jitted, as they are not in the code path of the traced iteration.

```

6 def test():
7     even = 0
8     oddLarge = 0
9     oddSmall = 0
10    c = 0
(*) 11    for i in range(1000000):
(*) 12        if i%2 == 1:
(*) 13            even += 1
(*) 14            elif i * i <= 100:
15                oddSmall += 1
(*) 16            elif i * i > 100 and i < 1000000:
(*) 17                oddLarge += 1
18            else:
19                c += 3
20    return c

```

Fig. 1 A sample PyPy code snippet with the jitted code marked as “(\*)”

## 2.2 PyPy's JIT Configuration

The types and the values of the JIT configuration parameters vary depending on the programming languages and the compilers. For example, there are more configuration parameters for PyPy's JIT compiler than Java's JIT compiler. Even within the same programming language, different language implementations may use different configuration parameters. For example, in Java, the configuration parameter which indicates the threshold value for the number of invocations for a method before code jitting is called *-XX:CompileThreshold* in Oracle's HotSpot JVM (Oracle Java 8 Advanced JIT Compiler Options 2017), and *-Xjit:count* for IBM's JVM (IBM Java 8 JIT and AOT command-line options 2017). In this paper, we have selected PyPy's JIT configuration parameters as our case study subject, due to the popularity of the Python programming language (Gewirtz 2017) and the fast execution under PyPy with its efficient JIT compiler (Bolz et al. 2009; PyPy speed center 2017). The list of JIT configuration parameters can be obtained through running the `ppypy --jit help` command. For PyPy version 5.7.1, which is the PyPy version used in this paper, there are 19 of them.

## 3 Exploratory Study

To motivate the importance of this work, we have conducted an exploratory study on the performance impact of PyPy's JIT configuration settings. We seek to answer the following three research questions:

- **RQ1:** *How different is the system performance before and after its code has been jitted?*  
When the system initially starts up, all of its code is executed under the interpreted mode. The code jitting process will not start, until certain regions of code have been repeatedly executed many times. In this RQ, we want to quantify the performance differences between the warmup and the warmed up phases.
- **RQ2:** *What is the performance impact by varying JIT configurations?*  
The system after the warmup phase would achieve its peak performance. But would the peak performance be different among different JIT configurations settings (e.g., the default config, random configurations, or disabling JIT)? In this RQ, we seek to find the performance impact of different JIT configurations.
- **RQ3:** *Do systems containing more jitted lines yield better performance?*  
Different JIT configuration settings would result in different amount of source code been jitted. Intuitively, a higher portion of the jitted code could lead to more code being executed natively, and hence result in better performance. However, the code jitting process is very resource heavy, which involves profiling system executions and compiling the hot code path into the binary executable format. In addition, the systems may need to constantly switch between the two running modes (the interpreted vs. the native execution mode). The goal of this RQ is to examine whether there is any relation between the portions of the jitted code and the system performance.

The remaining three subsections in this section will address the above three research questions. For each research question, we will first explain the experimentation process. Then we will describe the data analysis techniques, present the result findings, and discuss their implications.

### 3.1 (RQ1) How Different is the System Performance Before and After its Code has been Jitted?

During the benchmarking and the performance testing processes, it is considered as a common practice to wait for a period of time (a.k.a., the warmup phase) for the system to stabilize (Bondi 2007; Java Microbenchmark Harness (JMH) 2017), before starting the actual benchmarks or the performance tests. During the warmup phase, various regions of the code are getting jitted and the system caches are slowly being filled up. Hence, the performance of the warmup phase is generally considered as suboptimal and is discarded during the subsequent performance analysis. In this RQ, we want to quantitatively compare the system performance during and after the warmup phase.

#### 3.1.1 Experiment

To tackle this research question, we selected the following two microbenchmark suites, which assess the performance of different software systems:

- *The PyPy benchmark suite* is run daily on PyPy’s nightly builds and is mainly used to compare the performance of various Python implementations (PyPy vs. cPython). The benchmark suite consists of about 60 small Python programs, which perform various computation tasks like the n-queens solver, HTML table building, etc. For each run of the benchmark, the same benchmark programs will be run under PyPy and cPython (the default Python implementation). The performance results are uploaded and visualized in the PyPy’s Speed Center (PyPy speed center 2017). In this exploratory study, we randomly selected seven benchmark programs as shown in Table 1 for experimentation. We further instrumented these benchmark programs to gather additional performance information (e.g., individual request response time).
- *The TechEmpower Web Framework Benchmark suite*<sup>1</sup> (TechEmpower Web Framework Benchmarks 2017), whose main objective is to evaluate among various web frameworks, consists of more complicated web application-related tasks like JSON serializations, database accesses, and server-side template compositions. Different from the PyPy benchmark suite, whose programs are usually short-lived and computation intensive, the TechEmpower benchmark suite executes on long running web application servers built with various frameworks. For example, the benchmark includes Java-based web frameworks (e.g., Jetty), as well as Python-based web frameworks (e.g., Tornado and Flask). In this paper, we only focus on the Django web application frameworks.

We ran the two microbenchmark suites under the default PyPy configuration setting. To avoid measurement bias and errors (Georges et al. 2007), for the PyPy benchmark, in which the studied programs are short-lived and computational intensive, we repeated the benchmark for 30 times. For the TechEmpower benchmark, which examines the performance of processing web requests for long-running servers, we set the duration for each benchmark task to be two hours. During the benchmarking process, resource utilizations (e.g., CPU, memory, and disk) for the servers were monitored and recorded using pidstat (Performance monitoring tools for Linux 2015). We also added additional instrumentation using the JIT logging function from PyPy’s *jitlog* module to record the JIT logs to the disk. The recorded JIT logs can be further parsed with VMPprof (vmprof - a statistical program profiler 2017)

<sup>1</sup>To ease explanation, we will call this the TechEmpower benchmark in the rest of this paper.

**Table 1** PyPy benchmark programs description

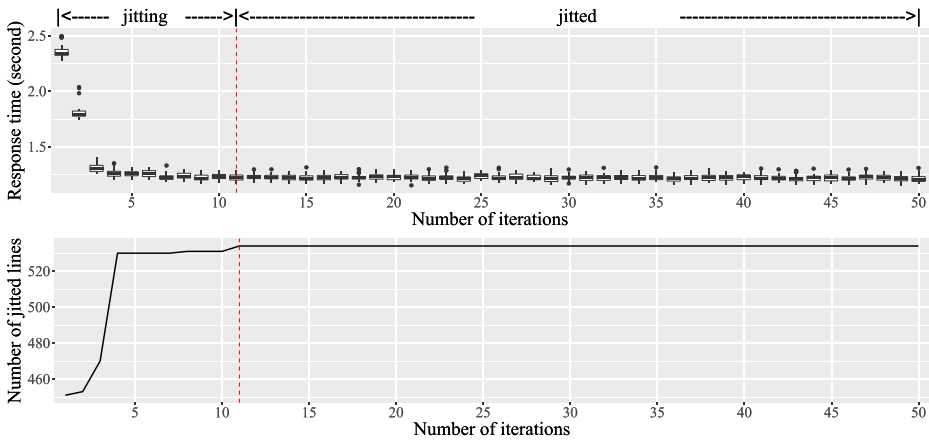
Program	Description
ai	Test the performance of simple AI solvers.
bm_mako	Benchmark for test the performance of Mako templates engine.
chaos	Test the performance of the Chaos benchmark. Create chaosgame like fractals.
django	This will have Django generate a 100x100 table as many times as you specify.
rietveld	This will have Django render templates from Rietveld with canned data as many times as you specify.
html5lib	Test the performance of the html5lib parser.
pidigits	Test the pidigit calculation performance.

to obtain the exact lines of code that were jitted. However, the recorded JIT logs do not have timestamps to mark when a code snippet is jitted. To estimate the exact timing when individual code snippets are jitted, we decided to periodically take snapshots of the JIT log files. We took snapshots of the JIT log files after each iteration of the PyPy benchmark and every minute for the TechEmpower benchmark. These snapshots would help us gain insights on the time and the location of the jitted code regions. Finally, we also archived the benchmarking logs, so that we can extract the response time for each individual request.

### 3.1.2 Data Analysis

We parsed the JIT logs using VMProf to obtain the regions of the jitted code during each snapshot period. We also processed the benchmarking logs to extract the response time for each iteration of the PyPy benchmark and the response time for individual requests in the TechEmpower benchmark.

To understand the performance of the systems during and after the warmup phase, we need to determine the duration of the warmup phase. For the PyPy benchmark, we kept track of the amount of the jitted code during each iteration. We considered the warmup phase to be completed, when the amount of the jitted code remains stable during the remaining of the benchmarking run. Figure 2 shows the result for the *html5lib* program from the PyPy benchmark. The topper subgraph of Fig. 2 shows the how the response time evolve over different number of iterations. Since each program within the PyPy benchmark is repeatedly executed 30 times, we aggregated the response time for that iteration across the 30 runs using boxplots. For example, the first boxplot contains all the response time values for the first iteration during the 30 runs. The bottom subgraph of Fig. 2 shows the evolution of number of jitted lines across different iterations. The red dotted lines in both subgraphs indicates the iteration when the number of the jitted lines becomes stable. Hence, we considered the first 11 iterations as the warmup phase (“jitting”) and the remaining iterations (a.k.a., the 12th to the 50th iterations) as the warmed up phase (“jitted”). The response time is the highest during the initial iteration, and gets slowly improved when the amount of the jitted code increases. After the 11th iteration, the response time stabilizes. For the TechEmpower benchmark, we used a similar approach as the PyPy benchmark and divided response time into the warmup phase and the warmed up phase based on the time when the number of jitted lines gets stabilized.



**Fig. 2** Number of jitted lines and response time over 50 iterations for the *html5lib* program from the PyPy benchmark suite

We applied statistical techniques to rigorously compare and quantify the differences between the response time distributions from the two phases. Statistical test like the Wilcoxon Rank Sum (WRS) test would give us a rigorous measurement if the distributions of the performance data from these two phases are different. However, in some cases, even if the distributions are different, the differences between the two distributions can be small. For example, if the response time for one request is long (e.g., more than five minutes) and the differences of the response time between the two experiments are very small (e.g., one millisecond), such performance differences would not be useful for our study, as it will not be noticed by the end-users. Hence, we also need to quantify the strength of the differences between the two distributions, called the *effect size* (Kampenes et al. 2007). In this paper, we used Cliff’s Delta (CD) as our effect size measures. Both CD and WRS are non-parametric techniques. Hence, they do not hold any assumptions regarding the distributions of the data. We consider two datasets as statistically different, when the *p*-value from the WRS test is lower than 0.05. The strength of the differences and the corresponding range of CD values (Romano et al. 2006) are shown below:

$$\text{effect size} = \begin{cases} \text{trivial} & \text{if } |CD| < 0.147 \\ \text{small} & \text{if } 0.147 \leq |CD| < 0.33 \\ \text{medium} & \text{if } 0.33 \leq |CD| < 0.474 \\ \text{large} & \text{if } 0.474 \leq |CD| \end{cases}$$

We used the following criteria to judge if the response time from the warmed up phase (denoted as B) is getting better (>), worse (<), or relatively the same (~) as the warmup phase (denoted as A):

$$\text{difference} = \begin{cases} A > B & \text{if } CD \leq -0.33 \text{ and } p - \text{value} < 0.05 \\ A \sim B & \text{if } |CD| < 0.33 \text{ or } p - \text{value} \geq 0.05 \\ A < B & \text{if } CD \geq 0.33 \text{ and } p - \text{value} < 0.05 \end{cases}$$

The *p*-values shown above are calculated from the WRS test. In other words, B improves over A ( $B > A$ ) when the WRS test and the CD value satisfy the following three conditions: (1) the two distributions are statistically significantly different ( $p$ -value < 0.05), and (2) the differences between the two distributions have medium or large effect sizes, and (3)



CD value is positive, indicating B is smaller than A. The conditions for B degrades from A ( $B < A$ ) is similar, except the CD value is negative, indicating B is bigger than A. If A and B are relatively the same ( $B \sim A$ ), when there is no statistical difference between the two distributions (a.k.a.,  $p\text{-value} \geq 0.05$ ) or the effect size between A and B is small or trivial.

We extracted the performance data from the warmup and the warmed up phases based on the time when the number of jitted lines stabilizes for all the runs of the two microbenchmark suites. We compared the response time between the two phases for each run. Table 2 shows the results.

