

A Systematic Mapping Study on Architectural Approaches to Software Performance Analysis

Yutong Zhao, Lu Xiao, Chenhao Wei, Rick Kazman, and Ye Yang

Abstract—Software architecture is the foundation of a system’s ability to achieve various quality attributes, including software performance. However, there lacks comprehensive and in-depth understanding of why and how software architecture and performance analysis are integrated to guide related future research. To fill this gap, this paper presents a systematic mapping study of 109 papers that integrate software architecture and performance analysis. We focused on five research questions that provide guidance for researchers and practitioners to gain an in-depth understanding of this research area. These questions addressed: a systematic mapping of related studies based on the high-level research purposes and specific focuses (RQ1), the software development activities these studies intended to facilitate (RQ2), the typical study templates of different research purposes (RQ3), the available tools and instruments for automating the analysis (RQ4), and the evaluation methodology employed in the studies (RQ5). Through these research questions, we also identified critical research gaps and future directions, including: 1) the lack of available tools and benchmark datasets to support replication, cross-validation and comparison of studies; 2) the need for architecture and performance analysis techniques that handle the challenges in emerging software domains; 3) the lack of consideration of practical factors that impact the adoption of the architecture and performance analysis approaches; and finally 4) the need for the adoption of modern ML/AI techniques to efficiently integrate architecture and performance analysis.

Index Terms—Software architecture, Software performance, Secondary study, Systematic mapping study



1 INTRODUCTION

Software architecture is defined as the set of software elements, their properties, and the relationships among them [1]. It is the foundation of a system’s ability to achieve various quality attributes [2]. Software performance is, for many systems, the most important quality attribute driving the design [3], [4]. Performance is the ability of a software system to perform its duties according to time constraints within its allowance of resources [3]. Poor performance can result in long execution times, unhappy users, and even system crashes [3], [4]. Performance problems may be introduced and incubate in the system in early software design and architecture decisions [5]. These decisions may lead to severe performance compromises in the system and may require significant effort to detect and fix later, often during software maintenance. Thus, architects should consciously address performance design considerations in their earliest software architecture design activities.

Due to the intrinsic connections between software architecture and performance, there is a substantial body of literature that integrates architecture and performance analysis [6], [7], [8], [9], [10], [11], [12]. The most thoroughly studied aspect is model-based performance prediction. That is, practitioners leverage software architecture models to predict the performance measures of a system in the early stages of the software development life cycle. A number of secondary studies have already summarized the state-of-the-art in this direction [7], [8], [9]. However, it remains unclear what are the other pillars of software architecture modeling and analysis that can facilitate software performance in practice. In particular, we are interested to study if there are methodologies that leverage

architectural information to address software performance in the later stages of software development life-cycle, such as during software maintenance.

The high-level objective of this study is to help researchers and practitioners gain a comprehensive and in-depth understanding of why and how software architecture and performance analysis are integrated in the current literature. Toward this objective, we searched and retrieved a total of 24901 potentially relevant studies, out of which 109 studies are retained in our final dataset. The process is guided by the well-established systematic mapping study process [13], [14]. We provide a landscape of the current state in the integration of software architecture and performance analysis. More specifically, provided answers to three research questions:

RQ1: What are the different research purposes and facilitated development activities of integrating software architecture and performance analysis? The 109 studies that integrate software architecture and performance analysis serve four objectives 1) model-based performance analysis (74 studies); 2) performance anti-pattern detection and resolution (16 studies); 3) profiling and comparison of architectural alternatives (11 studies); and 4) self-adaptive architecture for dynamic performance optimization (8 studies). We also provide a mapping of software development activities, including *Implementation*, *Deployment*, *Operation*, and *Maintenance*, with these four research purposes. The mapping provides guidance if one would like to identify the most relevant technique for their focused activity. It also highlights gaps that could be filled in future research.

RQ2: How are software architecture and software performance analysis integrated for different research purposes? We summarized the typical study templates of the four fundamental purposes from the 109 studies, which offer

a guide for reproducing existing studies and advancing the state-of-the-art. Guidance regarding the key analysis components in the templates, including the architecture models, model annotation, performance models, and anti-pattern detection rules are provided.

RQ3: *To what degree is the architecture performance analysis automated by available tools and instruments?* We provided a catalog of tools for architecture analysis and for performance profiling curated from the 109 studies. This can help researchers and practitioners identify resources for their own research, or for reproduction studies.

Finally, we discussed the limitations of existing research, from which a tentative future research roadmap with related research questions is shared for researchers and practitioners to refer to and build upon, including: 1) how to enhance the reproducibility and interoperability of existing research based on the techniques, tools, and datasets analyzed from existing studies? 2) How well do existing software architecture and performance analysis techniques apply to emerging domains, such as Web 3.0 and VR? What are the unique challenges with these new domains in the direction of architecture and performance analysis? 3) What are the most significant challenges and impacting factors on the practical usage of architecture and performance analysis techniques? And, 4) how to enable AI in architecture and performance analysis, in particular for mitigating the architecture analysis's reliance on experts and conquering the complexity and uncertainty of performance analysis?

In summary, the key contributions of this study include:

- A comprehensive overview of the current state of research that integrates software architecture and performance analysis.
- A summary of available tools and instruments that support the integration of software architecture and performance research, including tools for constructing and analyzing architecture models and performance models, as well as instruments for collecting performance metrics.
- A summary of limitations of the prior work, which motivates future directions that need more attention.

The rest of this paper is organized as follows: Section 2 introduces the background of software architecture and software performance; Section 3 illustrates our methodology for preparing, retrieving, selecting, annotating, and synthesizing the relevant literature; Section 4.1 introduces the overall statistics of our dataset; Section 4 presents the study results; Section 5 discusses the potential future research directions; Section 6 discuss similar secondary studies; and Section 7 concludes.

2 BACKGROUND

In this section, we introduce background information about software architecture and software performance.

2.1 Software Architecture

2.1.1 Software Architecture Definition

Bass *et al.* define software architecture as “*the structure or structures of the system, which comprise software elements,*

the externally visible properties of those elements, and the relationships among them.” [1]. When describing software architecture, practitioners and researchers often treat software components as the composing elements [5]. A software component is a modular, portable, replaceable, and reusable set of well-defined functionality that encapsulates its implementation, and the component is exported as a high-level interface [5]. The externally visible properties of an architectural element refer to the assumptions that other elements can make of it, such as services it provide, its performance characteristics, shared resource usage, and so on [15]. The architecture of a software system serves as a blueprint of the system, which involves the high level structure of software system abstraction [5]. The design of software architecture must be compatible with the functionality of the system, but is driven by its quality attribute requirements such as performance, reliability, scalability, availability, and flexibility, etc [5].

2.1.2 Software Architecture Models

Different stakeholders view software architecture differently, as discussed in the classic “4+1” view [16]. To describe software architecture, practitioners use a variety of models and their representations, each serving distinct purposes. For instance, the Unified Modeling Language (UML) is an object-oriented language primarily used to visualize, specify, construct, and document the high-level structure and behavior of a software system [17], [18]. Architecture Description Languages (ADLs) are another class of models that provide precise syntactic and semantic descriptions for software architecture [19]. They support the formal specification of software features such as processes, threads, data, and sub-programs, as well as hardware components like processors, devices, buses, and memory. Other software architecture models include the Palladio Component Model, Message Sequence Charts, and the Descartes Modeling Language [S1], [S2], [S3], [S4]. The Palladio Component Model aims to predict performance, reliability, and maintainability of software architectures in early development stages, while Message Sequence Charts display interactions between processes or objects in real-time systems, and the Descartes Modeling Language is used for performance prediction and system quality evaluation for self-adaptive software systems.

2.2 Software Performance

2.2.1 Software Performance Engineering

According to Smith *et al.*, performance concerns how well a software system or its components meets the requirements for timeliness [20]. Timeliness encompasses two important dimensions: response time and throughput [21], [22]. Response time is the time required to respond to events, while throughput is the number of events processed during a time interval [22]. Zaman *et al.* further extends the performance concerns by considering resource utilization as another aspect of software performance [3], [23]. The resources of interest usually include: 1) hardware resources, such as CPU, disk I/O, and memory; 2) logical resources, such as buffers, locks, and semaphores; and 3) processing resources, such as threads and processes [3], [4], [24],

[25]. Moreover, recent studies also address other system characteristics, such as the data transmission loss ratio [S5], [S6], energy consumption [S7], and network latency [S8].

To ensure these performance concerns are met throughout the development process, practitioners must incorporate performance considerations from the outset, rather than deferring them until testing—a practice referred to as the “fix-it-later” approach. To alleviate the problems caused by this delayed approach, Smith *et al.* proposed *Software Performance Engineering* [19], a systematic, quantitative approach initiated in the early stages of the software development life cycle, and continuing through architecture design, implementation, and maintenance [25]. This performance engineering life cycle mirrors the software development life cycle [26]. It emphasizes the analysis of requirements regarding performance, scalability, stability, and reliability during requirement gathering and analysis [25]. During the design phase, it deploys modelling techniques like *Queueing Networks*, *Place/Transition (Petri) Nets*, and *Stochastic Process Algebra* [9] to validate design assumptions concerning performance. When it comes to implementation, this approach provides guidelines for validating and reporting on the performance of the code [25]. Lastly, in the maintenance phase, it focuses on monitoring the performance metrics of a system, thoroughly evaluating their compliance with the performance requirements, and proposing potential optimization recommendations [20], [S13].

2.2.2 Software Performance Testing and Profiling

Performance testing is one of the most thoroughly studied approaches to addressing performance concerns. The purpose of performance testing is to identify bottlenecks in software systems [27]. Performance testing executes a system and constructs a profile of it, in terms of responsiveness and stability under various workloads [20]. Avritzer and Weyuker made notable contributions to this area with their work on the automatic generation of load test suites and the assessment of the resulting software [28]. There are four major types of performance testing methodologies: load testing, stress testing, endurance testing, and spike testing [1], [4]. Load testing evaluates the behavior of a software system under specific workloads. Stress testing executes and profiles the system under extreme workloads to discover the maximum capacity of the system. Endurance testing executes and profiles the system under continuous workload. The purpose is to determine whether the system can scale up to support enduring and increasing workloads. Spike testing determines whether a system can sustain a sudden increase in workload [25].