Table 2 shows almost all the programs/scenarios (except one) for both microbenchmark suites exhibit better performance during the warmed up phase. However, in the PyPy benchmark suite, the *ai* program is not showing significant performance improvement. This is because at the end of the first iteration while running the *ai* program, the majority of the code jitting process has already been completed. Only a few lines from the test library, which does test setup, got jitted during the benchmarking process (at the 34th iteration). These additional jitted lines have no impact on the actual performance of the benchmark program. Thus, the performance differences between the two phases are very small.

**Findings:** For most of the studied programs/scenarios, the performance in the warmed up phase is statistically much better than the warmup phase. This clearly highlights the huge performance gain contributed by the JIT compilations.

**Implications:** Only performance data from the warmed up phase can be representative of the performance of systems due to the big difference in performance between the warmup and the warmed up phase. Thus, performance analysts should be careful when conducting the analysis and focus on the data after the warmup phase. In addition, the warmup phase for each system should be as short as possible, due to its inferior performance. Existing techniques for speeding up the jitting processes (Jantz and Kulkarni 2013; Gal et al. 2009; Lion et al. 2016) can be very useful in this aspect.

### 3.2 (RQ2) What is the Performance Impact by Varying JIT Configurations?

In RQ1, we have found that the system performance significantly improves after the warmup phase. Hence, the data during the warmup phase is normally discarded during the performance analysis phase for a benchmark or a performance test. In this RQ we only focus on the system performance after the warmup phase. We study the performance impact of various JIT configuration settings. In particular, we would like to (1) verify whether the system configured with the default configuration setting would yield the optimal performance amongst other configuration settings, and (2) measure the performance impact of the jitting process (a.k.a., comparing the performance against completely disabling the jitting process).

**Table 2** Comparing the response time between the warmup phase (A) and the warmed up phase (B)

Performance Difference	# of scenarios in the PyPy benchmarks	# of scenarios in the TechEmpower benchmarks
$B < A$	0	0
$B \sim A$	1	0
$B > A$	6	6
Total # of scenarios	7	6

### 3.2.1 Experiment

Similar to RQ1, we still used the same two microbenchmark suites as our experimental subjects. However, instead of keeping the default JIT configuration setting, we varied the values for the following six JIT configuration parameters: *decay*, *function\_threshold*, *threshold*, *loop\_longevity*, *trace\_eagerness*, *trace\_limit*. Table 3 shows the detailed information about these six configuration parameters. We picked these six parameters because we think they are tunable and can have an impact on where and when certain regions of the source code will be jitted. Other parameters like *enable\_ops*, *inlining*, and *off* need to be kept as default to enable the jitting process; whereas other parameters: *vec*, *vec\_all*, *vec\_cost* are not included in our study, as they are not relevant to the selected microbenchmark suites. Since there can be many possible combinations of these parameter settings, due to time constraints, we decided to run the two microbenchmark suites under the following eleven configuration settings from three different groups:

1. *Group 1 (Varying Default Configurations)* consists of five configuration settings ( $\frac{1}{4}X$ ,  $\frac{1}{2}X$ ,  $X$ ,  $2X$ , and  $4X$ ) by mutating the default configuration setting.  $X$  refers to the default configuration setting.  $2X$  means doubling the default configuration values, whereas  $\frac{1}{2}X$  means cutting the default configuration values by half and rounding to the nearest integer values. As shown in Table 4, the default PyPy JIT configuration setting,  $X$ , is: (1039, 1619, 40, 6000, 200, 1000), which corresponds to the configuration parameters (*threshold*, *function\_threshold*, *decay*, *trace\_limit*, *trace\_eagerness*, *loop\_loogevity*). Hence, the  $2X$  setting would be: (2078, 3238, 80, 12000, 400, 2000) and the  $\frac{1}{2}X$  setting would be: (519, 809, 20, 3000, 100, 500). To avoid PyPy command line parsing errors, when the value of parameter *trace\_limit* exceeds the upper bound, we just set it to be two times of the default value (12000).
2. *Group 2 (Randomly Generated Configurations)* consists of five randomly generated configuration settings ( $R_1$ ,  $R_2$ ,  $R_3$ ,  $R_4$ , and  $R_5$ ). For each parameter in the configuration setting, we randomly generated an integer value within the defined boundary. As we can see from Table 4, the randomly generated JIT configurations in Group 2 are very different from the JIT configurations from Group 1.

**Table 3** List of relevant PyPy's JIT configuration parameters and their information

Parameter	Range	Default	Descriptions
<i>decay</i>	[0, 1000]	40	Amount to regularly decay counters by
<i>function_threshold</i>	(0, $\infty$ )	1619	Number of times a function must run for it to become traced from start
<i>loop_longevity</i>	(0, $\infty$ )	1000	A parameter controlling how long loops will be kept before being freed
<i>threshold</i>	(0, $\infty$ )	1039	Number of times a loop has to run for it to become hot
<i>trace_eagerness</i>	(0, $\infty$ )	200	Number of times a guard has to fail before we start compiling a bridge
<i>trace_limit</i>	[0, 16385]	6000	Number of recorded operations before we abort tracing with ABORT_TOO_LONG

**Table 4** The JIT configurations chosen for performance evaluation

Group	Config	threshold	function_threshold	decay	trace_limit	trace_eagerness	loop_longevity
Default	$X$	1039	1619	40	6000	200	1000
Group1	$\frac{1}{4}X$	260	405	10	1500	50	250
	$\frac{1}{2}X$	520	810	20	3000	100	500
	$X$	1039	1619	40	6000	200	1000
	$2X$	2078	3238	80	12000	400	2000
	$4X$	4156	6476	160	12000	800	4000
Group2	$R_1$	64	101	120	375	200	2000
	$R_2$	519	809	5	375	200	2000
	$R_3$	519	101	20	1500	200	4000
	$R_4$	3117	1619	2	1500	25	2000
	$R_5$	259	4857	120	375	200	2000
Group 3	JIT Off	–	–	–	–	–	–

3. *Group 3 (JIT Off)* consists of only one configuration setting, which sets the parameter *off* to be true. This setting will completely disable the jitting process.

Similar to RQ1, to avoid the measurement errors and noise, we repeatedly executed each PyPy benchmark program for 30 times, and ran each TechEmpower benchmark scenario for two hours. We also collected the same kind of performance data (a.k.a., the resource utilization metrics, the JIT logs, and the benchmarking logs) for further analysis.

### 3.2.2 Data Analysis

For each experiment, we first parsed the JIT log snapshots. Based on the time when the number of jitted lines stabilizes, we divided the benchmark runs into the warmup and the warmed up phase. We extracted response time from the warmed up phase for further analysis. We used the same statistical analysis techniques as in RQ1 to compare the response time among all the runs. For each program inside the PyPy benchmark suite, we compared the performance between each pair of the JIT configuration settings and identify the best performing configuration setting. Similarly, we also located the best performing JIT configuration settings for each scenario inside the TechEmpower benchmark suite. Table 5 shows the results. There are ties when ranking the top performing configuration settings across different programs/scenarios. We noted with a “\*” besides a configuration setting if it shares the first place with other configuration settings in any programs/scenarios.

As shown in Table 5, Only 3 out of the 7 programs from the PyPy benchmark suite, where the default configuration setting yields the best performance. Furthermore, there are many configuration settings which perform best for some PyPy benchmarks (e.g.  $R_4$ ,  $R_5$ ) or share the top performance with other configurations (e.g.  $\frac{1}{4}X$ ,  $4X$ ). For the TechEmpower benchmark suite, all of the scenarios perform the best (with one scenario tied with  $2X$  and  $\frac{1}{2}X$ ) under the default configuration setting.

To quantify the performance impact of the jitting process, we also compared the performance of different JIT configuration settings against the JIT off setting. For each JIT enabled configuration setting, we measured the number of programs/scenarios that perform

**Table 5** Number of best performing programs/scenarios under each JIT configuration setting

Settings	# of best performing programs/scenarios	
	PyPy Benchmark	TechEmpoer Benchmark
$\frac{1}{4}X$	1*	0
$\frac{1}{2}X$	0	1*
$X$	3*	6*
$2X$	2*	1*
$4X$	1*	0
$R_1$	0	0
$R_2$	0	0
$R_3$	2*	0
$R_4$	1	0
$R_5$	1	0
<i>JIToff</i>	0	0

worse (<), similar ( $\sim$ ), or better (>) than the JIT off setting. Table 6 shows the result. The number of programs/scenarios whose performance under jit enabled configurations is worse or no different than jit off is marked as bold.

From Table 6, we can see that, among the PyPy benchmark suite, all the JIT enabled configuration settings perform better than the JIT off setting, except  $R_1$  in which the performance of two PyPy benchmark programs is even worse than completely disabling the jitting process (a.k.a., JIT off)! Similarly, in the TechEmpower benchmark suite,  $R_1$  is still the odd one, as none of its scenarios is better than the JIT off setting. In addition, only three of the TechEmpower scenarios under  $R_5$  are better than the JIT off setting.

**Table 6** Comparing the jitting performance against the configuration under JIT off

Configs	PyPy Benchmark			TechEmpower Benchmark		
	<	$\sim$	>	<	$\sim$	>
$\frac{1}{4}X$	0	0	7	0	0	6
$\frac{1}{2}X$	0	0	7	0	0	6
$X$	0	0	7	0	0	6
$2X$	0	0	7	0	0	6
$4X$	0	0	7	0	0	6
$R_1$	<b>2</b>	0	5	<b>5</b>	<b>1</b>	0
$R_2$	0	0	7	0	0	6
$R_3$	0	0	7	0	0	6
$R_4$	0	0	7	0	0	6
$R_5$	0	0	7	<b>2</b>	<b>1</b>	3

The number of programs/scenarios whose performance under jitted configuration is worse or no different than JIT off setting is highlighted in bold

**Findings:** Programs/scenarios running under the default configuration setting do not necessarily yield the best performance when comparing to other configuration settings. The optimal JIT configuration setting can vary depending on the programs/scenarios. The performance of some of the JIT enabled configurations can be worse than turning JIT off.

**Implications:** PyPy's JIT configuration settings have a big impact on the system performance. It is important to find the optimal JIT configuration setting for each system to achieve the best performance.

### 3.3 (RQ3) Do Systems Containing more Jitted Lines Yield Better Performance?

In RQ2, we have found that the default JIT configuration setting does not necessarily result in the optimal performance. Different JIT configuration settings would result in different portions of the code been jitted. However, does more jitted code always lead to better performance? In this RQ, we want to study the relationship between these two aspects.

#### 3.3.1 Experiment

We used the same data from RQ2 and did not run any additional experiment for this RQ.

#### 3.3.2 Data Analysis

We first selected the configuration setting that has the best performance for each program/scenario based on the results of RQ2. Then we also selected the configuration settings with the highest number of jitted lines. If there are two different configuration settings corresponding to the above two criteria, we further performed the WRS test and calculated the CD value between the performance data under those configuration settings.

Table 7 shows the effect size between the best performing and the most jitted configuration settings for each program/scenario. Since there can be ties in either category, we compared all pairs of configuration settings from the best performing category to the category of the largest portion of jitted code. We label *True* at the third column, if there is at least one common configuration setting in both categories for one program/scenario. In 71% ( $\frac{5}{7}$ ) of the programs in the PyPy benchmark suite and all the scenarios in the TechEmpower benchmark suite, the best performing configuration setting is different from the one that has the highest number of the jitted lines. When comparing the performance differences, we compared all the pairs of these configuration settings from the two categories. In the end, twelve of the programs/scenarios have a medium to large effect size differences. In other words, the results show that more jitted lines do not necessarily lead to better performance. For the PyPy benchmark 'html5lib', we marked the effect sizes as '-', since the best performing configuration settings and the highest amount of the jitted code configuration settings are exactly the same. Hence, we do not calculate the effect sizes for this case.

**Findings:** JIT configuration settings, which resulted in the highest number of the jitted lines, do not necessarily yield the best performance.

**Implications:** We cannot just arbitrary choose the configuration settings that favor more jitted lines of code while tuning the system performance. A more sophisticated approach is needed to locate the optimal configuration setting(s) for one system.

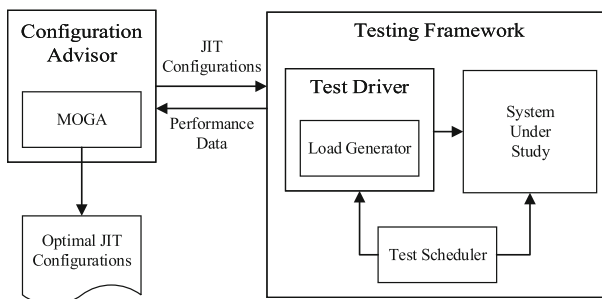
**Table 7** Comparison between the configuration settings yielded the best performance and the configuration settings resulted in the most jitted code

PyPy Benchmark			TechEmpower Benchmark		
Programs	Effect Size	Same	Scenarios	Effect Size	Same
ai	large	False	db	large	False
bm_mako	trivial,medium	False	fortune	large	False
chaos	large	False	json	large	False
django	large	False	plaintext	large	False
html5lib	–	True	query	large	False
pidigits	large	True	update	large	False
rietveld	large	False			

#### 4 Automatically Tuning the JIT Configuration Parameters

In the previous section, we have found that PyPy’s JIT configuration settings do have a significant impact on the system performance. Furthermore, there is no straightforward way to recommend a performance-efficient JIT configuration setting, since such setting can be application-dependent and a higher portion of the jitted code does not necessarily result in better performance. Hence, in this section, we will propose our automated approach, ESM-MOGA (Effect Size Measurement-based Multi-Objective Genetic Algorithm), to tuning the JIT configuration parameters for one system.

Figure 3 provides an overview of our tool which we called the PyPyJITTuner. It leverages a search-based technique called Multi-Objective Genetic Algorithm (MOGA) (Deb et al. 2000) and a statical measure called effect size, for the exploration of the JIT configuration space. Genetic Algorithm (GA) is a search-based method inspired by evolutionary biology, in which a population of solutions is evolved during each generation. The solutions from the next generation should be generally better than the previous generations evaluated based on some objective functions. MOGA is a type of GA, in which multiple objectives are being considered. We chose MOGA, as there can be multiple objectives associated with a system’s performance (e.g., optimizing the response time for multiple scenarios). One machine, which is deployed with the tailored-version of the MOGA, acts as the configuration advisor. When new solutions have been created, this machine continuously sends the

**Fig. 3** Overall process of PyPyJITTuner

JIT configuration settings (solutions) to the test scheduler machine, which will deploy and configure the system under study (SUS). Once the test scheduler machine receives these settings, it will reset the test environment (a.k.a., clean up the database, and remove the testing data from the previous run), and start up the SUS under the suggested JIT configuration setting. The same performance test (a.k.a., the same workload) will be executed. Once the test is completed, the performance data will be collected and sent to the configuration advisor machine for further analysis. The configuration advisor machine will evaluate the newly received performance data against the data from other configuration settings and leverage the MOGA methods to select the best solutions and generate the next generation. If the solutions in the next generation are good enough, the MOGA will stop the evolution and output one or multiple “optimal” configuration settings. Otherwise, the MOGA will continue with another round of iteration. The newly generated JIT configuration settings will be sent to the test scheduler machine for another round of testing.

The rest of this section is organized as follows: Section 4.1 explains the general idea behind the ESM-MOGA approach. Section 4.2 presents our performance testing framework, and Section 4.3 describes briefly our implementation.

### 4.1 Tailoring MOGA for JIT Configuration Tuning

GA is a search-based method inspired by evolutionary biology. GA encodes the candidate solutions into a set of values, called “chromosomes”. Inside the chromosomes, the set of values, which are to be optimized are called the “genes”. GA starts off with a population of the initial solutions and keeps iterating until any of the termination criteria is met. During each iteration, GA improves the population via crossover (combining existing solutions to produce new solutions), mutation (randomly changing some values in the solutions), and selection (picking the best candidate solutions). The termination criteria can either be the optimization conditions (e.g., the resulting solutions are better than a predefined threshold) or the maximum number of iterations. The MOGA, which is a type of GA, evaluates multiple objectives simultaneously. In general, as illustrated in Fig. 4, the MOGA consists of six phases: the problem formulation phase, the initialization phase, the tournament phase, the evolution phase, the selection phase, and the stopping phase. The process of going through the tournament, the evolution, and the selection phase can be repeated multiple times, with each iteration called one generation. At the end of each generation, a new population will be produced. This process will be repeated until any termination criteria described in the stopping phase is met.

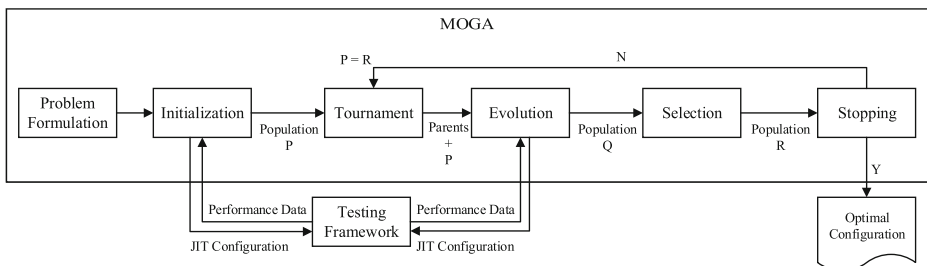


Fig. 4 The workflow for our tailored version of the MOGA method

In this subsection, we will explain the ESM-MOGA approach by using a running example. For illustration purposes, we assume the SUS in our running example is a simple e-commerce system, which consists of only three scenarios: login, browse, and purchase.

#### 4.1.1 Phase 0 - Problem Formulation

We formulated our problem of automated tuning of JIT configuration settings into a multi-objective optimization problem. Our objective is to find one or more optimal JIT configuration settings that yield the best performance in all the scenarios in the system. Below we define our solution encoding and objectives:

- **Solution Encoding:** The ESM-MOGA requires us to encode its solution into binary strings. As shown in Section 2.2, all the studied JIT configuration parameters are integers and have a large range (a.k.a., many possible values). Hence, we decided to select eight representative values from the input domain of each configuration parameter: ( $4X$ ,  $3X$ ,  $2X$ ,  $X$ ,  $\frac{1}{2}X$ ,  $\frac{1}{4}X$ ,  $\frac{1}{8}X$ ,  $\frac{1}{16}X$ ), where  $X$  refers to the default value for that configuration parameter. We chose the above eight levels, as these eight values cover a wide range of the input domain and each configuration parameter value can be easily encoded into a binary string of length three. In this way, the smallest configuration value ( $\frac{1}{16}X$ ), the default configuration value ( $X$ ), and the largest configuration value ( $4X$ ) for each parameter are encoded as 111, 011, and 000, respectively. We set the largest configuration values to be  $4X$ , as Section 3 shows that large JIT configuration settings usually do not yield good performance. We set the smallest configuration values to be  $\frac{1}{16}X$ , as Section 3 shows very small configuration settings could result in a high number of jitted code, but worse performance. To ease explanation, in our running example, there are only two configuration parameters. Hence, the default configuration setting can be encoded as a binary string: 011011.
- **Objectives:** For a real world system, there can be more than one aspect associated with the performance of the system. Examples of optimizing performance aspects can be optimizing the resource utilizations (e.g., CPU, memory, and disk) or the responsiveness of different scenarios in a system. Some of these concerns can be conflicting with each other. In our approach, we focus on optimizing the response time for different scenarios in a system. Each objective refers to a list of response time for each scenario during the warmed up phase, measured through performance testing. In our running example, the objectives are to optimize the response time for the above three scenarios in the e-commerce system.

#### 4.1.2 Phase 1 - Initialization

During the initialization phase, the ESM-MOGA defines an initial population ( $P$ ) consisting of  $n$  solutions. In our approach,  $P$  consists of the default configuration setting, and  $n - 1$  randomly generated configuration settings. The ESM-MOGA will intentionally included the default configuration setting in the initial population, as we want to ensure the default configuration setting is evaluated among its alternatives and the final “optimal” setting(s) will be at least as good or better than the readily available default configuration setting. Once the initialization process is completed, the ESM-MOGA enters the iterative process of going through the tournament, the evolution, and the selection phase to refine and improve its population until any termination criteria is met.



To ease explanation, we set  $n = 4$  in our running example and compose the initial population (P) with the following solutions:

$$P = \begin{cases} C_1 : 011011, \\ C_2 : 001010, \\ C_3 : 000101, \\ C_4 : 101110 \end{cases}$$

Once the initial population is generated, the solutions in the initial population will be sent to the test scheduler machine in Section 4.2. Multiple performance tests with the same workload will be conducted under each given configuration setting. The response time of the three scenarios under warmed up phase will be collected as performance data and assigned as the objectives for each solution.

### 4.1.3 Phase 2 - Tournament

During the tournament phase, the ESM-MOGA will first randomly select two solutions from a pool which contains all solutions of the current population. Then, a pairwise comparison is done to recognize the better solution from the two. The better solution will be used as one of the parents for the next phase. These evaluated solutions will not be put back to the pool for efficient concerns. The process will be repeated until all solutions in the pool have been evaluated pairwise. In our approach, the pairwise comparison is done using a pre-defined dominant comparison function and the dominating configuration setting (a.k.a., the better solution) will be selected. In this dominant comparison function, the configuration setting A dominates the configuration setting B, if the response time distributions under the two configuration settings satisfy the following two criteria:

1. **The response time for all the scenarios under A are statistically no worse than under B:** The response time under configuration A for one scenario is statistically no worse than under B, if (1) the response time distributions for that scenario under the two settings are not statistically significantly different under the WKS test, or (2) they are statistically significantly different under the WKS test, but there is only a trivial to small effect size calculated by the CD. This relation has to be held for all the scenarios when comparing the two settings.
2. **There is at least one scenario whose response time under A is statistically better than under B:** The response time under configuration A for one scenario is statistically better than B, if the response time distributions for that scenario under the two settings are statistically significantly different under the WKS test and there is a medium to large effect size calculated by the CD.

In other words, one configuration setting (A) only dominates the other one (B), if (1) the performance of all the scenarios under A is at least as good as B, and (2) there will be at least one scenario under A whose performance is better than B.

The dominance comparison among all the pairs of the configuration settings are shown in Table 8. Each row in Table 8 corresponds to the comparison results of one configuration setting pair. For example, the second row shows the comparison results between configuration setting  $C_1$  and  $C_3$ . The response time for the login scenario is statistically better under  $C_1$  than  $C_3$ . The performance of the other two scenarios are statistically not different between the two settings. Hence, the configuration setting  $C_1$  dominates  $C_3$ . Assume, from the pool

**Table 8** Dominance relations among the four configuration settings in our running example

Config Pairs	Login	Browse	Purchase	Dominance
$(C_1, C_2)$	Better	Worse	Better	$\approx$
$(C_1, C_3)$	Better	Equal	Equal	$\succ$
$(C_1, C_4)$	Better	Better	Equal	$\succ$
$(C_2, C_3)$	Better	Equal	Equal	$\succ$
$(C_2, C_4)$	Better	Better	Equal	$\succ$
$(C_3, C_4)$	Equal	Worse	Worse	$\prec$

“ $\succ$ ” means the the left configuration setting dominates the right configuration setting, “ $\prec$ ” means the right configuration setting dominates, and “ $\approx$ ” means there is no dominance relation

that contains all solutions of population (P), we selected  $(C_1$  and  $C_3)$ ,  $(C_2$  and  $C_4)$  for pairwise comparison. The two configuration settings,  $C_1$  and  $C_2$ , will be selected as parents for the next phase.

#### 4.1.4 Phase 3 - Evolution

During the evolution phase, the parents from the Tournament phase will undergo the following two actions to produce new solutions:

- **Crossover:** Two solutions from the Parents are randomly selected as parents. A new solution will be created by randomly selecting some bits from one solution and the remaining bits from another solution. In our running example,  $C_1$  and  $C_2$  will be selected as parents. A new solution  $C_5$  (011010) can be created by inheriting the first three bits from  $C_1$  and the remaining bits from  $C_2$ .
- **Mutation:** Some of the newly produced solutions will be mutated by randomly flipping some bits (a.k.a., turning 0s into 1s and 1s into 0s). For our running example, after flipping the first and the last bits of  $C_5$ , it becomes 111011.

Similar as the Initialization phase, a performance test with the same workload but a new configuration setting ( $C_5$ ) will be conducted. Once completed, the performance data will be sent back as objectives for the new configuration setting. The overall population (Q) at the end of this phase will consist of the new solutions produced after the crossover and the mutation operations as well as existing solutions from P.

#### 4.1.5 Phase 4 - Selection

In this phase, the “best”  $n$  solutions in Q will be selected with NSGA-II selection (Deb et al. 2000) (Evaluation tools in Python (DEAP) 2017). It will first use the non-dominated sorting algorithm to sort the solutions into different levels ( $L_0, L_1, \dots$ ). Solutions that were dominated by the smallest number of solutions will be assigned to the top level ( $L_0$ ). Solutions in  $L_0$  are the “best” solutions in this iteration, followed by the solutions in  $L_1$ , and then  $L_2$ , and so on. The solutions within the same levels (e.g.,  $L_0$ ) are not dominant over each other. For example, if configurations settings A and B are both within  $L_0$ , it means that A and B are not dominant over each other. In other words, some scenarios are better performed under A and whereas some other scenarios are better performed under B. When selecting the top  $n$  solutions, we will first start picking solutions from the top level ( $L_0$ ), followed by solutions

from  $L_1$ , and so on. If there are more solutions in a level than we needed (a.k.a., exceeding the total  $n$  solutions), we will rank solutions within that level with crowding distance sorting and select the top ranked solutions in that level.