During performance testing, practitioners leverage profiling tools to keep track of metrics such as response time, throughput, and resource utilization [29]. Avritzer and Weyuker also investigated metrics for architectural assessment in performance testing, enhancing the understanding of the metrics necessary for assessing software architecture [30]. These tools are available for many platforms, programming languages, and execution environments, with different advantages [31], [32]. For example, *WebLoad* can generate real-life and reliable workload scenarios for testing complex systems [33].

LoadNinja has the highest coverage for performance testing [34]. *LoadView* can be applied to real-life browsers and web applications [35]. *StresStimulus* can detect hidden concurrency errors by measuring performance metrics, such as network latency and data transmission loss ratio [36]. *Apache JMeter* [37] is the most widely-used performance testing and profiling tool for *Java* projects and it has been integrated to many IDEs. *SmartMeter* can automatically generate a performance assessment report [38]. *Rational Performance Tester* is a powerful performance testing tool developed by *IBM* [39], which supports load testing that involves multiple users and generates a comprehensive performance assessment report.

3 STUDY PROCESS

This study followed the guidelines for *systematic literature review* studies in software engineering [13], [14]. We followed the process suggested by Kitchenham *et al.* [14] with study planning (i.e. defining the RQs), literature searching, study selection, extracting study data, and data aggregation and synthesis. In addition, we followed the guidelines in [13] for ensuring stability and reliability. We recorded detailed data from each step, including initial search results, notes for the selection process, and the data annotation and synthesis, to maintain a clear chain of evidence for our findings. The detailed data is available on a public repository: <https://sites.google.com/view/arch-perf-data-repo/>

Figure 1 illustrates the five steps in this study:

- 1) **Study Preparation:** This step defined the research questions.
- 2) **Study Retrieval:** We defined the search strings and retrieved the initial set of papers from major digital libraries.
- 3) **Study Selection:** We conducted three rounds of selections to identify the most relevant literature.
- 4) **Snowballing Search:** We further examined the references of the literature retrieved from Step 2 to identify additional studies.
- 5) **Data Annotation and Synthesis:** We read each paper in the final set, annotated information on each paper that was relevant to our research questions, and synthesized the results to answer the research questions.

3.1 Study Preparation: Defining Research Questions

The first step was defining the research questions, which guided the rest of the study.

RQ1: What are the different research purposes facilitated development activities of integrating software architecture and performance analysis? Model-based performance prediction, leverages software architecture models for predicting performance metrics of a system, is the most well-known motivation for integrating architecture and performance analysis. However, it remains unclear whether exist and what are the other important themes that motivate the integration of architecture and performance analysis. In addition, this RQ investigate the development activities that related studies target at and facilitate. A prior study [9] reported that the model-based prediction techniques were

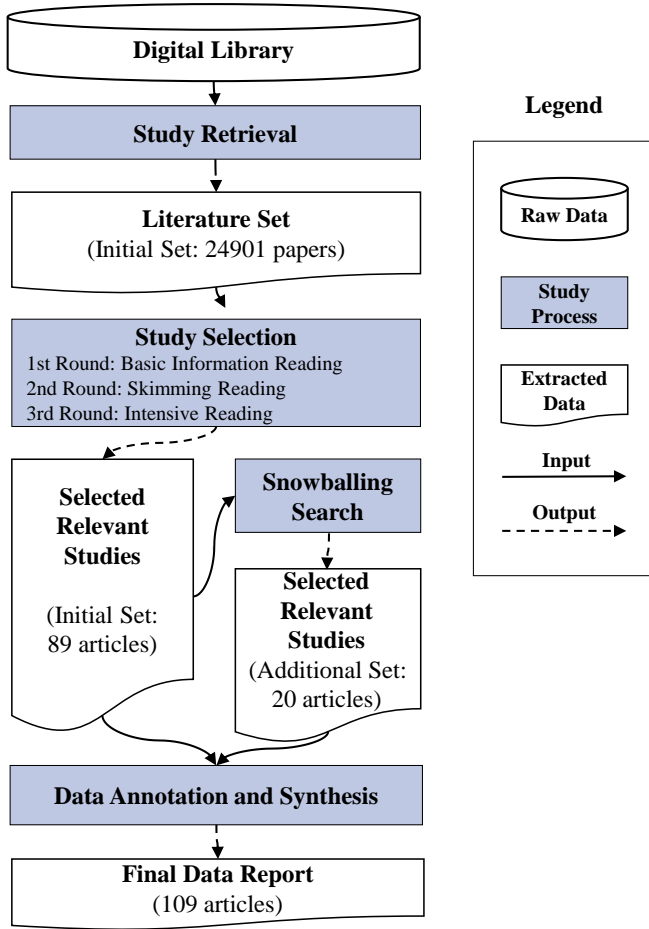


Fig. 1: Overview of Study Process

mostly applied during the software architecture design. This provides early performance assessment and enables optimization of the system architecture design before the implementation starts. However, it was not clear whether the different research purposes usually maps to different development activities, especially activities other than software architecture design reported in the prior study. The motivation of this RQ is to provide practitioners a comprehensive categorization and mapping of existing studies based on their objectives and the specific development activities that they concern. As such, this mapping can guide future researchers and practitioners to the studies that best fit their interests.

RQ2: How are software architecture and software performance analysis integrated for different research purposes? This RQ aims to address “how” software architecture and performance analysis are integrated for different objectives. We aim to provide reference research templates of the different research purposes identified in RQ1. Each template should summarize the common study steps and their workflow from papers of the same purpose.

RQ3: To what degree is the architecture performance analysis automated by available tools and instruments? This RQ aims to summarize tools and instruments that are used to automate the architecture and performance analysis. In particular, we aim to reveal what is the extent of manual effort reported in the literature. For summarizing

the tools, we present two parts: 1) RQ-3.1: What are the available architecture and performance modeling and analysis tools? and 2) RQ-3.2: What are the performance profiling instruments for collecting dynamic performance metrics, and which metrics are collected? This RQ helps practitioners to quickly identify relevant tools or techniques for their purposes.

3.2 Study Retrieval

Our search string was defined as follows:

(software) AND (architecture OR architectural OR architecting OR design OR structure) AND (performance)

The first keyword, “software”, establishes the overall focus on the field of software engineering. Originally, we included terms such as “architecture”, “architectural”, “architecting”, “design”, “structure”, and “behavior” to encompass various aspects of software architecture.

The first keyword, “software”, establishes the overall focus on the field of software engineering. Originally, we included terms such as “architecture”, “architectural”, “architecting”, “design”, “structure”, and “behavior” to encompass various aspects of software architecture. However, introducing “behavior” retrieved an overabundance of unrelated papers primarily covering human behavior studies, hence, it was excluded to maintain specificity. Further exploration of architecture-related terms, including “component-based system”, “distributed system”, and “service-oriented system”, did increase the total number of papers, but our analysis revealed that all these papers were already included within the scope of the terms “architecture”, “design”, or “structure”. Therefore, these additional terms were deemed redundant and excluded. This streamlined search string demonstrated its effectiveness and efficiency in our subsequent experiment, where it successfully retrieved relevant papers without introducing unnecessary noise or omitting key studies. However, the term “behavior” could also represent the broader human behavior related topics and led to large noise. We removed “behavior” out of the final search string.

We conducted the search on seven digital libraries: 1) *IEEE Xplore*, 2) *ACM Digital Library*, 3) *Web of Science*, 4) *Springer Digital Library*, 5) *Elsevier ScienceDirect*, 6) *Wiley Online Library*, and 7) *Google Scholar*. For each library, we searched the document title, keywords, and abstract, to comprehensively retrieve potential relevant studies. The search range was January 2010 to December 2020. The guideline in [14] specifies the first six libraries for software engineering literature review. We included *Google Scholar* following the practice in [6], [7], to be comprehensive. To the best of our knowledge, these seven libraries should be comprehensive in covering literature from major software engineering venues.

3.3 Study Selection

Next, we conducted three rounds of selection to determine papers in our final study dataset. Table 1 lists the inclusion and exclusion (I&E) criteria.

TABLE 1: Inclusion and Exclusion Criteria

ID	Inclusion Criteria
I1	The paper is peer-reviewed.
I2	The paper is written in English.
I3	The paper contains six or more pages.
I4	The paper is published in an international conference, journal or symposium.
ID	Exclusion Criteria
E1	A previous version of the paper whose extended version has been included.
E2	The paper is a secondary study (literature review) of existing techniques/approaches.
E3	The paper focus on performance but does not consider software architecture.
E4	The paper relates to architecture but does not focus on software performance.
E5	The paper focuses on hardware-related performance.
E6	The paper is not a full paper, e.g. missing evaluation.

3.3.1 1st Round: Basic Information Reading

This round examined the basic information—the publication venue, the number of pages, and the title of each paper.

We first applied a set of objective inclusion criteria, including I1, I2, I3, and I4 (defined in Table 1) to retain papers that met certain quality specifications. I1 ensured that only peer-reviewed papers were included. For example, white papers, technical reports, and thesis/dissertations were not included. I2 ensured that only papers written in English were included. I3 ensured that papers with six or more pages were included. The motivation was to only include full research papers with elaborated research methodology and evaluation. Although most premier software engineering venues, such as ICSE, FSE, and ASE, have a ten page requirement for full research papers, we loosen this to six or more pages to make sure that relevant papers from venues with a different requirement were also included. Based on our observations, papers less than six pages were not likely to contain sufficient detail. I4 ensured that we only included papers published in an international conference, journal or symposium. Studies published at regional venues, such as “4th India Software Engineering Conference” [40] and “Journal of Zhejiang University” [41] were removed by this criterion.

Next, we scanned the title of each paper, and applied two exclusion criteria, E2 and E5. E2 excluded papers that were secondary studies, such as literature reviews or surveys. For example, we excluded the paper [42] “A systematic literature review on methods that handle multiple quality attributes in architecture-based self-adaptive systems”. E5 excluded papers that focused on hardware performance such as multi-core processors. For example, “Performance Evaluation of a Hybrid Computer Cluster Built on IBM POWER8 Microprocessors” [43] apparently focused on hardware. Note that it was not always possible to apply E2 or E5 simply based on the title of the paper. As discussed later, these two criteria were also used in the next two rounds with greater understanding of a paper.