Suppose for our running example, after sorting the solutions in  $Q$  using the non-dominated sorting algorithm, they are divided into the following levels:

$$Q = \begin{cases} L_0 : C_5, \\ L_1 : C_1, C_2, \\ L_2 : C_4 \\ L_3 : C_3 \end{cases}$$

Since at the end of the selection phase, only  $n$  solutions will be kept. Hence, our resulting population ( $R$ ) will be  $C_1, C_2, C_4, C_5$ .

#### 4.1.6 Phase 5 - Stopping

During this phase, the resulting population  $Q$  formed during this generation will be evaluated to decide whether its solutions are good enough comparing to the previous generation. The main idea is to decide whether any progress has been made during this generation. In other words, we want to check whether there are any better solutions produced during this generation. We used the Mutual Dominance Rate (MDR) (Martí et al. 2009) to measure the improvements made between the current population  $B$  and the population  $A$  from the previous generation:

$$MDR(B, A) = \frac{dom(B, A)}{\|B\|} - \frac{dom(A, B)}{\|A\|},$$

where  $dom(A, B)$  is defined as the number of solutions in population  $A$  that are dominated by at least one solutions in  $B$ . Hence, for our running example,  $dom(P, R)$  would be 4, since all four solutions in  $P$  are dominated by at least one solutions in  $R$ .  $dom(R, P)$  would be 1, since only  $C_4$  in  $R$  is dominated by the solutions in  $P$ . Hence,  $MDR(R, P) = \frac{1}{4} - \frac{4}{4} = -\frac{3}{4}$ . The closer the MDR value gets to  $-1$ , the larger the improvement has been made in the current generation. If MDR is close to 0, it means little progress has been made to the population. The iteration should be stopped if the improvement between two generations is insignificant (a.k.a.,  $|MDR|$  is smaller than some threshold values), or it has been running for too long (e.g., over 100 iterations).

When the termination criterion is met, we will output the top configuration settings for the current generation with the NSGA-II selection. Since there is only one solution ( $C_5$ ) at  $L_0$  for our running example,  $C_5$  will be the optimal configuration setting outputted.

## 4.2 Our Performance Testing Framework

For any newly generated solutions, we need to measure their performance using a performance test. Each solution, which is sent to the test scheduler machine inside the Testing Framework, will first be parsed into the corresponding JIT configuration setting. The test scheduler machine will then start the system with the new configuration setting and measure the system performance under a predefined workload. At the end of each performance test, performance data during the warmed up phase will be collected and sent to the configuration advisor machine, so that they can be used as objectives to evaluate among solutions in the ESM-MOGA.

For our running example, a total of five performance tests with the same workload but different JIT configuration settings, which correspond to the four initial solutions in  $P$  and the new solution in  $Q$ , will be run. Once each test is completed, the test scheduler shuts

down the SUS, collects the performance data (response time for the individual scenarios, the resource utilization metrics, and JIT logs) and sends the data to the configuration advisor machine.

### 4.3 Implementation

We implemented the ESM-MOGA using the NSGA-II algorithm (Deb et al. 2000), which is a fast and efficient multi-objective genetic algorithm, from the DEAP framework (Distributed Evolutionary Algorithms in Python) (Distributed Evolutionary Algorithms in Python (DEAP) 2017). The framework contains the relevant library functions for NSGA-II, like assigning crowding distance, non-dominated sorting algorithm, and NSGA-II selection. We had to implement the dominance function, and input encoding ourselves to fit into NSGA-II algorithm. We re-implemented the non-dominate sorting and NSGA-II selection functions so that they can use our dominance function to compare among solutions.

We also implemented the automated performance testing framework, which leverages JMeter (Apache JMeter 2015) as the load generator. The performance testing framework can startup, initialize, execute, and stop a performance test under one particular configuration setting.

## 5 Case Study

In this section, we evaluated the performance of our automated approach to tuning the JIT configuration parameters on three Python-based open source systems: Saleor (Saleor - An e-commerce storefront for Python and Django 2017), Wagtail (Wagtail CMS: Django Content Management System 2017), and Quokka (Quokka CMS (Content Management System) - Python 2017). As shown in Table 9, these three systems vary from system sizes, application domains, and technology stacks. All three systems can be deployed on top of the Tornado WSGI server which uses Gunicorn for worker process management. And they all require a database to be functional. Saleor is an e-commerce system, built using the Django framework, and uses Postgres as its database. Wagtail shares the same technology stack as Saleor (a.k.a., Django and Postgres), but is from a different application domain: the Content Management System (CMS). Although Quokka is also a CMS, it is built with the Flask framework and uses MongoDB as its database.

The rest of this section is organized as follows. Section 5.1 describes the case study setup. Section 5.2 explains the case study results.

### 5.1 Case Study Setup

We deployed the above three systems on the same physical machines, which have the following hardware configurations: *Intel i7-4790* CPU, 16 GB memory, and 2 TB hard-drive.

**Table 9** An overview of the three Python-based systems under study

Name	Version	LOC	Application Domain	Technology Stack
Saleor	2017.07.0	48482	E-commerce	Gunicorn,Tornado,Postgres,Django
Wagtail	1.12.1	85006	Content Management	Gunicorn,Tornado,Postgres,Django
Quokka	0.2.1	34468	Content Management	Gunicorn,Tornado,MongoDB,Flask

JMeter was deployed on another physical machine with the following hardware configuration: *Intel(R) Core(TM)2 Duo* CPU, 4 GB memory, 160 GB hard-drive. And all machines have Ubuntu 14.04 deployed. The reason for the separate deployment of JMeter and the SUSs is to ensure no overhead caused by the load generator to the SUSs (Jiang and Has-san 2015). The version of the PyPy that we used for evaluation is 5.7.1 which corresponds to Python version 2.7.13. Similar to the exploratory study, we focused on the same six JIT configuration parameters. Hence, each solution (a.k.a., configuration setting) requires 18 bits to be encoded into our tailored MOGA method. For example, the default configuration setting would be encoded as 011011011011011011. As for the rest of the MOGA configurations, we set the initial population ( $P$ ) size as 40, and the mutation rate as 0.10 based on some small trials. We also configured our termination criteria to be either  $|MDR| \leq 0.1$  holds for two consecutive generations or the MOGA has iterated for 10 generations.

Table 10 shows the workload that we have set for the three systems. The workload tries to simulate how real users use the systems in the field. The overall workload intensity (10 requests/sec) is the same for all three systems and the workload mix for each system is shown below.

- For the e-commerce system, Saleor, the workload tries to mimic the purchasing workflow from a real customer. Hence, we divided this scenario into 12 actions, which correspond to 12 different webpage operations. The overall workload intensity (10 request/sec) corresponds to 10 different users performing the above 12 actions at the same time. Hence, all the actions in this workload are assigned with the same ratio in the workload mix.

**Table 10** Workload description for the three systems

System	Workload Mix	Workload Intensity
Saleor	load index page load login page login request view category view product add product to cart view cart check out cart select shipping method select shipping address payment payment confirm	10 req/sec
Wagtail	add blog(1) add event(1) edit blog(2) view blog page(3) view event page(3)	
Quokka	add blog(1) edit blog(2) view blog(7)	

- For the content management system, Wagtail, the workload tries to mimic users reading, posting, and editing blogs or events. Since, in the majority of the time, users will be viewing the blogs or events, we assigned a higher ratio for these two actions. The scenarios of adding a new blog or a new event happen the least frequently. Hence, they are assigned with the smallest weight in the workload mix.
- Quokka’s workload is similar to Wagtail’s, as they are in the same application domain. Since Quokka only supports reading/editing/adding blogs, we adjusted the workload mix accordingly.

As Section 3 indicates, only the performance data from the warmed up phase is representative of the actual system performance. Hence, in the case study, we want to make sure the system has been running long enough (a.k.a., finished the warmup up phase). To properly decide the test duration, we first did a test run with the default configuration setting, in which the predefined workload was executed for three hours. We leveraged a similar technique as Alghamdi et al. (2016) to test when the system’s performance behavior gets repetitive, so that we can identify the duration of the warmup phase. We divided the collected performance data into intervals of every 20 minutes. Then we performed statistical analyses with the WRS test and CD values on the response time between two adjacent time intervals. We considered the system to be fully warmed up when the response time from the two adjacent time intervals show insignificant difference for all scenarios (a.k.a., not statistically different by the WRS test or CD values show trivial to small effect sizes).

We performed the above process in all three case study systems. We found that all three systems finish the warmup phase in the first 40 minutes before their performance behaviors start to be repetitive. Hence, for consistency concerns, we set the test duration to be 50 minutes for each test and only used the data from the last 10 minutes (a.k.a., the data from the warmed up phase) for further analysis in the ESM-MOGA method.

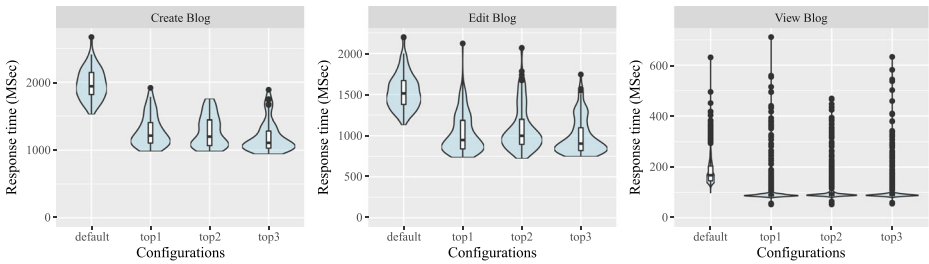
## 5.2 Case Study Results

For all the three systems, we applied our ESM-MOGA method with the aforementioned setup. Table 11 shows the runtime statistics for the ESM-MOGA method after running on the three systems. The search algorithm all terminates under termination criteria  $|MDR| \leq 0.1$ . For Saleor, it takes 36 hours and evaluated 100 solutions, 67 hours and 199 solutions for Wagtail, and 35 hours and 96 configuration settings for Quokka. For all three systems, the ESM-MOGA method found optimal solutions when it terminated.

For each system, we used the NSGA-II selection to select the top three configuration settings. We compared the response time and the resource utilization between the optimal configuration settings and the default configuration setting. Figure 5 visually compares the response time distributions between the default and the top three optimal configuration settings for Wagtail. Due to space limitations, it only contains the performance comparisons of

**Table 11** Statistics after running the MOGA approach on the three case study systems

System	# of Generations	# of Configurations Evaluated	Duration (hour)
Saleor	3	100	36
Wagtail	7	199	67
Quokka	3	96	35



**Fig. 5** Visualizing the response time distributions of different scenarios under different configuration settings for Wagtail

the three blog-related scenarios inside Wagtail. Each sub-figure corresponds to one scenario. Within each sub-figure, the four violin plots correspond to the response time distributions of that scenario under the default and the three optimal configuration settings. Among all the sub-figures, the response time under the default configuration setting is significantly much higher than the optimal configuration settings. A similar trend also holds for the two event-related scenarios in Wagtail.

To quantify the differences between the optimal and the default configuration settings, we performed the WRS test and CD test between the response time distributions under each configuration pair. We used the same criteria as in Section 3.1 to judge whether one response time distribution is better, or same, or worse than the other one. Table 12 shows the results. Each row corresponds to one optimal configuration setting for one particular system. There are no orderings among the top three optimal configurations ( $O_A$ ,  $O_B$ , and  $O_C$ ). The second to the fourth columns contain the number of scenarios which show better, equal, or worse difference when comparing this configuration setting against the default. For Saleor and Wagtail, all the scenarios performed better under the suggested optimal configuration settings. For Quokka, at least one scenarios performed better under the suggested optimal configuration settings, while the remaining scenarios performed no worse than the default configuration setting. The fifth and the sixth columns show the minimum and maximum percentage of differences when comparing the average response time under the suggested configuration setting with the average response time under the default for all scenarios. The percentage improvement in terms of the average response time can vary between 5% to 60%.

Since each system consists of a web server and a database, we further compared the CPU and the memory utilizations between the optimal and the default configuration settings. Last four columns in Table 12 shows the comparison results for the resource utilizations. The CPU usage for both components drops. The decrease in CPU is more significant in the web server, with the average improvement ranges between 12% to 33.7%. However, the memory usage for the web server dramatically increases (12.7% to 202.5%) across all the optimal configuration settings. We suspect this may be due to the storage of the compiled jitted code. For a more detailed discussion, please refer to Section 6.4.

## 6 Discussions

In this section, we discussed the findings based on the case study results and their implications.

**Table 12** Comparing the performance between the optimal configuration settings and the default configuration setting for the three systems. The optimal configurations are labelled as  $O_A$ ,  $O_B$ , and  $O_C$ . WS stands for the web server, and DB stands for the database. “-” means the ESM-MOGA suggested optimal configuration setting outperforms the default setting and “+” means otherwise

System	Top Configurations	Comparing scenarios			Average differences (%)		Average resource usage difference (%)			
		better	equal	worse	min	max	WS_CPU	WS_Memory	DB_CPU	DB_Memory
Saleor	$O_A$	12	0	0	-27.11	-56.17	-33.71	+24.80	-3.61	-0.24
	$O_B$	12	0	0	-27.32	-58.93	-25.68	+31.97	-7.59	+1.18
	$O_C$	12	0	0	-9.66	-60.28	-32.38	+12.74	-7.30	+0.84
Wagtail	$O_A$	5	0	0	-33.47	-45.10	-23.86	+158.71	-5.96	+1.55
	$O_B$	5	0	0	-29.59	-44.93	-18.81	+202.53	-2.90	-2.04
	$O_C$	5	0	0	-35.93	-44.28	-25.18	+157.94	-3.46	-3.34
Quokka	$O_A$	2	1	0	-9.39	-34.95	-15.39	+59.34	-13.80	-5.5688
	$O_B$	3	0	0	-13.45	-22.71	-16.51	+61.08	-21.89	+0.942
	$O_C$	1	2	0	-4.98	-25.44	-11.94	+62.03	-19.05	+2.21



**Table 13** Top three optimal configuration settings for the three studied system

System	Config.	decay	function_threshold	loop_longevity	threshold	trace_eagerness	trace_limit
Saleor	$O_A$	10	101	2000	129	25	12000
	$O_B$	2	1619	250	1039	12	12000
	$O_C$	2	1619	1000	3117	100	12000
Wagtail	$O_A$	10	101	1000	1039	12	12000
	$O_B$	2	404	250	259	12	12000
	$O_C$	2	101	4000	4156	50	12000
Quokka	$O_A$	80	404	500	519	12	12000
	$O_B$	2	3238	2000	4156	12	6000
	$O_C$	2	202	4000	4156	25	3000
Default	—	40	1619	1000	1039	200	6000

## 6.1 Optimal Configurations across Different Systems

As we can see from the previous section (Section 5), the optimal configuration settings obtained using ESM-MOGA significantly out-performed the default configuration. In this subsection, we would like to compare the optimal configuration settings against the default configuration setting in order to see if we can derive some rules or provide some guidance for PyPy users when tuning the JIT configuration settings for their systems.

For each of the studied systems above, we obtain its top three optimal configuration settings, whose values are shown in Table 13. We also included the default configuration setting in the table to ease comparison.

One common pattern as we can see from Table 13 is that *trace\_eagerness* is significantly smaller in all the optimal configuration settings when comparing to the default ones. *trace\_eagerness* refers to the eagerness to compile a non-jitted branch within a loop. A system can go through various branches within a loop. A smaller *trace\_eagerness* is preferred, so that the branch(es) corresponds to frequently executed scenarios will be jitted faster. There is no pattern found in the other five configuration parameters, as they can be either bigger or smaller than the default values among the nine optimal configuration settings.

We are also interested in understanding the correlation of JIT configuration parameters to the system performance. We collected all the evaluated configuration settings and the corresponding response time for different scenarios. We summed up the average response time for all scenarios as the overall response time under a JIT configuration setting. Then we calculated the Spearman's rank correlation between each configuration parameter and the overall response time. Spearman is a non-parametric correlation metric measuring the strength of the relation between the two variables. The scale of the Spearman's  $\rho$  correlation coefficient is indicated below (Hopkins 2016):

$$\rho = \begin{cases} \text{trivial} & \text{if } 0 \leq \rho < 0.1 \\ \text{small} & \text{if } 0.1 \leq \rho < 0.3 \\ \text{moderate} & \text{if } 0.3 \leq \rho < 0.5 \\ \text{large} & \text{if } 0.5 \leq \rho < 0.7 \\ \text{very large} & \text{if } 0.7 \leq \rho < 0.9 \\ \text{near perfect} & \text{if } 0.9 \leq \rho \leq 1.0 \end{cases}$$

Table 14 shows the result of the correlation between each configuration parameter and the overall system performance. We highlight the cell in bold if the correlation measure is “large” or “very large”. Each system has at least one configuration parameter which has a “large” or “very large” correlation measure. However, highly correlated configuration parameters vary among the three systems, except *trace\_limit*. A large *trace\_limit* enables the system to compile large frequently executed loops, which can subsequently improve the system performance. Meanwhile, the three JIT configuration parameters *function\_threshold*, *loop\_longevity* and *threshold* have very low correlations (small or trivial) with the overall system performance.

**Findings:** The configuration parameter *trace\_eagerness* should be generally set lower than the default values in order to obtain better performance. *trace\_limit* is highly correlated with the overall system performance, whereas *function\_threshold*, *threshold* and *loop\_longevity* have no or weak correlations.

**Implications:** There are some general guidance in terms of tuning the PyPy JIT configuration settings on web frameworks. However, the optimal configuration settings are still highly system dependent. In this paper, we only evaluated the impact of PyPy JIT configuration settings on benchmark programs or web applications. However, Python is also popular in data statistic analysis and machine learning (e.g. *scipy*, *tensorflow*), which could be time consuming to train a model. One of the interesting future research area would be to derive rules or general guidance to improve the performance of various machine learning or statistic analysis packages by tuning their JIT configuration parameters.

## 6.2 Top Configurations Across Different Workloads

In the previous study, we compared the optimal JIT configuration settings and its performance cross different applications under the same level of workload. In this subsection, we want to compare the optimal JIT configuration settings under different workloads. We selected Wagtail as our experiment subject.

In addition to the default Wagtail workload (10 req/sec), we generated two other workloads for comparison: 15 req/sec and 5 req/sec, while keeping the workload mixes. For each

**Table 14** Spearman correlation between configuration and response time

System	Correlation	decay	function_threshold	loop_longevity	threshold	trace_eagerness	trace_limit
Saleor	corr. coeff.	-0.5120	-0.1338	0.0291	0.1869	-0.2961	0.6293
	<i>p</i> - value	6.057e-08	0.1864	0.7748	0.0638	0.0029	3.028e-12
	scale	<b>large</b>	small	trivial	small	small	<b>large</b>
Wagtail	corr. coeff.	-0.3719	-0.2750	0.1506	0.1392	-0.5446	0.5790
	<i>p</i> - value	6.346e-08	8.452e-05	0.0336	0.0498	2.2e-16	2.2e-16
	scale	moderate	small	small	small	<b>large</b>	<b>large</b>
Quokka	corr. coeff.	-0.0403	-0.1932	0.1157	0.0568	0.0506	0.8024
	<i>p</i> - value	0.7554	0.1323	0.3703	0.6608	0.6956	4.476e-15
	scale	trivial	small	small	trivial	trivial	<b>very large</b>

The large and very large correlation measures are shown in bold

of newly generated workload, we ran PyPyJITTuner to derive the optimal JIT configuration settings. For 15 req/sec workload, the framework iterated for 5 generations before termination. And it takes 4 generations for the PyPyJITTuner to be terminated under 5 req/sec workload. The resulting configuration settings yield significant performance gain (20% - 50%) when compared to the default configuration setting.

Table 15 shows the actual performance for the top three configuration settings under each workload. It shows the average response time for each scenario in milliseconds, as well as various resource usage metrics like CPU and memory usage for the web server and the database, respectively. In addition, it also shows the number of jitted lines under each configuration setting. Although the workloads are different, all the nine top optimal configuration settings share similar performance in terms of response time for each scenario. The number of jitted lines are similar across different workloads. As workload increases, the CPU and memory consumption for both the web server and the database increases. Hence, in this case, the workload intensity is the main reason behind the increase of memory and CPU consumption of the two server components.

Table 16 shows the top three optimal JIT configuration settings under different workloads. The optimal JIT configuration settings under different workloads share very similar properties (e.g., large *trace\_limit* and small *decay* and *trace\_eagerness* values).

**Findings:** Varying the workload intensity would not impact the optimized JIT performances. And different workloads would result in different optimal configuration values. However, there are some common properties (e.g., large *trace\_limit* and small *decay* and *trace\_eagerness* values) shared across them.

**Implications:** Researchers could further improve the efficiency of the ESM-MOGA by developing machine learning algorithms to proactively eliminate some of the performance deficient configuration settings from each generation.

### 6.3 Code Jitting vs. Performance

In Section 3.3, we have shown that more jitted lines do not necessarily lead to better performance. Furthermore, the performance under some of the JIT configuration settings are even worse than turning the JIT completely off! In this subsection, we would like to perform a more in-depth study to find out the reasons.

In Fig. 6, we plotted the number of jitted lines with respect to system performance across all the runs we did for the Saleor system. And the red dotted line shows the overall average response time under the JIT off configuration setting. As the number of jitted lines increases, the average response time for the system gradually decreases. As we can see from the figure, there are a few JIT configuration settings which are even worse than turning the JIT off! We conducted further analysis to understand the reason why some jitted code would even lead to worse performance.

We focused on comparing the code structure under two configuration settings: configuration A, which is a configuration with JIT enabled. As shown in Fig. 6, the performance of configuration A is worse than turning JIT off. For brevity, we call the JIT off configuration as configuration B.

We first applied cProfile (The Python Profilers 2018) to gather the high-level performance numbers for Saleor. cProfile is a profiling tool for Python-based systems. It can provide information like the execution time, the number of execution for each of the executed functions. We enabled cProfile and ran the Saleor under the default workload twice:

**Table 15** Comparing the performance among the optimal configuration settings under different workloads for Wagtail

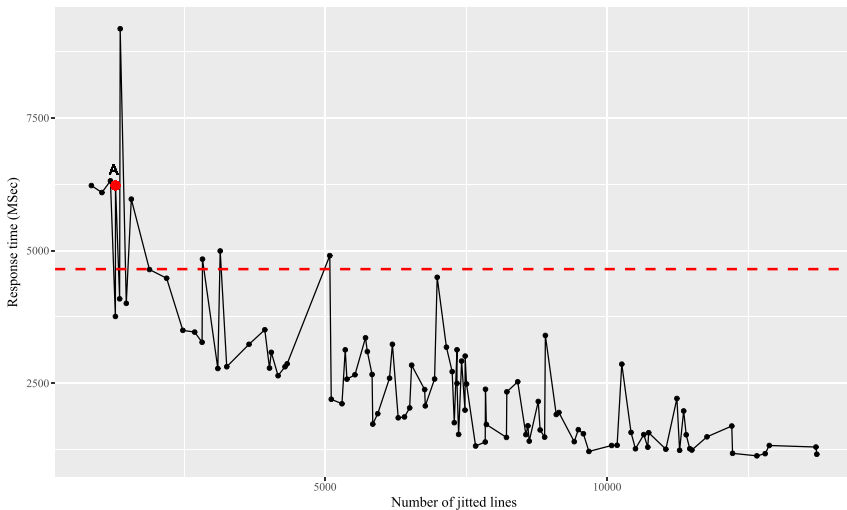
Experiments	Workload Top configs	5 req/sec			10 req/sec			15 req/sec		
		$O_A$	$O_B$	$O_C$	$O_A$	$O_B$	$O_C$	$O_A$	$O_B$	$O_C$
Response time per scenario (msec)	add blog	1316.00	1286.00	1237.00	1276.00	1275.00	1191.00	1379.00	1428.00	1366.00
	add event	1104.00	1073.00	1090.00	1033.00	1190.00	1133.00	1326.00	1409.00	1375.00
	edit blog	1077.00	1068.00	920.00	1025.00	1085.00	987.40	1080.00	1143.00	1100.00
	view blog	102.50	104.60	111.00	101.30	109.00	104.90	119.00	115.40	117.10
	view event	89.55	94.81	95.03	89.98	89.06	90.10	98.10	101.0	98.17
Resource usage	WS_CPU (%)	27.89	25.62	23.82	46.16	48.04	44.27	61.63	62.19	63.27
	WS_Memory (MB)	2532.00	2229.00	1848.00	3337.00	3770.00	3214.00	2998.00	3123.00	2944.00
	DB_CPU (%)	0.21	0.22	0.21	0.41	0.40	0.40	0.61	0.59	0.59
JIT	DB_Memory (MB)	87.49	87.85	86.96	90.58	93.90	95.09	97.84	97.90	95.20
	# of jitted lines	12987	9993	8780	11045	13065	12298	11002	11451	11402

**Table 16** Top three optimal configuration settings for Wagtail under different workloads

Workload	decay	function_threshold	loop_longevity	threshold	trace_eagerness	trace_limit
5 req/sec	2	404	250	129	12	12000
	40	101	250	519	12	12000
	20	404	500	1039	25	12000
10 req/sec	10	101	1000	1039	12	12000
	2	404	250	259	12	12000
	2	101	4000	4156	50	12000
15 req/sec	2	1619	125	3117	25	12000
	2	1619	62	2078	25	12000
	10	101	62	2078	25	12000
Default	1039	1619	40	6000	200	1000

one run under configuration A, and the other run under configuration B. After the profiling, we extracted the total execution time and the frequency of the executions for each function. Although the cProfile can provide us with function level profiling, it cannot provide information on which lines are executed during runtime. Hence, we implemented a simple tracer based on Python’s tracing library (Tracing a Program As It Runs 2018). Since both runs executed exactly the same workload, the lines of the executed source code should be the same. Hence, we only ran our tracer once. Finally, we parsed the jitted logs we collected for configuration A in order to know the exact lines of source code that were jitted.