3.3.2 2nd Round: Skimming

In this round, we skimmed through the abstract, introduction, and conclusion of each remaining paper, which helped us to grasp the key theme of the study. We applied the exclusion criteria, E1, E3, E4, and E5 based on the skim reading. E1 excluded papers which had an extended version in our search results. As an example, [S9] indicated that “The paper is an extension of our previous work [44]”, and thus [44] was removed. Criteria E3 and E4 excluded papers that focused either on architecture or on performance, but not the integration of the two. For example, the study [45] was excluded by E3 because it only profiled the performance of big data applications. The study did not investigate the

software architecture of the system. As an example for E4, the study [46] was excluded since it investigated the effort for the architecture design of software systems, while the performance of the systems was not a key focus. Criterion E5 removed papers that focused on hardware performance. For example, we realized that “Modeling performances of concurrent big data applications” [47] analyzed the speed and power consumption of multi-core processors in data center infrastructures.

3.3.3 3rd Round: Intensive Reading

In the final round, we carefully reviewed each remaining paper. Here, we applied all the exclusion criteria (E1 to E6). Note that some of the criteria (E1, E2, E3, E4, and E5) were already applied in previous rounds based on partial reading. However, in the previous two rounds we always preferred inclusion over exclusion for cases lacking strong confidence. The goal was to make sure that we did not exclude papers recklessly. Thus, these five criteria were evaluated again based on the updated understanding of the full content in the paper. For E1, some papers (e.g., [S10], [S9]) clearly stated that they were extensions to previous versions in the introduction, and thus E was applied in skim reading. While, some papers, e.g. [S11], clarified this in the background or even study approach sections, and thus were excluded by E1 in intensive reading. In comparison, E6 was evaluated for the first time, which excluded papers that did not contain any sort of evaluation or case study. As examples, [48], [49] were excluded due to E6, since they only presented conceptual solutions without any evaluations or experiments.

3.4 Snowballing Search

As shown in Figure 1, the papers selected after the three rounds were called the *Initial Set*. Next we performed backward snowballing based on papers in this *Initial Set*. Our goal was to retrieve additional papers that were missed in the search process. For example, there may exist papers that were not included in the seven digital libraries used in our study. As discussed later, we were able to identify an *Additional Set* with 699 papers by tracing the references of each paper in the *Initial Set*. Of note, this was after excluding papers published before January 2010 and papers that were already in our initial search. We repeated the three rounds of selection as described above on the *Additional Set* from snowballing. This ensured that all the retrieved papers, through initial search or snowballing, were evaluated following the same process and criteria. Therefore, it is reasonable for us to claim that our study captures related studies comprehensively.

3.5 Data Annotation and Synthesis

In this step, we carefully read each paper and annotated information that correspond to the RQs discussed in Section 3.1. Table 2 lists the detailed data items we annotated. There are two types of data items: 1) general information, such as the title, publication type, authors, publication years and the number of (#) pages—listed as D1 to D5 in the first half of the paper; and 2) specific information that maps to the RQs — listed as D6 to D18.

For RQ1, we first extracted text describing the motivation of each paper as D7. Next, two authors independently

TABLE 2: Annotated Data Items

ID	Data Item	Description	Correspond to
D1	Document Title	The title of the article.	General
D2	Publication Venue	Published in which Conference, journal, or workshop, etc.	
D3	Authors	The author names.	
D4	Publication Year	The year when the study is published.	
D5	# Pages	The number of pages in the paper.	
D6	Purpose Category	Mapping of this paper to the hierarchical categorization of all papers.	RQ1
D7	Motivation	The extracted description of the motivation of the paper.	
D8	Purpose Code	The codes for key study purposes used by two authors independently by open coding.	
D9	Development Activity	The software development activities the proposed techniques aim to facilitate	
D10	Study Process	The summary of key study steps	RQ2
D11	Architecture Model	The architecture model employed in the study.	
D12	Performance Model	The performance model employed in the study.	
D13	Annotated Parameters	The performance parameters being annotated to architecture models	
D14	Performance Metrics	The metrics used to measure software performance in a study.	
D15	M2M Integration	The integration of architecture and performance model analysis.	
D16	Modelling Tools	Tools that are used for creating and analyzing architecture and performance models	RQ3
D17	Profiling Tool	The profiling tool for collecting the performance metrics.	
D18	Manual Effort	The manual effort involved in the research steps.	
D19	Future Work	The discussed limitations and future work directions.	Future Directions

TABLE 3: Examples of Manually Open Coding for Study Purpose Categorization

Reference	Extracted Motivation (D7)	Author-1 Coding (D8)	Author-2 Coding (D8)
[S12]	Uncertainty is particularly critical in the performance domain when it relates to values of parameters such as workload, operational profile, resource demand of services, service time of hardware devices, etc. The goal of this paper is to explicitly consider uncertainty in the performance modelling and analysis process.	Recognizing and managing the uncertainties in software models	Model-based performance evaluation to handle uncertainty.
[S13]	Current approaches to modeling the context of software components are not suitable for use at run-time. In this paper, we aim to propose a performance meta-model (DMM) for run-time performance prediction, specifically designed for use in online scenarios.	(1) Model-based prediction; (2) Run-time predictions.	Model-based run-time prediction using DMM.

performed open coding, D8, which summarized the key objective of the study based on D7. Finally, we merged the codes from D8 to derive a hierarchical categorization of the papers, and mapped each paper to a category recorded as D6. Table 3 shows examples of the annotated motivation (D7) and the codes (D8) generated by the two authors.

After the annotation and open coding were done, the two authors constructed a hierarchical categorization of the papers by merging their codes. They first came up with the most high-level categorization of the papers based on the most fundamental objectives. For example, the codes of [S13] suggested the key theme of predicting and evaluating performance based on architecture/performance modeling; thus they were merged to derive the category called *Model-based Performance Prediction*. Similarly, [S14] derived another category called *Performance Anti-pattern Detection and Solution*. Next, the two authors further derived sub-categories based on the specific focuses (if available) of each paper. For example, both [S15] and [S13] specifically focused on “run-time” prediction reflected in the codes generated by the two authors. Thus they led to a sub-category, named *Run-time Prediction*.

In the process of deriving the sub-categories, two authors had disagreements on nine papers. All authors were invited to read these papers in full, and each author summarized the key theme of the paper independently. Then, the team employed the Delphi technique in a team discussion to reach consensus on these nine papers.

In addition, we annotated information (D9) about the development activities that the proposed approach was intended to facilitate. For example, paper [S16]

indicated that “*we proposed a modelling approach that provides fast evaluation to support design choices*”. Thus, this paper [S16] aimed at facilitating the software design. Some papers contributed techniques that may facilitate multiple development activities.

For RQ2, we summarized the study process (D10) of each paper. This was usually based on the description of the proposed technique in the approach section of the paper. For example, for paper [S17], the study process was annotated as “*1) Model Acquisition: modeling the system architecture using UML diagrams; 2) Model Annotation: use UML-MARTE to annotate the properties such as workload and resource demand; 3) Transformation: use ArgoSPE to translate UML into GSPN; 4) Performance Analysis: calculate response time, scalability, and resource utilization by perf model simulation.*”. In addition, we also recorded the detailed information in the study process, including the employed architecture model (D11), the performance model (D12), performance parameters in the modeling (D13), performance metrics being analyzed (D14), and the technique (D15) to integrate the analysis of architecture and performance models.

The study templates of different research purposes were synthesized by identifying and aggregating common research steps in the process annotations. For example, model-acquisition was an obvious common step for all the studies that leverage architecture modeling for performance prediction. Some papers retrieved architecture models based on the design documents [S1], [S3], [S4], [S8], [S16], [S17], [S18], [S19], [S20], [S21], [S22], [S23], [S24], [S25], [S26], [S27], [S28], [S29], [S30], while other papers recovered software architecture models from the existing systems by using

reverse engineering techniques [S12], [S13], [S15], [S31]. Of particular note, papers may use different terms for the same process.

Furthermore, each paper may contain more fine-grained sub-steps. As an example, [S26] presented a classification of performance-related parameters for annotating architecture models of a system. This sub-step categorized the performance-influencing factors as workload-specific characteristics (e.g., number and time interval of requests) and system-specific characteristics (e.g., resource demand of a specific system component), for the higher-level purpose of architecture model annotation. We considered this classification as a part of the model annotation step in the high-level study template. In addition, some papers may contain unique steps that did not recur in other studies. We did not include such unique steps as a part of the high-level research template.

RQ3 focused on collecting the information of tools that automate architecture and performance analysis. We recorded any tools that were used or developed for the analysis (D16). We also recorded any profiling tools (D17) that were used to collect the software performance metrics at run-time. We noted any explicitly mentioned manual effort (D18) in each study.

Finally, to explore future research directions, as discussed in Section 5, we annotated the limitations and future directions (D19) discussed in the papers.

3.6 Threats to Validity

In this section, we discuss the threats to validity of this study, following the guidelines provided in [50].

3.6.1 Threats to Construct Validity

First, we cannot guarantee that the search string as constructed captured 100% of all possible related studies — especially those that used less common terms for software architecture or performance. This may impact the comprehensiveness of the results in this paper. To mitigate this problem, we included synonyms, alternatives and hypernyms in the search string to represent the concept of architecture. For instance, the keywords that are related to architecture include “*architectural*”, “*architecting*”, “*design*”, and “*structure*”. Plus, we also applied an additional snowballing search step to include the potentially missed studies in the first round of search. Thus, it is reasonable to assume that our study has covered most of the related literature in this field.

In addition, we acknowledge that our study process did not include the formal quality assessment step suggested in [14]. This step was to ensure that the literature review was based on high quality studies. This required the definition of a formal assessment checklist and a scoring schema based on the checklist, to be performed on each individual study. However, we would like to argue that some of the inclusion and exclusion criteria, including I1, I3, I4, and E6 in Table 1, served as filters to ensure the selected papers met reasonable quality standards. In particular, E6 filtered out preliminary studies whose validity was not evaluated. We did not define any further assessment criteria, such as the soundness of the proposed approaches and experiment

results. The reason was that the retrieved literature was rich in architecture models, performance models, model transformation, and performance profiling tools/techniques. An in-depth assessment of study quality was beyond the scope of this literature review.

3.6.2 Threats to Conclusion Validity

Although we followed the guidelines in [13] for ensuring stability and reliability requirements, we cannot guarantee that our study results were free of any potential personal biases. The study results were potentially impacted by our biases in conducting the selection, data annotation, and data aggregation and synthesis. This is a threat to conclusion validity. However, we believe that following our replication package <https://sites.google.com/view/arch-perf-data-repo/>, another researcher would be able to replicate our study following the same process and would reach consistent results and conclusions.