Based on the cProfiling results, we calculated the average execution time for each function and computed their differences between the two configuration settings. We sorted the differences in decreasing order and selected top 30 functions whose performance is worse in configuration A for manual examination.



**Fig. 6** Number of jitted lines and overall average response time among all evaluated JIT configuration settings in Saleor. The red dotted line shows the overall average response time with JIT turned off

We found that the main reason behind the worse performance is configuration A is due to the overhead of time switching between the two execution modes (interpreted vs. native execution). When a function is executed each time, PyPy will be running under the default interpreted mode. When a code region is marked as jitted, PyPy will switch from the interpreted mode to the native execution mode and executed the compiled binary code. After the binary code is executed, PyPy will have to switch back to interpreted mode to execute the rest of the function. Executing the compiled binary code is much faster than running the same code under the interpreted mode. However, the switching between the two executing mode takes time. Figure 7 shows two such examples. The different text styles are defined as follows:

- **grey**: not executed;
- **bold**: executed but not jitted;
- **bold & highlighted**: jitted under configuration A.

As we can see from the Fig. 7, only a single line is jitted in both functions. Both lines are related to Python list comprehension, which internally execute a *for* loop. In both cases, PyPy has to switch the execution mode twice: starting from the interpreted mode to the native execution mode and back. The time saved under the native execution mode is much smaller than the time takes for switching between the two modes, which causes the performance degradation in configuration A (enabling JIT) when comparing against configuration B (JIT off).

Such jitting behavior can be explained using the configuration settings. The configuration A is shown in Table 17. For reference, we also included the default configuration values in the table. The ‘trace\_limit’ in configuration A is set to be a very small value, which would only allow a small region of code to traced and jitted. The smaller the jitted code region is, the less the performance gain code jitting can bring. In this case, the overhead of frequently switching between the two modes out-weights the gain from code jitting.

**Findings:** Enabling code jitting does not necessarily lead to good performance. Some JIT configuration settings can perform worse than the disabling the JIT completely. This is mainly due to the overhead of switching between the interpreted and the native execution mode.

**Implications:** Programming language researchers may look into adaptive JIT compilation techniques, which can disable inefficient code jitting behavior during runtime.

## 6.4 JIT vs. Memory Usage

The case studies have shown that by using our automated approach, we are able to locate JIT configuration settings whose performance significantly outperform the default configuration setting. The CPU for all the components are better or no worse in the optimal configuration settings than the default. This is mainly because the CPU can process the same amount of work much more efficiently when more code is compiled into efficient machine code. However, the memory usage for the Tornado web servers are much worse. The memory usage for the worker processes for Wagtail even tripled in the optimal configuration settings. Hence, we want to investigate whether there is any relation between the amount of code jitted and the amount of memory used in the worker processes.

For all the performance tests in each case study, we processed the JIT logs to obtain the amount of the jitted source code and extract the memory usage at the end of the test. Then

```

def patch_vary_headers(response, newheaders):
    """
    Adds (or updates) the "Vary" header in the given HttpResponse object.
    newheaders is a list of header names that should be in "Vary". Existing
    headers in "Vary" aren't removed.
    """
    # Note that we need to keep the original order intact, because cache
    # implementations may rely on the order of the Vary contents in, say,
    # computing an MD5 hash.
    if response.has_header('Vary'):
        vary_headers = cc_delim_re.split(response['Vary'])
    else:
        vary_headers = []
    # Use .lower() here so we treat headers as case-insensitive.
    existing_headers = set(header.lower() for header in vary_headers)
    additional_headers = [newheader for newheader in newheaders
                          if newheader.lower() not in existing_headers]
    response['Vary'] = ', '.join(vary_headers + additional_headers)

```

(a) Code snippet 1

```

def get_javascript_catalog(locale, domain, packages):
    app_configs = apps.get_app_configs()
    allowable_packages = set(app_config.name for app_config in app_configs)
    allowable_packages.update(DEFAULT_PACKAGES)
    packages = [p for p in packages if p in allowable_packages]
    paths = []
    # paths of requested packages
    for package in packages:
        p = importlib.import_module(package)
        path = os.path.join(os.path.dirname(upath(p.__file__)), 'locale')
        paths.append(path)

    trans = DjangoTranslation(locale, domain=domain, locale_dirs=paths)
    trans_cat = trans._catalog

```

(b) Code snippet 2

**Fig. 7** Two code snippets showing the executed code and the jitted code under the two configuration settings: A vs. B. Configuration A is a jit-enabled configuration shown in Fig. 6. It has worse performance than configuration B, which is JIT off. The colour scheme is defined as follows: grey coloured code is for not executed code; bolded black coloured code is for executed but not jitted code; and highlighted bold coloured code is for jitted code under configuration A

we calculated the Spearman's rank correlation between the lines of the jitted code and the amount of memory usage.

As shown in Table 18, there is a very large correlation between the memory usage of the web server processes and the amount of the jitted code. In other words, the larger the amount of the jitted code, the higher the memory usage for the worker processes. As more code is

**Table 17** Configuration A and the default configuration

JIT Configuration Parameter	A	Default
decay	5	40
function_threshold	404	1619
loop_longevity	1000	1000
threshold	4156	1039
trace_eagerness	100	200
trace_limit	374	6000

jitted, a larger amount of compiled machine code is generated. The generated machine code will be kept in the memory during the system execution.

**Findings:** The improvement in response time using the JIT compilation process is at the cost of higher memory usage.

**Implications:** More jitted code can generally lead to more responsive system. However, the number of jitted lines should be kept in a moderate range as: (1) more jitted code means higher memory usage, and (2) the configuration settings with the highest amount of jitted lines will not guarantee the best performance. One of the interesting future research work would be incorporating memory usage as one of the objectives in the ESM-MOGA while searching for the optimal configuration settings.

## 6.5 Termination Criteria

The above case studies show that among the three studied systems, all top three configuration settings significantly outperform the default configuration setting. We set ( $\mathbf{T}_0$ ;)  $|MDR| \leq 0.1$  in the hope that there is a higher chance to obtain the optimal configuration settings, as solutions in the previous generation are good enough so that little optimization can be made during the last generation before termination. However, such conditions may be too strict and cost too much time ( $\geq 35$  hours as shown in Table 11). Many systems nowadays need to be updated more frequently (e.g., daily or even a few times a day) under the continuous integration/continuous delivery processes. Hence, during the case studies, we also examined the following two termination conditions: (1) ( $\mathbf{T}_1$ ;)  $|MDR| \leq 0.25$ , and (2) ( $\mathbf{T}_2$ ;)  $|MDR| \leq 0.50$ , in the hope that the search process terminates earlier (a.k.a., saving the time for searching), while we are still able to locate the optimal solutions.

Table 19 shows the runtime statistics for ESM-MOGA under three different termination criteria:  $T_0$ ,  $T_1$ , and  $T_2$ . All termination criteria stopped at the same number of generations for Saleor and Quokka. However, for Wagtail, the less strict termination criteria  $T_1$  and

**Table 18** Spearman correlation between number of jitted line and memory usage for each system

System	correlation	<i>p</i> – value	scale
Saleor	0.8244878	2.2e-16	very large
Wagtail	0.882144	2.2e-16	very large
Quokka	0.8549971	2.2e-16	very large



**Table 19** Runtime statistics for ESM-MOGA under different termination criteria

Stoppage Criteria	Saleor			Wagtail			Quokka		
	$T_0$	$T_1$	$T_2$	$T_0$	$T_1$	$T_2$	$T_0$	$T_1$	$T_2$
# of Generation	3	3	3	7	3	3	3	3	3
Evaluated Configurations	100	100	100	199	97	97	96	96	96
Duration (hours)	36	36	36	67	34	34	35	35	35

$T_2$  are met after the third generation. Hence, we also extracted the top three configuration settings for Wagtail at the end of the third generation for further comparison.

We performed a pairwise comparison using the dominant comparison function between the top three configuration settings under termination criteria  $T_1$  and  $T_2$  (a.k.a., stopped after the third generation) and the top three configuration settings under termination criterion  $T_0$  (a.k.a., stopped after the seventh generation). The results show that two out of the three top configuration settings under  $T_0$  dominate all top three configuration settings under  $T_1$  and  $T_2$ . The other remaining top three configuration settings under  $T_0$  show no dominance when comparing against one of the top configuration settings under  $T_1$  and  $T_2$ . The system performance under the top three configuration settings under  $T_1$  and  $T_2$  also shows large improvement for all scenarios when comparing against the default setting.

**Findings:** We have evaluated the ESM-MOGA under three different termination criteria: (1)  $MDR \geq 0.50$ , (2)  $MDR \geq 0.25$ , and (3)  $MDR \geq 0.10$ . The ESM-MOGA can terminate successfully (a.k.a., finding the optimal solutions) under all three criteria for the three case study systems. And the less strict termination criteria  $T_1$  and  $T_2$  can obtain some configuration settings which are as good as some top configuration settings under termination criterion  $T_0$ .

**Implications:** There are various configuration parameters within ESM-MOGA. It requires further research to systemically tune ESM-MOGA configuration parameters in order to achieve the best performance (finding the optimal solutions within the shortest amount of time). Furthermore, it would be beneficial to compare ESM-MOGA against other hyper-parameter tuning techniques (e.g., Google Viser (Golovin et al. 2017)).

## 7 Threats to Validity

In this section, we will discuss the threats to validity.

### 7.1 Construct Validity

To avoid measurement errors and noise, we repeated each experiment 30 times (Georges et al. 2007).

We leveraged techniques from Alghmadi et al. (2016) to determine the duration of performance tests, when the performance behavior becomes repetitive. We have found that after 40 minutes under the default configuration settings for the three systems, the performance behavior becomes repetitive. For consistency concerns, we took additional 10 minutes of the performance data for our performance tests for the warmed up period. We assumed the

warmup period would be similar under other configuration settings. To verify this threshold, we randomly sampled five performance test from all performance testing runs in each case study. For each sampled test, we divided the performance data into intervals of ten minutes and compared the performance behavior among the adjacent intervals. Our analyses confirmed that the performance behavior also became repetitive after 40 minutes under these five sampled configuration settings.

Since the JIT logs do not contain timestamps, the only way to monitor the jitting progress for PyPy is to periodically take snapshots of the existing JIT logs. However, regularly taking snapshots of the JIT logs would bring huge performance overhead for a server-based system. Hence, to minimize the measurement impact, we did not take snapshots in the middle of the performance tests in our case study in Section 5. Instead, we estimated the jitting progress by judging whether the system performance behavior stabilizes (a.k.a., becoming repetitive).

Our approach, ESM-MOGA, is a tailed version of the Multi-Objective Genetic Algorithm, which uses effect size measures to compare the results of different test runs. MOGA is an efficient search-based technique, which automatically explores the solution space. It has been used widely in various software engineering research areas (e.g., test case generation (Abdessalem et al. 2018; Fraser and Arcuri 2011; Shamshiri et al. 2015), software architecture (Henard et al. 2015), and bug prediction (Canfora et al. 2013)) and is shown to be highly effective. Although our case study results show that the configurations derived from the ESM-MOGA approach yield much better performance than the default configuration, the ESM-MOGA might not be the most efficient approach to locate the optimal JIT configuration setting(s). Furthermore, the search time that it takes to find the optimal solutions using ESM-MOGA varies depending on the systems and their associated workload. One of the future areas of research is to evaluate the effectiveness of various hyper-configuration tuning techniques in the context of tuning JIT configuration parameters.

## 7.2 Internal Validity

We kept all the other factors (e.g., the versions of the systems, the deployment infrastructure, and the workload) the same, while varying the JIT configuration settings for each performance tests.

In this paper, we assumed the systems which undergo the JIT tuning process, can handle the exercised workload. In other words, the systems are not in a bottleneck state when we tune their JIT configurations. We feel this is a valid assumption, as the top priority for bottlenecked systems would be performance diagnosis and migration actions instead of tuning their JIT configuration settings.

We used the WRS test and the values from CD to implement our dominance functions in the ESM-MOGA. WRS is a statistical test which compares the distributions of two datasets. CD is an effect size measure, which indicates the strength of the difference between two datasets. Both the WRS test and the CD are non-parametric tests, which do not hold any assumptions regarding the underlying distributions of the data. The two techniques have been used together in previous works (Gao et al. 2016; Gao and Jiang 2017) to evaluate the system performance between two alternatives.

## 7.3 External Validity

In this paper, we have conducted a case study on the performance impact of the JIT configuration settings from PyPy. The experimented PyPy version is PyPy 5.7.1, which corresponds

to Python version 2.7.13. The empirical findings in the exploratory studies may not be generalizable to other Python versions (e.g., Python 3), other Python implementations (e.g., Jython), or other programming languages (e.g., Java or C#).

Our case study results have shown that the optimal JIT configuration settings vary from systems to systems. Our search-based configuration tuning framework can be used to automatically search for configuration settings, which are much better than the default. Our automated tuning technique can also be used to tune the JIT performance of other programming languages, whose parameters are integer types. We plan to extend our approach in the future to accommodate other types of configuration parameters (e.g., string, float and boolean types).

As each experiment requires starting the Python-based applications with a different set of configuration parameters, it is not yet practical to apply ESM-MOGA techniques into the continuous integration and continuous delivery process. One of the future areas of research is to look into techniques like transfer learning (Jamshidi et al. 2017b) to infer the optimal configurations for the newer releases of the same systems or even configuration parameters other systems.

## 8 Related Work

In this paper, we discuss two areas of research works which are related to this paper: (1) JIT compiler; and (2) techniques used to assess and optimize the configuration settings of a software system.

### 8.1 JIT Compiler

JIT is introduced as a technique to improve the system behavior during runtime by compiling the frequently executed (a.k.a., “hot”) code snippets into binaries (Bolz et al. 2009; Oaks 2014). Currently, there are two general approaches on recognizing and compiling hot code: (1) the method-based jitting approach (Cramer et al. 1997), which compiles the whole hot method; and (2) the trace-based jitting approach (Bolz et al. 2009), which only compiles the frequently executed code path(s) within one method. Both techniques have their pros and cons and are adopted by different programming languages. The code jitting process takes a while to recognize and compile the hot code snippets (Barrett et al. 2017). Hence, various techniques have been proposed to speed up the JIT compilations (Jantz and Kulkarni 2013; Gal et al. 2009; Lion et al. 2016). Since only portions of the source code are jitted, during runtime, depending on the actual execution path(es), the system may switch between the native mode (a.k.a., executing the compiled binaries) and the interpreter mode. Gong et al. (2015) developed a technique to detect performance anti-patterns that prohibit the system to execute certain portions of the code natively. Our paper differs from the above works, as it focuses on the configuration settings of the JIT compiler. The closest work was done by Hoste et al. (2010), which proposed a search-based technique to automatically tune the Java compiler. Although the two programming languages differ in their jitting techniques (method-based JIT for Java and tracing-based JIT for PyPy), both Hoste et al. (2010) and this paper reported the need to automatically tune JIT configurations, as the optimal JIT configurations are system and workload dependent. Hoste et al. (2010) even found that the JIT configuration tuning is hardware dependent. In this paper, we further studied the characteristics of the PyPy jitting behavior and tried to derive general patterns/guidelines on tuning the JIT configurations. For example, we have found a high correlation between the

amount of jitted code and memory utilization. Generally, the configuration parameter *decay* should be set with a small value and a large value in the *trace\_limit*.

## 8.2 Understanding and Tuning the Configuration Settings

Software configuration settings play an important role in the performance of a software system. However, there can be many configuration parameters, each of which has various possible settings. Hence, the overall configuration space for one system can be huge. In this subsection, we will discuss the related works in the area of assessing and optimizing the configuration settings for one system. We have divided the existing techniques into the following three categories:

- **Understanding the Performance Impact of Various Configuration Parameters Through Experimentation:** Not all configuration parameters can impact the system performance. Hence, researchers have devised a set of experiments with various combinations of the configuration settings to assess the impact of the configuration parameters (Brecht et al. 2006; Sopitkamol and Menascé 2005). Various experimental design techniques (e.g., screening design (Yilmaz et al. 2007), and covering array Hoskins et al. 2005) have been applied to assess the performance impact of various configuration parameters. These techniques require a much smaller set of experiments than exhaustively enumerating all the possible combinations of the configuration settings, while still able to identify the high performance impacting configuration parameters.
- **Modeling System Performance Under Different Configuration Settings:** Instead of isolating the impact of various configuration parameters, another approach to assess the performance impact of configuration settings is through performance modeling. Siegmund et al. (2012) predicted the system performance by detecting performance-relevant feature interactions. In their later work (Siegmund et al. 2015), Siegmund et al. leveraged machine learning and sampling heuristics to build performance models, which can describe the performance influences among different configuration options and their interactions, from a small set of experiments. Libič et al. (2014) used queuing theory to model the performance of the JVM garbage collectors (GCs). Singer et al. (2007) built machine learning models using the data from the existing configurations of the GCs. Recently, Jamshidi et al. (2017a) proposed to use the transfer learning technique to model and infer the system performance under different configuration settings.
- **Automated Tuning of the Configuration Settings:** There are two general approaches to automatically tuning the configuration settings of a software system: (1) through reduction of the possible candidates of optimal configurations (e.g., based on the similarities among configuration settings (Osogami and Kato 2007) or through iterative experimentations Thonangi et al. 2008); (2) through the use of the search-based algorithms (e.g., hill-climbing (Xi et al. 2004; Wang et al. 2012), ParamILS (Lengauer and Mössenböck 2014), or multi-objective genetic algorithms Hoste et al. 2010; Singh et al. 2016).

In this paper, we assessed the performance impact of PyPy's JIT configuration settings through experimentation. Furthermore, we have provided some insights into why certain code regions are jitted or not jitted. For locating the optimal JIT configuration settings for each system, we adopted the multi-objective genetic algorithms. It can be an interesting future work to compare the effectiveness and efficiency in terms of tuning the JIT configuration parameters among the various approaches described above.

## 9 Conclusion and Future Work

The JIT compilation is introduced to improve the runtime performance of software systems. During the system execution, various regions of the systems are compiled into the binary executable format, so that they can be executed more efficiently. In this paper, we have performed an empirical study on the performance impact of PyPy's JIT configuration settings. In particular, we have compared the performance differences between the default and some other configuration settings. We have shown that systems running under the default configuration setting does not necessarily yield the best performance. In addition, we have shown that there is no strong connection between the JIT coverage and the system performance and the optimal JIT configuration settings are system dependent. To cope with such findings, we have developed a search-based approach, called ESM-MOGA, which automatically tunes the JIT configuration settings for a given system. Case studies have shown that systems running under the resulting configuration settings are significantly faster than the default configuration setting.

In the future, we plan to further expand the ESM-MOGA to accommodate more objectives (e.g., memory and network efficiency) during its search process. We also would like to apply the ESM-MOGA on other Python-based applications or frameworks (e.g., machine learning based frameworks like TensorFlow). In addition, we plan to explore the use of data mining or experimental design techniques to further reduce the number of performance tests conducted during the search process. Finally, we would like to evaluate the effectiveness of various hyper-parameter tuning techniques in the context of tuning JIT configurations.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

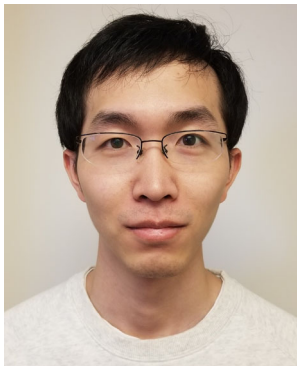
## References

- Abdessalem RB, Panichella A, Nejati S, Briand LC, Stifter T (2018) Testing autonomous cars for feature interaction failures using many-objective search. In: Proceedings of the 33rd ACM/IEEE international conference on automated software engineering (ASE)
- Alghamdi HM, Syer MD, Shang W, Hassan AE (2016) An automated approach for recommending when to stop performance tests. In: 2016 IEEE international conference on software maintenance and evolution (ICSME), pp 279–289
- Apache JMeter (2015) <http://jmeter.apache.org/>, visited 2015-10-23
- Barrett E, Bolz-Tereick CF, Killick R, Mount S, Tratt L (2017) Virtual machine warmup blows hot and cold. In: Proceedings of the ACM Programming Language 1(OOPSLA), pp 52:1–52:27. <https://doi.org/10.1145/3133876>
- Bolz CF, Cuni A, Fijalkowski M, Rigo A (2009) Tracing the meta-level: Pypy's tracing jit compiler. In: Proceedings of the 4th workshop on the implementation, compilation, optimization of object-oriented languages and programming systems (ICOOOLPS), pp 18–25
- Bondi AB (2007) Automating the analysis of load test results to assess the scalability and stability of a component. In: Proceedings of the 2007 computer measurement group conference (CMG), pp 133–146
- Brecht T, Arjomandi E, Li C, Pham H (2006) Controlling garbage collection and heap growth to reduce the execution time of java applications. *ACM Trans Program Lang Syst* 28(5):908–941. <https://doi.org/10.1145/1152649.1152652>
- Candan KS, Li WS, Luo Q, Hsiung WP, Agrawal D (2001) Enabling dynamic content caching for database-driven web sites. In: Proceedings of the 2001 ACM SIGMOD international conference on management of data (SIGMOD), pp 532–543
- Canfora G, Lucia AD, Penta MD, Oliveto R, Panichella A, Panichella S (2013) Multi-objective cross-project defect prediction. In: Proceedings of the 2013 IEEE 6th international conference on software testing, verification and validation (ICST)