To best mitigate personal bias in the selection process we employed two researchers to conduct the selection in *Skimming Reading* and *Intensive Reading* following the inclusion and exclusion criteria in Table 1. The two researchers agreed on most papers. The *Inter-Rater Reliability* of the final decision after these two rounds of selection was 86%. Next, to avoid inaccurate or incomplete data annotation, the data items in Table 2 were provided with detailed descriptions in our data repository. The data annotation was initially completed by the first author, and then checked by the second and third authors to high-light disagreements. The three authors held group discussions through online chat and real-time meetings to resolve disagreements. Any remaining disagreement after the review of the three authors was managed through group discussion. Furthermore, to best mitigate potential biases in the data synthesis, two researchers worked on analyzing the raw annotation data, and presented and discussed the derived results with all authors. In particular, the categorization of research purposes in RQ1 emerged from the data following an open coding approach, conducted by two authors independently, and then extensively discussed and iteratively refined by the entire research team.

3.6.3 Threats to Internal Validity

The validity of our study results in RQ3 (research templates), RQ4 (available tools and instruments), and RQ5 (evaluation) are subject to a threat to internal validity associated with the selected papers. The results for RQ3, RQ4, and RQ5 were synthesized from data annotation of the information items from the literature. If the authors of the original study failed to accurately provide the details of their study, this would impact our study results. For example, if the authors of a paper did not describe certain steps in their study which were not the key study contribution, this may impact our aggregation of the typical research templates in RQ3. Similarly, if the authors did not report the tools or instruments they used for architecture modeling and performance analysis, our results in RQ4 would be impacted. This threat was associated with the retrieved literature, and was largely out of our control, but is presumably mitigated by the large number of papers studied.

3.6.4 Threats to External Validity

This study potentially suffers from threats to external validity. First, we focused research published between 2010 and 2020. We cannot guarantee that the findings derived from this range could be generalized to studies published before 2010 or after 2020. In addition, we cannot guarantee that the tools and techniques identified for software architecture modeling and analysis were representative of the tools and techniques used in a more general setting without performance analysis. Likewise, we summarized instruments for performance profiling in RQ4, but we cannot guarantee that these findings were representative of all the performance profiling tools that were used independent of architecture analysis.

4 STUDY RESULTS

4.1 Overview of Study Dataset

After searching the seven digital libraries using the search terms, we retrieved a pool of 24901 papers in total. In the 1st round, i.e., basic information reading, we retained 1039 articles. In the second round, i.e., skimming, we retained 313 articles. In the third round, intensive reading, we retained 89 articles as our *Initial Set*. Next, we applied a snowballing search process starting from these 89 articles. We collected a pool of 793 additional articles, after excluding papers before 2010 and removing duplicates. After repeating the three selection rounds on these articles, we identified an additional 20 articles—195 articles retained after the first round; 111 articles retained after the second round; and 20 articles after the third round. Thus, there were a total of 109 (89+20) papers in the final set.

The majority, 68%, of the studies were peer-reviewed full conference papers. 22% of the studies were journal articles. 10% of the studies were peer-reviewed workshop papers. Table 4 lists the details of the publication venues. Those contributing only one study were combined as “Others”. The venues that had the largest number of related papers were either specialized in software architecture, such as QoSA¹, ICSA, and ECSA or specialized in software performance, such as ICPE.

TABLE 4: Publication Venues

Publication Venue	Abbreviation	Type	No.	%
Inter. Conf. on Performance Engineering	ICPE	Conference	13	12%
Inter. Conf. on Quality of Software Architectures	QoSA	Conference	12	11%
Inter. Conf. on Software Architecture	ICSA	Conference	10	10%
European Conf. on Software Architecture	ECSA	Conference	7	7%
Inter. Conf. on Software Engineering	ICSE	Conference	3	3%
Journal of Systems and Software	JSS	Journal	4	4%
IEEE Transaction on Software Engineering	TSE	Journal	3	3%
Inter. Workshop on Software and Performance	WOSP	Workshop	2	2%
Others			43	41%

4.2 RQ-1: Study Purpose Categorization

Figure 2 illustrates our categorization of the 109 studies based on objectives and specific focuses. Overall a majority, 73 (67%), of the studies leveraged architecture modeling for performance analysis, including performance prediction, root cause analysis, or architectural optimization. Given the large number of studies in this category, they were further sub-categorized based on the focus of each study.

The other 36 (34%) studies were categorized into three purposes: performance anti-pattern detection and resolution, performance profiling and comparison of alternative architectural solutions, and constructing software systems with self-adaptive architectures.

Next, we discuss each category in detail.

4.2.1 Category 1: Model-based Performance Analysis

This is the largest category, with 73 (67%) studies. These studies were sub-categorized into: 1.1 Model-based Performance prediction; 1.2 Retrospective Root Cause Analysis of performance issues; and 1.3 Model-based Architectural Optimization.

Category 1.1 Model-based Performance Prediction:

These studies aimed at predicting the performance metrics of a system, such as response time, resource utilization, and throughput, based on architecture models of the system; they employed architecture models to capture the abstract behaviors or component structures of a system, enriched with performance properties, such as resource and workload. The architecture models were then transformed into mathematical performance models, such as *Queueing Networks* and *Stochastic Petri Nets*, or simulation code to derive performance measures. Practitioners could benefit from such techniques by assessing performance and addressing concerns early in the development life cycle.

Category 1.2 Retrospective Root Cause Analysis:

Four studies [S32], [S33], [S34] provided retrospective analysis of the root causes of performance issues. [S33] revealed that a significant portion (around 30%) of real-life performance issues require design-level optimization, i.e., simultaneously revising a group of related source code files. This study used a new modeling technique, named *Diff Design Space Matrix (D-DSM)*, to capture the essential design structure changes in code revisions for performance issues. Another study [S34] proposed an attributed graph model to investigate performance anomaly propagation across micro-services. [S32] proposed a fine-grained (i.e. method-level) architectural modeling approach, named *Butterfly Space* modeling, that combined the analysis of dynamic profiling metrics and static method-call dependencies to understand architectural patterns behind performance issues.

Category 1.3 Model-based Architectural Optimization: There are three studies [S2], [S35] in this sub-category which automatically recommend optimal architectural solutions by exploring the design space of a system. [S2] proposed an evolutionary algorithm for rule-based software performance optimization at the software architecture level. The proposed approach addressed the limited search space and uncertain application of rules in traditional rule-based performance optimization. Similarly, [S35] integrated performance monitoring data and architectural modelling to detect performance problems and suggest architectural changes.

Given the large volume of the studies in Category 1.1 we further mapped them into five clusters based on the focus of each study. Of particular note, one study [S36] compared different model-based prediction methods in terms of the prediction accuracy and the required manual effort,

1. QoSA is merged into ICSA after 2016.

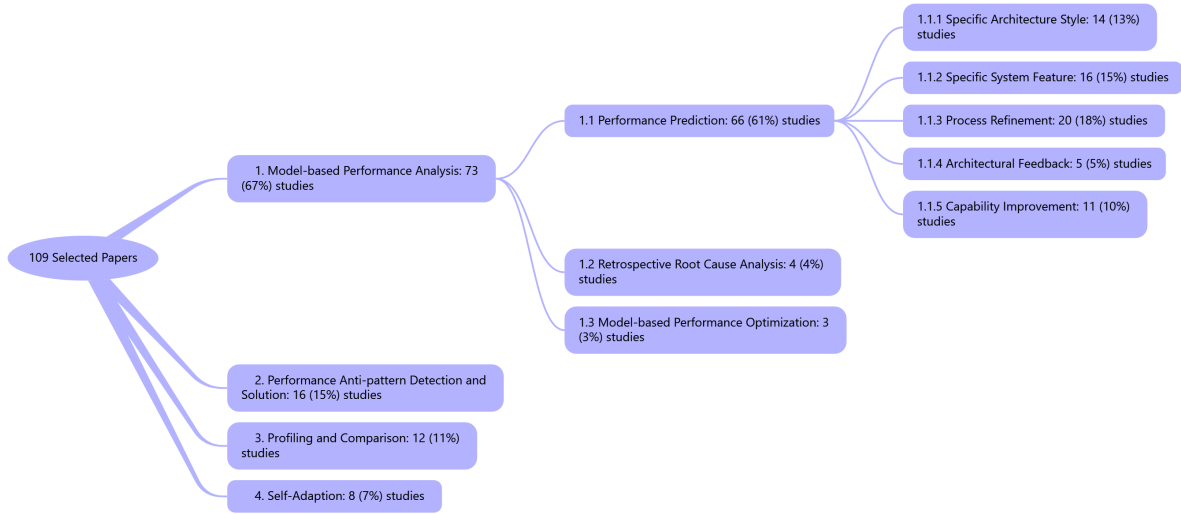


Fig. 2: Study Categorization based on Purposes (RQ1)

for providing insights to practitioners. This paper was not classified into any of the five clusters described below:

Category 1.1.1 contains studies that contributed prediction techniques for specific architecture styles. Nine studies focused on distributed architectures [S9], [S22], [S30], [S37], [S38]. They employed modeling techniques with rich notations for capturing the behavior and dependencies of distributed system components, such as *Directed Acyclic Graph (DAG)* [S30], *Discrete-Time Markov Chain (DTMC)* [S22]. In particular, [S30] focused on service-oriented architectures. [S17] focused on interoperable architecture, which enabled diverse interaction modes captured as performance parameters in the model.

Category 1.1.2 contains studies that targeted the performance prediction of specific systems. The most well studied were data intensive systems [S4], [S16], [S18], [S19], [S20], [S21], [S25], [S26], [S39]. Computational resources were critical due to the high demand for data storage, transmission, and manipulation. These studies all considered parameters for representing resource demand in the modeling process.

Three studies focused on performance prediction for systems featured by hardware and software co-design, including cyber-physical system [S22], real-time embedded system [S23], and Internet of Things system [S8]. Hardware-related properties, such as communication delay among hardware components, were considered in the modeling process. In addition, architecture patterns for the hardware and software co-design are assessed for deriving the system performance [S22], [S23]. [S24] focused on performance prediction for systems with dynamic and varying processing requests.

Category 1.1.3 contains studies that refined common processes in performance prediction techniques for automation, efficiency, and accuracy improvement [S12], [S27], [S28], [S29], [S40], [S41], [S42]. [S27], [S40] focused on developing automated methods for extracting and constructing architecture and performance models [S27], [S41]. [S12] annotated performance properties and parameters on architecture models. In particular, practitioners often faced uncertain performance properties, such as operational profile and resource demand, which

compromised the accuracy of performance prediction. [S12] focused on explicitly capturing and managing such uncertainties in the modeling process. [S42] made the performance model solution and simulation faster, since the time required to solve or simulate a model increases exponentially with system size. Finally, [S40] simplified and automated the overall prediction process, by proving a unified interface for performance queries, independent of the employed modeling formalism.

Category 1.1.4 contains studies [S31], [S43] that focused on providing traceable architectural feedback based on performance prediction results. [S43] focused on visualizing possible architectural refactorings based on predicted results. [S31] focused on building backward traceability links from predicted results to the performance and architecture models to help in interpreting the prediction results.

Category 1.1.5 contains studies that enabled performance prediction under specific settings. Four studies [S3], [S13], [S15], [S44] enabled run-time or online performance prediction. This pushed the boundary of performance prediction that traditionally only occurred at design time, to satisfy the needs of continuous system reconfiguration in modern applications. These studies featured flexible architecture models, such as *Descartes Modeling Language* [S13], [S15], [S44], for annotating dynamic performance parameters in the architecture models. 3 studies [S1], [S11], [S45] contributed prediction techniques that allow continuous calibration of architecture model and performance parameters to reflect continuous system evolution.

4.2.2 Purpose 2: Performance Anti-pattern Detection & Resolution

The second high-level study purpose is focused on architectural anti-patterns that lead to performance problems. We identified 16 studies in this category [S6], [S10], [S14], [S46], [S47], [S48], [S49], [S50], [S51], [S52], [S53], [S54]. Performance anti-patterns are recurring issues in design or implementation that lead to negative performance consequences. The three studies by Smith *et al.* [51], [52], [53] summarize 12 anti-patterns from prior studies.

Each performance anti-pattern is usually paired with a problem and its solution. For example, *Blob/God Class* was characterized by a high number of messages transmitted between a single big component and a large number of other components. It could be resolved by breaking down the responsibilities of the *Blob/God Class* to smaller parts, to split the “busy” message traffic.

Related studies contributed automated approaches for detecting and refactoring such performance anti-patterns. The detection of anti-patterns usually relied on a set of rules that combine the analysis of performance metrics and architectural characteristics. The detection rules usually compared the performance metrics with a certain threshold, which were either obtained from model-based prediction [S10], [S14], [S49], [S50], [S51], [S52], [S54], [S55] or dynamic profiling [S6], [S46], [S47], [S48], [S53]. The detection rules also matched typical characteristics in the architecture model that formed anti-patterns. For example, the *MessageSize Rule* checked if the deviation of service throughput associated with different components exceeds the respective threshold. Once matched, the *MessageSize Rule* identified a potential bottleneck in the *Pipe and Filter* pattern.

In addition to anti-pattern detection, 12 studies [S6], [S10], [S14], [S14], [S48], [S49], [S50], [S51] provided resolutions to eliminate the anti-patterns. For example, “*changing scheduling algorithms to enable concurrent execution*” resolves *Concurrent Processing Systems*; and “*alleviate the congestion and use Shared Resource Principle to minimize conflicts*” resolves *One-Lane Bridge*.

4.2.3 Purpose 3: Profiling & Comparison

12 (11%) studies focused on dynamically profiling and comparing the run-time performance of alternative architectural solutions for implementing the same applications [S5], [S56], [S57], [S58], [S59], [S60], [S61], [S62], [S63], [S64], [S65]. The architectural variables considered in these studies span a broad spectrum, including: 1) the types of data storage systems employed, such as *Cassandra*, *MongoDB*, *NoSQL*, and *MySQL* [S57], [S58]; 2) architectural patterns such as *Leader/Follower* pattern and *Half-Sync/Half-Async* pattern [S59] and *Workload Smoother* pattern [S62]; 3) web-service architecture solutions [S5], [S61]; and 4) deployment architectures [S56], [S60], [S63], [S64].

Alternative architectures for data storage we compared in [S57], [S58]. [S58] compared the performance of two widely-used database management systems with different architectures, *Cassandra* and *MongoDB*. *Cassandra* was for storing large amount of unstructured data. It had a multi-master architecture where data was distributed evenly across clusters to guarantee load balancing. *MongoDB* was a document-oriented NoSQL data storage system that automatically partitioned data across multiple servers—known as sharding. It had a master-slave architecture where all writing operations must be directed to a single master cluster. The comparison showed that under heavy-reading workload, *MongoDB* slightly outperformed *Cassandra*. However, under heavy-writing workload, *Cassandra* substantially outperformed *MongoDB*. Similarly, [S57] compared different implementations of a web application with three alternative databases—the traditional *Hibernate + MySQL*, and two *NoSQL* data storage

systems, *MongoDB* and *Neo4j*. The comparison showed that both *MongoDB* and *Neo4j* outperformed *MySQL + Hibernate* in scenarios of processing document-oriented and graph-oriented data.

Two studies compared and contrasted various architectural patterns that can impact system performance. [S59] compared two widely-used concurrent processing architectures, *Half-Sync/Half-Async (HS/HA)* and *Leader/Follower (LF)*. The *HS/HA* pattern consisted of an asynchronous layer for the high priority requests in a message queue to avoid packet loss, and a synchronous layer for other requests and necessary I/O and computational operations. *HS/HA* had the advantage of decoupling the asynchronous and synchronous services to avoid multi-thread blocking. In comparison, *LF* had only one leader thread at all times, and when this thread finished processing a request, it became a follower thread and returned to the thread pool. The advantage of *LF* was to minimize overhead and prevent race conditions. The experimental results showed that *LF* performed slightly better than *HS/HA* under a small workload; while *HS/HA* performed much better with a large number of threads. The study [S62] evaluated the implementation of a workload smoother design pattern on Function as a Service (FaaS) cloud platforms, such as Amazon AWS Lambda, IBM, and Azure Cloud Function. A workload smoother helps distribute traffic evenly over time, thus reducing bursts and ensuring smoother performance. The results indicated that different FaaS platforms adopted unique scaling strategies, and by applying a workload smoother, software engineers could achieve a success rate of 99 - 100%, compared to a 60 - 80% rate when the FaaS system was saturated.

The comparison of different web service architecture solutions has been the focus of the two studies [S5], [S61]. [S5] explores two different MVC implementations - the *Laravel Framework* and the *Slim Framework*. The *Laravel Framework*, being a full-stack MVC framework, enabled programmers to concentrate more on implementing business logic. Conversely, the *Slim Framework* contained only a View, with the Controller and Models added as needed. They concluded that for small projects, the *Slim Framework* was more suitable, whereas, for larger projects with numerous functions and libraries, the *Laravel Framework* was the better choice. [S61] compares the performance between 2-tier and 3-tier cloud applications, specifically Ghost and WordPress. Ghost, a blog cloud application without a database, was pitted against WordPress, a similar application but employing RDS (Relational Database Service) to host a Maria database. Their findings indicated that WordPress, a 3-tier blogging system, not only performed better but also exhibited greater stability than the 2-tier Ghost blogging system.

Four studies [S56], [S60], [S63], [S64] focused on comparing deployment architectures. Four of these studies concentrated on profiling and comparing deployment architectures, specifically between monolithic and micro-service architectures [S56], [S60], [S63], [S64]. [S56] compared micro-service architecture deployment environments, taking into account factors like memory allocation, CPU fraction in the deploying servers, and the number of Docker container replicas assigned to each

micro-service. [S60] delved into the migration of a large object-oriented system from a monolithic to a modular, microservice architecture, removing circular dependencies between modules. Performance metrics before and after the architecture migrations were compared, revealing variations in processing times, throughput, and other parameters. [S63] contrasted the performance and scalability of monolithic and microservice architectures within a reference web application, concluding that a monolith excels in single-machine performance compared to a microservice-based architecture. Lastly, [S64] proposed a quantitative approach for the performance assessment of microservice deployment alternatives. Their four-step methodology incorporated operational profile data analysis, experiment generation, baseline requirements computation, and experiment execution, thereby providing a comprehensive lens to evaluate deployment configurations. [S95] compared the performance of web applications deployed on two server architectures, i.e., 1) the *Active/Active* environment that consisted of multiple redundant servers where the workload was balanced by a load-balancer installed between web services and web applications; and 2) *Active/Passive* environment that consisted of a single server that provided service at any time while the other servers were used as standby devices. The study showed that the *Active/Active* deployment architecture could provide more stable processing time and CPU utilization.

4.2.4 Purpose 4: Self-Adaption

Researchers have proposed a number of frameworks that monitor the performance of the system, and dynamically switch to the most efficient architectural solution based on run-time workload [S7], [S66], [S67], [S68], [S69], [S70].

The study [S66] presented a decentralized, autonomic architecture for managing multi-tier transaction applications deployed on the cloud environment. For each tier in the application, an autonomic controller (in form of a *Java* class) monitored the resource utilization of the tier, then enacted a set of elasticity policies that dynamically scale the acquisition and release of cloud resources to and from the tier. The experiments demonstrated the effectiveness of the proposed architecture.

[S67] proposed a framework named *SAFCA*, which contained three well-known concurrent thread processing architectures—*Dynamic-thread-creation (DTC)*, *Half-Sync and Half-Async (HS/HA)* and *Leader-Followers (LF)*. The framework automatically switches to the most efficient mechanism, based on adaptation policies. The framework activates *HS/HA* during stable, continuous workload. When a sudden burst of requests is detected, the framework automatically switches to *DTC*. When the system experiences failures, the framework switches to *LF* because it ensures higher reliability. After the system recovers, the framework switches back to *HS/HA*.

[S69] proposed a self-adaptive framework was proposed to optimize tail latency in event-driven microservices within large-scale cloud applications. This adaptive dual-mode autoscaling mechanism features: 1) a reactive mode that actively monitors event arrival rates and adjusts consumer microservice instances to respond to load changes promptly, and 2) a proactive mode that employs an

autoregressive prediction model to foresee future event loads, thus enabling preemptive scaling. This proactive approach, which continuously learns online using an exponentially weighted recursive least squares algorithm, prevents latency spikes by facilitating a seamless transition between modes once a configured accuracy level is achieved. This swift mode-switching capability enhances system adaptability to fluctuating loads, thereby minimizing latency. Experimentally, the framework demonstrated its effectiveness in maintaining tail latency Service Level Agreement (SLA) guarantees while operating with a minimum number of replicas.