- Clark M (2017) How the BBC builds websites that scale. <http://www.creativebloq.com/features/how-the-bbc-builds-websites-that-scale>. Last accessed 10/06/2017
- Cramer T, Friedman R, Miller T, Seberger D, Wilson R, Wolczko M (1997) Compiling java just in time. *IEEE Micro* 17(3):36–43
- Deb K, Agrawal S, Pratap A, Meyarivan T (2000) A fast elitist non-dominated sorting genetic algorithm for multi-objective optimization: Nsga-ii. In: International conference on parallel problem solving from nature. Springer, pp 849–858
- Distributed Evolutionary Algorithms in Python (DEAP) (2017) <https://github.com/DEAP/deap>. Last accessed 10/06/2017
- Duan S, Thummala V, Babu S (2009) Tuning Database Configuration Parameters with iTuned. *Proceedings of the VLDB Endowment* 2(1):1246–1257. <https://doi.org/10.14778/1687627.1687767>
- Eaton K (2017) How One Second Could Cost Amazon \$1.6 Billion In Sales. <https://www.fastcompany.com/1825005/how-one-second-could-cost-amazon-16-billion-sales>. Last accessed 10/06/2017
- Evaluation tools in Python (DEAP) (2017) <http://deap.readthedocs.io/en/master/api/tools.html?highlight=dominance>. Last accessed 10/06/2017
- Fraser G, Arcuri A (2011) Evosuite: automatic test suite generation for object-oriented software. In: Proceedings of the 19th ACM SIGSOFT symposium and the 13th european conference on foundations of software engineering (ESEC/FSE)
- Gal A, Eich B, Shaver M, Anderson D, Mandelin D, Haghghat MR, Kaplan B, Hoare G, Zbarsky B, Orendorff J, Ruderman J, Smith EW, Reitmaier R, Bebenita M, Chang M, Franz M (2009) Trace-based just-in-time type specialization for dynamic languages. In: Proceedings of the 30th ACM SIGPLAN conference on programming language design and implementation (PLDI), pp 465–478
- Gao R, Jiang ZMJ (2017) An exploratory study on assessing the impact of environment variations on the results of load tests. In: Proceedings of the 14th international conference on mining software repositories (MSR)
- Gao R, Jiang ZMJ, Barna C, Litoiu M (2016) A framework to evaluate the effectiveness of different load testing analysis techniques. In: 2016 IEEE international conference on software testing, verification and validation (ICST)
- Georges A, Buytaert D, Eeckhout L (2007) Statistically rigorous java performance evaluation. In: Proceedings of the 22nd international conference on object-oriented programming, systems, languages and applications (OOPSLA)
- Gewirtz D (2017) Which programming languages are most popular (and what does that even mean)? <http://www.zdnet.com/article/which-programming-languages-are-most-popular-and-what-does-that-even-mean/>. Last accessed 10/06/2017
- Golovin D, Solnik B, Moitra S, Kochanski G, Karro J, Sculley D (2017) Google vizier: a service for black-box optimization. In: Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining (KDD)
- Gong L, Pradel M, Sen K (2015) Jitprof: Pinpointing jit-unfriendly javascript code. In: Proceedings of the 2015 10th joint meeting on foundations of software engineering (ESEC/FSE), pp 357–368
- Grigorik I (2017) Optimizing Encoding and Transfer Size of Text-Based Assets. <https://developers.google.com/web/fundamentals/performance/optimizing-content-efficiency/optimize-encoding-and-transfer>. Last accessed 10/06/2017
- Hashemi M (2014) 10 Myths of Enterprise Python. <https://www.paypal-engineering.com/2014/12/10/10-myths-of-enterprise-python/>. Last accessed 10/06/2017
- Henard C, Papadakis M, Harman M, Traon YL (2015) Combining multi-objective search and constraint solving for configuring large software product lines. In: Proceedings of the 37th international conference on software engineering (ICSE)
- Hopkins WG (2016) A new view of statistics. [Online accessed 2017-10-14] <http://www.sportsci.org/resource/stats/index.html>
- Hoskins DS, Colbourn CJ, Montgomery DC (2005) Software performance testing using covering arrays: Efficient screening designs with categorical factors. In: Proceedings of the 5th international workshop on software and performance (WOSP)
- Hoste K, Georges A, Eeckhout L (2010) Automated just-in-time compiler tuning. In: Proceedings of the 8th annual IEEE/ACM international symposium on code generation and optimization (CGO), pp 62–72
- IBM Java 8 JIT and AOT command-line options (2017) [https://www.ibm.com/support/knowledgecenter/SSYKE2\\_8.0.0/com.ibm.java.aix.80.doc/diag/appendixes/cmdline/commands\\_jit.html](https://www.ibm.com/support/knowledgecenter/SSYKE2_8.0.0/com.ibm.java.aix.80.doc/diag/appendixes/cmdline/commands_jit.html). Last accessed 10/06/2017
- Insights GP (2017) Remove Render-Blocking JavaScript. <https://developers.google.com/speed/docs/insights/BlockingJS>. Last accessed 10/06/2017

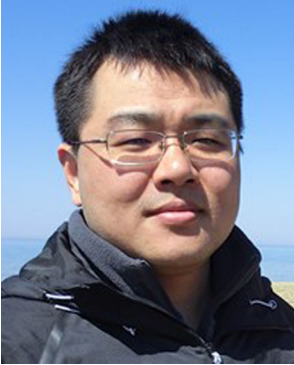
- Jamshidi P, Siegmund N, Velez M, Kästner C, Patel A, Agarwal Y (2017) Transfer learning for performance modeling of configurable systems: an exploratory analysis. In: Proceedings of the international conference on automated software engineering (ASE)
- Jamshidi P, Siegmund N, Velez M, Kästner C, Patel A, Agarwal Y (2017) Transfer learning for performance modeling of configurable systems: an exploratory analysis. In: Proceedings of the 32nd IEEE/ACM international conference on automated software engineering (ASE)
- Jantz MR, Kulkarni PA (2013) Exploring single and multilevel jit compilation policy for modern machines 1. *ACM Trans Archit Code Optim (TACO)* 10(4):22:1–22:29
- Java Microbenchmark Harness (JMH) (2017) <http://openjdk.java.net/projects/code-tools/jmh/>. Last accessed 10/06/2017
- Jiang ZM, Hassan AE (2015) A survey on load testing of large-scale software systems. *IEEE Trans Softw Eng* 41:1-1. <https://doi.org/10.1109/TSE.2015.2445340>
- Kampenes VB, Dybå T, Hannay JE, Sjøberg DIK (2007) Systematic review: A systematic review of effect size in software engineering experiments. *Inf Softw Technol* 49(11-12):1073–1086
- Komorn R (2016) Python in production engineering. <https://code.facebook.com/posts/1040181199381023/python-in-production-engineering/>. Last accessed 10/06/2017
- Lengauer P, Mössenböck H (2014) The taming of the Shrew: increasing performance by automatic parameter tuning for java garbage collectors. In: Proceedings of the 5th ACM/SPEC international conference on performance engineering (ICPE), pp 111–122
- Libič P, Bulej L, Horký V, Tůma P (2014) On the limits of modeling generational garbage collector performance. In: Proceedings of the 5th ACM/SPEC international conference on performance engineering (ICPE)
- Lion D, Chiu A, Sun H, Zhuang X, Grcevski N, Yuan D (2016) Don't get caught in the cold, warm-up your jvm: Understand and eliminate jvm warm-up overhead in data-parallel systems. In: Proceedings of the 12th USENIX conference on operating systems design and implementation (OSDI), pp 383–400
- Martí L, García J, Berlanga A, Molina JM (2009) An approach to stopping criteria for multi-objective optimization evolutionary algorithms: The mgbm criterion. In: IEEE congress on evolutionary computation, 2009. CEC'09. IEEE, pp 1263–1270
- Oaks S (2014) Java performance: the definitive guide, 1st. O'Reilly Media, Inc, Sebastopol
- Osogami T, Kato S (2007) Optimizing system configurations quickly by guessing at the performance. In: Proceedings of the 2007 ACM SIGMETRICS international conference on measurement and modeling of computer systems (SIGMETRICS)
- Oracle Java 8 Advanced JIT Compiler Options (2017) <https://docs.oracle.com/javase/8/docs/technotes/tools/windows/java.html#BABDDFFII>. Last accessed 10/06/2017
- Performance monitoring tools for Linux (2015) <https://github.com/sysstat/sysstat>, visited 2015-10-23
- PyPy speed center (2017) <http://speed.pypy.org/>. Last accessed 10/04/2017
- Quokka CMS (Content Management System) - Python Flask and MongoDB (2017) <http://quokkaproject.org/>. Last accessed 10/06/2017
- Replication package (2018) <https://github.com/seasun525/PyPyJITtuner>
- Romano J, Kromrey JD, Coraggio J, Skowronek J (2006) Appropriate statistics for ordinal level data: Should we really be using t-test and Cohen's d for evaluating group differences on the NSSE and other surveys? In: Annual meeting of the Florida Association of Institutional Research
- Saleor - An e-commerce storefront for Python and Django (2017) <http://getsaleor.com/>. Last accessed 10/06/2017
- Shamshiri S, Rojas JM, Fraser G, McMinn P (2015) Random or genetic algorithm search for object-oriented test suite generation? In: Proceedings of the 2015 annual conference on genetic and evolutionary computation (GECCO)
- Siegmund N, Grebhahn A, Apel S, Kästner C (2015) Performance-influence models for highly configurable systems. In: Proceedings of the 2015 10th joint meeting on foundations of software engineering, ESEC/FSE 2015. ACM
- Siegmund N, Kolesnikov SS, Kästner C, Apel S, Batory D, Rosenmüller M, Saake G (2012) Predicting Performance via Automated Feature-Interaction Detection. In: Proceedings of the 34th international conference on software engineering (ICSE)
- Singer J, Brown G, Watson I, Cavazos J (2007) Intelligent selection of application-specific garbage collectors. In: Proceedings of the 6th International Symposium on Memory Management, ISMM '07
- Singh R, Bezemer CP, Shang W, Hassan AE (2016) Optimizing the performance-related configurations of object-relational mapping frameworks using a multi-objective genetic algorithm. In: Proceedings of the 7th ACM/SPEC international conference on performance engineering (ICPE), pp 309–320

- Sopitkamol M, Menascé DA (2005) A method for evaluating the impact of software configuration parameters on e-commerce sites. In: Proceedings of the 5th international workshop on software and performance (WOSP), pp 53–64
- TechEmpower Web Framework Benchmarks (2017) <https://www.techempower.com/benchmarks/>. Last accessed 10/04/2017
- The Python Profilers (2018) <https://docs.python.org/2/library/profile.html>. Last accessed 10/28/2018
- Thonangi R, Thummala V, Babu S (2008) finding good configurations in High-Dimensional spaces: doing more with less. In: 2008 IEEE international symposium on modeling, analysis and simulation of computers and telecommunication systems
- Tracing a Program As It Runs (2018) <https://pymotw.com/2/sys/tracing.html>. Last accessed 10/28/2018
- vmprof - a statistical program profiler (2017) <http://vmprof.com/>. Last accessed 10/06/2017
- Wang K, Lin X, Tang W (2012) Predator - An experience guided configuration optimizer for Hadoop MapReduce. In: 4Th IEEE international conference on cloud computing technology and science proceedings
- Wagtail CMS: Django Content Management System (2017) <https://wagtail.io/>. Last accessed 10/06/2017
- What is Load Balancing? (2017) <https://www.nginx.com/resources/glossary/load-balancing/>. Last accessed 10/06/2017
- What is python used for at Google? (2017) <https://www.quora.com/What-is-python-used-for-at-Google>. Last accessed 10/06/2017
- Wimmer C, Brunthaler S (2013) Zippy on truffler: a fast and simple implementation of python. In: Proceedings of the 2013 companion publication for conference on systems, programming, & applications: software for humanity (SPLASH)
- Würthinger T, Wimmer C, Humer C, Wöß A, Stadler L, Seaton C, Duboscq G, Simon D, Grimmer M (2017) Practical partial evaluation for high-performance dynamic language runtimes. In: Proceedings of the 38th ACM SIGPLAN conference on programming language design and implementation, PLDI 2017. ACM, New York, pp 662–676. <https://doi.org/10.1145/3062341.3062381>
- Würthinger T, Wimmer C, Wöß A, Stadler L, Duboscq G, Humer C, Richards G, Simon D, Wolczko M (2013) One vm to rule them all. In: Proceedings of the 2013 ACM international symposium on new ideas, new paradigms, and reflections on programming & software (Onward!)
- Xi B, Liu Z, Raghavachari M, Xia CH, Zhang L (2004) A smart hill-climbing algorithm for application server configuration. In: Proceedings of the 13th international conference on world wide web (WWW). ACM, New York, pp 287–296. <https://doi.org/10.1145/988672.988711>
- Xu T, Jin X, Huang P, Zhou Y, Lu S, Jin L, Pasupathy S (2016) Early detection of configuration errors to reduce failure damage. In: Proceedings of the 12th USENIX conference on operating systems design and implementation (OSDI)
- Yilmaz C, Porter A, Krishna AS, Memon AM, Schmidt DC, Gokhale AS, Natarajan B (2007) Reliable effects screening: a distributed continuous quality assurance process for monitoring performance degradation in evolving software systems. *IEEE Trans Softw Eng (TSE)* 33(2):124–141. <https://doi.org/10.1109/TSE.2007.20>



**Yangguang Li** is a graduate student at the Department of Electrical Engineering and Computer Science, York University in Toronto, ON, Canada. He received his Bachelor of Engineering degree from the School of Computer Science and Technology at Beijing University of Posts and Telecommunications in Beijing, China. His research interests are software performance engineering, mining software repositories, and debugging and monitoring of distributed systems.





**Zhen Ming (Jack) Jiang** is an associate professor at the Department of Electrical Engineering and Computer Science, York University in Toronto, Canada. His research interests lie within software engineering and computer systems, with special interests in software performance engineering, software analytics, source code analysis, software architectural recovery, software visualizations, and debugging and monitoring of distributed systems. Some of his research results are already adopted and used in practice on a daily basis. He is the cofounder of the annually held International Workshop on Load Testing and Benchmarking of Software Systems (LTB). He received several Best Paper Awards including ICST 2016, ICSE 2015 (SEIP track), ICSE 2013, and WCRE 2011. He received the BMath and MMath degrees in computer science from the University of Waterloo, and the PhD degree from the School of Computing at the Queen's University.

## Affiliations

Yanguang Li<sup>1</sup>  · Zhen Ming (Jack) Jiang<sup>1</sup>

Zhen Ming (Jack) Jiang  
zmjiang@cse.yorku.ca

<sup>1</sup> Software Construction, AnaLytics and Evaluation (SCALE) lab, York University, Toronto, ON, Canada