In [S70], a performance enhancement in the assembly and optimization of parallel component programs is achieved by employing a set of distinct software agents that collaboratively implement four adaptive strategies. These strategies are: 1) Dynamically adjusting the parallelism of components, which adapts to the varying resource availability and load, thus ensuring efficient utilization of resources; 2) Altering data partitioning to balance the computational load across components, enhancing overall system throughput; 3) Enabling component migration, shifting load away from high-usage nodes to those with low load, thereby reducing potential performance bottlenecks; 4) Modifying implementation by switching between different versions of component implementation codes based on resource and performance requirements, thereby maintaining optimal system performance. The strategies' effectiveness was verified through experimentation on heterogeneous computer clusters, exhibiting significant performance benefits and flexibility over traditional methods.

In recent years, researchers started to employ machine learning techniques in self-adaptive systems [S7], [S68], rather than relying on pre-defined adaption policies or algorithms. Machine learning was used in making decisions about switching among alternative architectural solutions [S68], and in predicting the performance behavior of a system under execution environments [S7]. [S68] exploited an evolutionary algorithm, the "*NSGA-II genetic algorithm*", to suggest near-optimal alternative architectures in terms of mean response time for different system operational modes. [S7] used a classic machine learning model, "*Long-Short Term Memory (LSTM) Networks*", to forecast the behavior of each component of a system to facilitate adaption.

4.2.5 Facilitated Activities

We also the development activities that different studies aimed to facilitate, and maps the identified activities to the research objectives from RQ1. This provides guidance for practitioners to identify the most relevant techniques for the development activities they are engaged in.

Five software development activities (or phases) emerged from our analysis: 1) *Software Design* when the developers focus on the high-level and preliminary architecture and design of their systems; 2) *Software Implementation* when developers implement and test their systems; 3) *Software Deployment* when developers launch their systems in the physical environment; 4) *Software Operation* when developers focus on the run-time behavior of their systems in the product environment; and 5) *Software Maintenance* when developers

TABLE 5: Proportion of Research Purposes according to Development Activities (RQ1)

Development Activity	Design	Implementation	Deployment	Operation	Maintenance
1. Model-based Performance Analysis	50 (48%)	4 (4%)	1 (1%)	8 (8%)	9 (9%)
1.1. Performance Prediction	46 (44%)	4 (4%)	1 (1%)	8 (8%)	6 (6%)
1.2 Retrospective Root Cause Analysis	0	0	0	0	4 (4%)
1.3 Architectural Optimization	3 (3%)	0	0	0	0
2. Performance Anti-patterns Detection and Solution	10 (10%)	1 (1%)	1 (1%)	0	5 (5%)
3. Profiling and Comparison	4 (4%)	0	4 (4%)	2 (2%)	1 (1%)
4. Self-Adaption	1 (1%)	0	2 (2%)	6 (6%)	0
All	67 (64%)	6 (6%)	8 (8%)	16 (15%)	15 (14%)

keep updating and supporting existing systems to solve arising issues.

In Table 5, rows 2 to 5 show the number and percentage of studies in each purpose that addressed each development phase. The cells colored grey highlight the largest numbers in each category. The bottom row in Table 5, (“All”), shows the total number (and percentage) of studies for each phase. Next we elaborate the details of each study purpose.

4.2.5.1 Purpose 1: Model-based Performance Analysis: In *Category-1.1 Performance Prediction*, most studies (XX%) aimed at facilitating software design [S4], [S8], [S9], [S16], [S17], [S18], [S19], [S20], [S21], [S22], [S23], [S24], [S25], [S26], [S28], [S29], [S30], [S39], [S40], [S41], [S42], [S43]. These studies conducted early performance assessments based on design models. The performance metrics predicted based on these models helps to evaluate early design choices, to avoid rework before a system is implemented. Eight (13%) studies targeted run-time performance prediction [S3], [S13], [S15], [S44], [S45]. They calibrated performance models with resource usage profiles while systems were running. Six (10%) studies targeted software maintenance [S1], [S11]. They derived software models from existing systems and provided performance evaluation considering architectural evolution.

The four studies in *Category-1.2 Retrospective Root Cause Analysis* [S32], [S33], [S34] target the software maintenance phase. They perform a root cause analysis of performance issues in existing systems, providing guidance for developers doing maintenance.

All three studies in *Category-1.3 Model-based Architectural Optimization* aimed at facilitating better software design solutions [S2], [S35]. They provided support for stakeholders in performing architectural refactoring on design models.

4.2.5.2 Purpose 2: Performance Anti-pattern Detection and Resolution: The majority of the studies in *Purpose-2* (ten out of sixteen, or 63%) [S10], [S14], [S47], [S49], [S50], [S51], [S52] focus on the identification and refactoring of performance anti-patterns during the software design phase. These studies employ UML diagrams and PCM as their analytical tools.

Conversely, five studies (31%) [S6], [S46], [S48], [S53] emphasize the maintenance phase of software development. They apply dynamic profiling to existing software applications and use performance metrics to identify and eliminate anti-patterns, thus helping maintain the health and efficiency of the systems. Additionally, paper [S54] extends this scope to include the implementation and deployment phases. This study introduces a novel approach for Software Performance Anti-pattern (SPA) characterization and detection, designed to bolster continuous integration/delivery/deployment (CI/CD) pipelines. This inclusion addresses the computational

efficiency gaps in current detection algorithms, enhancing performance across all development stages.

4.2.5.3 Purpose 3: Profiling and Comparison: In this category, four (36%) studies, [S5], [S57], [S58], [S59], concentrated on the software design phase by providing empirical experience regarding design choices like databases, web frameworks, and programming paradigms. Focusing on the deployment phase, four (36%) studies, [S56], [S61], [S64], offer quantitative assessment and guidance for deployment decisions, including architectural comparisons and performance assessment of deployment configurations in micro-services and web-services environments. Two (18%) studies, [S62], [S63], emphasized the operation phase by exploring the effectiveness of workload smoother patterns and comparing the performance and scalability of different architectural styles in operational conditions. One study, [S60], with a focus on the maintenance phase, presented a detailed refactoring process, from a large monolithic to a modular system, analyzing the performance impact and refactoring effort.

4.2.5.4 Purpose 4: Self-Adaption: Self-adaptive systems dynamically adjust their configuration at run-time to meet performance goals in changing environments. The majority of studies (77%) focused on system operation at run-time [S7], [S66], [S67], [S68], [S69]. Notably, two studies, [S66] and [S70], went a step further and not only reconfigured the software components but also re-deployed resources in the deployment environment. Consequently, these studies also offer valuable insights for the deployment phase.

RQ1 Summary: The 109 studies that integrate software architecture and performance analysis serve four objectives 1) model-based performance analysis; 2) performance anti-pattern detection and resolution; 3) profiling and comparison of architectural alternatives; and 4) self-adaptive architecture for dynamic performance optimization. These studies tackle challenges covering five major phases of the software development life-cycle, namely *Design*, *Implementation*, *Deployment*, *Operation*, and *Maintenance*. Overall, *Design* is the most extensively covered.

4.3 RQ-2: Classic Study Templates

RQ2 provides “cheat-sheets” for practitioners to quickly familiarize themselves with study templates that integrate architecture and performance analysis. This can help them quickly get started if they are interested in replicating or extending existing techniques. This RQ reveals how software architecture and performance analyses are combined for different study purposes. It offers an in-depth understanding of: 1) which architecture models are used and how they are analyzed; 2) which performance models or metrics

are used, and how they are collected and analyzed; and most importantly, 3) how the architecture and performance analysis are integrated for addressing performance concerns. We summarize typical study templates for different study purposes, which reveal the common study processes. To further facilitate comprehension, we've color-coded the common steps across the templates, which visually highlights their shared elements and helps practitioners quickly discern the commonalities and differences among them.

4.3.1 Model-based Prediction Template

This section presents the study template summarized from sub-category 1.1, *model-based performance prediction*, which contains the majority (70%) of the studies in the model-based analysis. It is not meaningful to summarize a study template for the other two sub-categories (1.2 and 1.3) with only three studies each.

As depicted in Figure 3a, the model-based performance prediction studies encompass five phases: 1) Model Acquisition, in which architecture models of a system are procured; 2) Model Annotation, involving the annotation of performance-related information on the architecture models; 3) Model Transformation, which transforms the annotated architecture models into executable performance models; 4) Model Solution or Simulation, that involves solving or simulating the performance models to predict performance metrics; and 5) Result Interpretation and Feedback, providing interpretation of the predicted results and feedback to stakeholders.

Note that, not every study fits 100% to this template, which captures the major steps we observed from the studies. In particular, most studies stopped at Step 4, i.e., obtaining the performance metrics. Among them, six studies [S1], [S4], [S22], [S34] did not have Step 3 since they used graph-based architectural models which do not need transformation before solution or simulation. Nine studies [S9], [S17], [S24], [S31], [S43] extended the process with Step 5, which traced back from the predicted results for providing interpretation of the prediction results and architectural feedback for stakeholders.

Step 1: Model Acquisition: We observed two methods to acquire a system model: 1) retrieving from project documentation and 2) recovering via reverse-engineering (introduced in Section 2). Most of the studies used the first method, where the architecture models were created in the design phase by the project team. Six studies [S1], [S11], [S27] used the second method.

We summarized three types of architecture models: 1) *Behavioral Models (BM)* capture system architecture based on the collaborations among system objects, their internal state changes, the interactions, e.g., *UML Sequence Diagrams*, *UML Activity Diagrams*, *UML Use Case Diagrams*, and *Business Process Execution Language (BPEL)* are the behavioral models. 2) *Component Models (CM)* capture the architecture of a system emphasizing the static characteristics as components, connectors, and their compositions, such as *Palladio Component Model (PCM)*, *Descartes Modeling Language (DML)*, *Method Call Graph*, *UML Class Diagram*, *UML Component Diagram*, and *UML Deployment Diagram* are used for this purpose. And 3) *Hybrid Models (HM)* capture the system architecture as combinations of component structure

and behavioral interactions, including the *Web Modeling Language (WebML)*, the *Service-Oriented Architecture Modeling Language (SoaML)*, and *Architectural Description Languages (ADLs)* are hybrid architecture models

Step 2: Model Annotation: The second step is to annotate the architecture models with performance-related information in preparation for performance prediction. We summarized five types of annotation data from the literature, including: 1) *Workload*, which describes the run-time workload on a system during execution. The workload can be specified as a status or the intensity in numerical values. 2) *Resource Demand*, which is the estimation of the resource demand for jobs, such as memory, CPU, and disks. 3) *Probability*, which is the probability of executing a specific conditional process in a system. Processes with high probability have greater impact on performance than low probability processes. 4) *Timing Specification*, which is the timing behavior of a system, such as the time intervals between requests, the time to acquire and release a job. And 5) *Scheduling Policy*—the policy for scheduling workload requests, such as “*First In First Out (FIFO)*”, “*Last In First Out (LIFO)*”, and “*Processor Sharing (PS)*”, etc.

Step 3: Model Transformation: The third step is to transform an architecture model, annotated with performance related parameters, to a form that can be used to predict performance metrics. We observed two types of transformation methodologies—*Model-to-Model (M2M) Transformation* and *Model-to-Text (M2T) Transformation*.

In *M2M Transformation*, the architecture model is transformed to a performance model, which captures the dynamics of a system in operation. We found three types of performance models in the studies, including *Queueing Networks* with four variations, *Stochastic Petri Nets* with two variations, and *Stochastic Process Algebras*. *Queueing Networks* are composed of *service centers* representing system resources and *jobs* representing the consumers of the resources [S4]. *Jobs* travel through a network of *service centers* following probabilistic routes. This can be used to estimate the response time of the system by considering waiting times, queue lengths, and server utilization. *Stochastic Petri Nets* are composed of two types of elements: *places* and *transitions* [S19]. A *place* can contain any number of tokens, which enables a *transition* if all other *places* connected to it contain at least one token. The transitions fire after a probabilistic delay determined by a random variable. *Stochastic Petri Net* are often used to model step-wise processes in distributed systems, which include choice, iteration, and concurrent execution [S19], [S22]. Finally, *Stochastic Process Algebras* are derived from classical process algebras, which are abstract languages used for the specification of concurrent systems [S30]. Systems are modelled as collections of entities, called agents, which execute atomic actions. The atomic actions describe the behavior of the entities in concurrent execution.

The *M2M Transformation* is usually based on graph transformation. There are third-party tools to automate the process, including, *PUMA* [S29], *PRIMA-UML* [S52], *ArgoSPE* [S17], *KLAPER* [S40], *TwoTower* [S49], and *S-PMIF+* [S8], [S41]. Some architecture models have default model transformation techniques, such as *PCM2LQN* [S28] and *PCM2QPN* [S28] for the *PCM* model, and

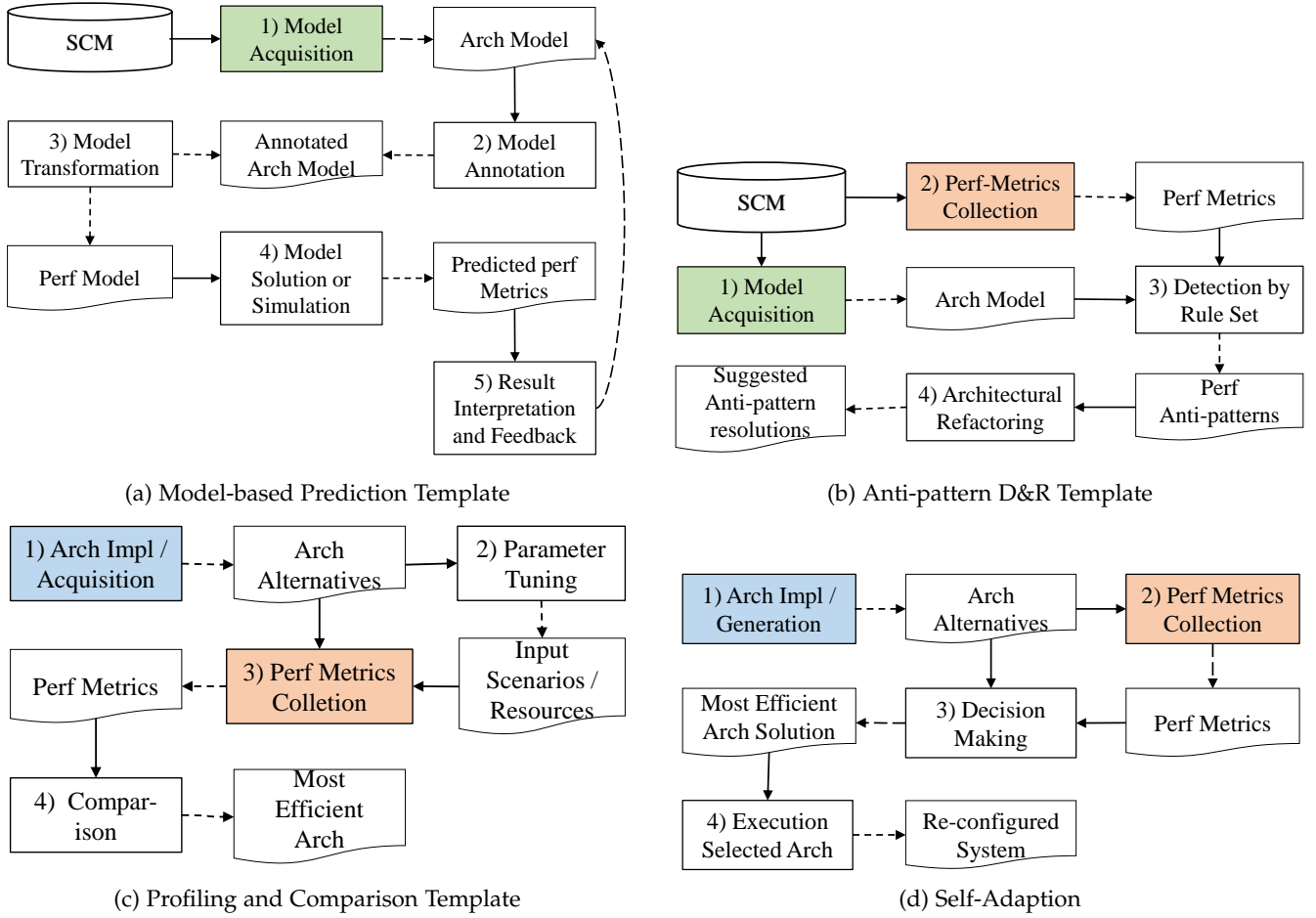


Fig. 3: Classic Study Templates

TABLE 6: Summary of Annotated Information with Applied Architecture Models

Model	Workload	Resource Demand	Probability	Timing Specification	Scheduling Policy
UML Diagrams	[S3], [S17], [S24], [S29], [S41], [S71], [S72], [S73], [S74], [S75], [S76], [S77], [S78], [S79], [S80]	[S3], [S12], [S17], [S19], [S29], [S71], [S72], [S73], [S74], [S76], [S77], [S79], [S81], [S82]	[S3], [S12], [S80]	[S12], [S17], [S75], [S77], [S82]	[S19], [S80], [S81]
PCM	[S11], [S18], [S20], [S25], [S26], [S28], [S30], [S83]	[S9], [S11], [S18], [S20], [S25], [S28], [S83], [S84], [S85], [S86], [S87], [S88]	[S9], [S28], [S84], [S88]	[S30]	[S20], [S85], [S86]
DML	[S13], [S27]	[S13], [S15], [S27], [S42]	[S13], [S15], [S42]		
Calling Dependency Graph	[?], [S21], [S89]	[?], [S4], [S21]	[S89]		[S4], [S16]
BPEL	[S72]			[S72]	
WebML		[S90]	[S90]		
SoaML	[S71], [S91]	[S91]	[S91]	[S91]	
ADLs	[S23]	[S23], [S92]			[S23]

DML2LQN [S44] for *DML* model. Some studies use declarative model transformation languages, including *Query/View/Transformation (QVT) language* [S41], *Janus Transformation Language (JTL)* [S35], and *Gra2Mol* [S49]. A few studies leverage meta-models to improve the accuracy while mapping components in the source model to the entities in the target model, including *Core Scenario Model (CSM)* [S29] and *Graph Transformation System (GTS)* [S30].

In *M2T Transformation*, an input architecture model is transformed to a textual artefact, which serves as executable simulation code [5]. The most commonly used *M2T Transformation* technique is *PCM2SimuCom*, which translates *PCM* into executable simulation code [S9], [S11], [S20], [S25], [S26], [S28], [S51]. *SimuCom* is the default analysis

tool for *PCM*. It packs the translation and execution function of simulation code together. Similarly, *SimuLizar* [S27] is a simulation tool for *PCM*. In addition, the tools *SLAstic.SIM* and *SimSo* generate simulation code for *Business Process Execution Language (BPEL)* and *ADLs* respectively.

Step 4: Model Solution or Simulation: This step obtains performance metrics, such as response time, resource utilization, and system throughput, by mathematically solving the performance model or executing the simulation code generated from the previous step. The advantage of simulation is its high accuracy, but it is time consuming, since the simulation can be complicated when evaluating real-world systems [5], [54]. In comparison, solving the performance models is faster but less accurate.

Performance models usually have analytical solvers or simulators. For example, *Queueing Networks* have multiple analytical solvers, such as *LQN Solver (LQNS)* [S24], [S43], *JMT Solver* [S12], [S31]. Similarly, *Petri Nets* has the analytical solver *GreatSPN* [S17], [S19], [S29] and the simulator *SimQPN* [S27], [S40], [S42], [S44]. In comparison, the simulation code imitates the execution of a system over a period of time [S18] to obtain performance metrics. The most commonly used simulation tool is *SimuCom* [S9], [S11], [S20], [S25], [S26], which focuses on executing the simulation code from *PCM*. Another simulation tool, *SimuLizar* [S27], is more advanced than *SimuCom*, with higher efficiency and reliability.

There are two model resolution techniques: analytical techniques and simulation techniques. The analytical resolution technique addresses performance models generated by *Model-to-Model* transformation. Analytical techniques are faster but less accurate than simulation. Different performance models usually have their own analytical resolution tools [S8], [S17], [S19], [S22], [S44], [S72], [S90], [S93]. Simulation is the imitation of the operations of a system over a period of time [S18].

Step 5: Result Interpretation and Feedback: As mentioned earlier, nine studies contain a step to interpret the prediction results and provide feedback to stakeholders [S9], [S17], [S24], [S31], [S43]. [S17], [S31] focused on locating performance anomalies in architecture models based on predicted metrics. They detected anomalies by comparing the predictions—response time [S31], resource utilization [S17], [S31], and throughput [S31]—with pre-defined performance specifications. [S9], [S24] mapped the fluctuation in the predicted performance metrics with the respective workload scenarios to understand problems behind the prediction results. [S43] visualized possible architectural refactorings based on the predicted results.

4.3.2 Performance Anti-Patterns D&R Template

Figure 3b illustrates the typical workflow of studies for *Purpose-2: Performance Anti-pattern Detection and Solution*. There are four steps: 1) Model Acquisition, which acquires the architecture model of a system; 2) Performance Metrics Collection, based on either dynamic profiling or using model-based prediction; 3) Detection by Rule Set, which detects performance anti-patterns through a set of detection rules that combine the architecture model and performance metrics; and 4) Refactoring, which suggests refactoring to resolve the identified anti-patterns. Of note, two studies in *Purpose-2* do not have the fourth step, since they only detect performance anti-patterns without providing solutions [S52], [S53].

Step 1: Model Acquisition: The first step is to acquire the architecture model of the system under investigation. 9 related studies retrieved architecture models, such as UML diagrams, *PCM*, and *ADL*, from project designs [S10], [S14], [S47], [S49], [S50], [S51], [S52], [S54], [S94]. The other five studies [S6], [S46], [S48] recovered architecture models by reverse engineering existing software systems.

Step 2: Performance Metrics Collection: This step collects performance metrics by either dynamic profiling or model-based prediction. Six (40%) studies used performance profiling tools, such as *JMeter* [S6], *YourKit* [S6], *Dynamic*

Spotter [S53], or self-built profiling programs [S47], [S48], to collect performance metrics, including response time [S6], [S10], [S14], [S47], [S48], [S49], [S50], [S51], [S52], [S53], resource utilization [S6], [S10], [S14], [S49], [S50], [S51], [S52], [S53], throughput [S6], [S10], [S49], [S50], [S51], [S52], [S53], latency [S6], and data transmission loss ratio. Several studies also collected some metrics relevant to anti-pattern detection, such as the number of calls of a function [S53], the queue length of requests, and the size of transferred messages among system components. The remaining studies [S10], [S49], [S50], [S52] collected performance metrics using model-based prediction methods. These studies followed the process summarized in template-1 (see Section 4.3.1) to transform the architecture model into performance models, *Queueing Network* [S10], [S14], [S50], [S52], or executable simulation code [S49], [S51], [S54].

The collected metrics (profiled or predicted) were used in the same way in follow-up steps. But the anti-patterns may impact how extensively the performance metrics were collected. The detection of *Blob / God Class*, *Pipe and Filter*, *Concurrent Processing Systems*, *Extensive Processing*, *Circuitous Treasure Hunt*, *Empty Semi Trucks*, *One-Lane Bridge*, and *Excessive Dynamic Allocation*, relied on sampling performance metrics during a period of time. That is, researchers only need to sample the average, maximal, or minimal value of related performance metrics during the profiling period [S53], [S55], [S94]. In comparison, the detection of *Traffic Jam*, *The Ramp*, and *More is Less* relied on collecting the trend and fluctuation of the performance metrics over time [S55].

Step 3: Detection by Rule Set: The detection of an anti-pattern usually relied on the combination of multiple rules to match both the dynamic metrics and architectural features [S52], [S53].

Table 7 lists the detection rules from the literature. These rules are of two types: 1) rules based on performance metrics (row 1 to 8), and 2) rules based on architecture features (row 9 to 12). Most rules work by comparing a metric to a threshold. For example, if the average size of transmitted message with a component is less than a threshold, the *MessageSize Rule* is matched. Some rules work by checking the *deviation* of the metric or feature over time or associated with different components. For instance, the *ServiceTh Rule* checks the deviation of service throughput associated with components to identify whether there is a bottleneck in a *Pipe and Filter* architecture. The *DiskMoreUtilized Rule* checks the difference between disk utilization and CPU usage.

Table 8 illustrates the mapping between each anti-pattern and its detection rule(s). Note that four anti-patterns, *Tower of Babel*, *Excessive Dynamic Allocation*, *The Ramp*, and *More is Less* [S1], [S2], [S3], cannot be automatically detected by any rule yet. The other anti-patterns can be automatically detected by one or more rules [S53]. For example, *Blob/God Class* can be detected by matching three rules: 1) *Usage Rule*: a component has a dependency on a large number of other components; 2) *Interaction Rule*: the component is called by a large number of external requests; and 3) *Utilization Rule*: the component has high memory utilization [S53].

Step 4: Architectural Refactoring: Some studies further contributed refactoring solutions to fix the anti-patterns, automatically or with manual effort [S6], [S10], [S14], [S46], [S47], [S48], [S49], [S50], [S51].

TABLE 7: Performance Anti-pattern Detection Rules

Measurement-based Rule	Description
<i>MessageSize Rule</i>	Avg size of message transmitted among a component \leq Threshold
<i>Interaction Rule</i>	Number of external calls of a component \geq Threshold
<i>Utilization Rule</i>	Resource (e.g. CPU, hard disk, memory) \geq Threshold
<i>ServiceTh Rule</i>	Deviation of service throughput associated with different components \geq Threshold
<i>WaitingTime Rule</i>	(1) Avg response time \geq Threshold; (2) Response time increases/decreases continuously.
<i>QueueLength Rule</i>	Queue length of requests \geq Threshold
<i>Unbalanced Rule</i>	Deviation of a resource (e.g. CPU) utilization with different components \geq Threshold
<i>DiskMoreUtilized Rule</i>	(Hard disk utilization - CPU utilization) \geq Threshold
Model-based Rule	Description
<i>RemoteCommunication Rule</i>	Number of communications between local and remote components \geq Threshold
<i>Structural Rule</i>	Scheduling policy of a component is a specific one, e.g. FIFO, LIFO, etc.
<i>Usage Rule</i>	Number of dependencies of a component with components \geq Threshold
<i>Probability Rule</i>	Deviation of the probabilities of executing different processes in a component \geq Threshold

TABLE 8: Performance Anti-patterns vs. Detection Rules

Perf Anti-pattern	Detection Rule
<i>Blob / God Class</i>	<i>Interaction Rule, Utilization Rule, Usage Rule</i>
<i>Concurrent Processing Systems</i>	<i>Utilization Rule, Unbalanced Rule</i>
<i>Pipe and Filter</i>	<i>ServiceTh Rule</i>
<i>Extensive Processing</i>	<i>Utilization Rule, Structural Rule, Probability Rule</i>
<i>Circuitous Treasure Hunt</i>	<i>Interaction Rule, Message Size Rule, DiskMoreUtilized Rule</i>
<i>Empty Semi Truck</i>	<i>Message Size Rule, Interaction Rule, Remote Communication Rule</i>
<i>One-Lane Bridge</i>	<i>Waiting Time Rule, QueueLength Rule</i>
<i>Traffic Jam</i>	<i>Waiting Time Rule</i>

Three performance anti-patterns—*Concurrent Processing Systems*, *Extensive Processing*, and *One-Lane Bridge*—can be automatically refactored. *Concurrent Processing Systems* can be fixed by ensuring a balanced request distribution—i.e. through modifying the probability of branch actions in the problematic components [S50]. *Extensive Processing* can be resolved by using a different scheduling policy that grants short-running requests higher priority. *One-Lane Bridge* can be refactored by optimizing the processing paths to alleviate the congestion of requests [S51].

The refactoring of *three* anti-patterns, *Blob / God Class*, *Circuitous Treasure Hunt*, and *Empty Semi Trucks*, are “semi-automated”. The resolutions of these performance anti-patterns requires expert knowledge for optimizing the logic in the impacted components. The *Blob / God Class* anti-pattern can be resolved by 1) manually optimizing the “Blob / God Class” component by delegating its business logic to the components that it communicates with, and 2) decreasing the number of calls to the *Blob / God Class* component. The *Circuitous Treasure Hunt* anti-pattern can be resolved by 1) decreasing excessive communications with the database, and 2) manually restructuring the database tables to decrease excessive communication.

For the other anti-patterns, we did not find refactoring solutions in the literature. The refactoring of these anti-patterns requires changing the code of the impacted components. This requires case-by-case analysis by an expert. For example, the refactoring of *Pipe and Filter* and *Towel of Babel* anti-patterns requires manual optimization of the business logic of the filter components and the translation component. The refactoring of the *Ramp* anti-pattern requires that developers choose a more efficient data structure to process large amounts of data. The solutions to *Traffic Jam* and *More is Less* anti-patterns need developers to expand resource capacity to improve the processing power of the system.

4.3.3 Profiling and Comparison Template

Figure 3c illustrates the workflow that captures all 12 studies for *Profiling and Comparison* [S5], [S56], [S57], [S58], [S59], [S60], [S61], [S62], [S63], [S64], [S65]. It contains four steps: 1) Alternative Architecture Acquisition, where the architectural solutions for an application are implemented or acquired; 2) Scenario Analysis, where typical workload scenarios for evaluating architectural solutions are identified and analyzed; 3) Dynamic Profiling, where run-time performance metrics in the identified workload scenarios are collected; and 4) Comparison, where the researchers make comparisons of the architectural solutions to provide empirical knowledge.

Step 1: Alternative Architecture Solution Acquisition:

In seven of the studies, researchers implemented alternative architecture solutions for the same application [S56], [S57], [S59], [S60], [S63], [S64]. This approach provided control over the variables and facilitated the comparison of the impact of architectural changes on performance. For example, in [S57], a three-tier web application, *E-Health Record*, was implemented using three different data storage systems. Similarly, the studies [S60], [S63], [S64] explored the same project’s performance in monolithic and micro-service architectures.

In four of the studies, researchers used off-the-shelf systems as their subjects [S5], [S58], [S61], [S62]. In these cases, systems providing similar business functions but employing different architectures were selected for analysis. For instance, the study [S58] compared two different data storage architectures, i.e., *Apache Cassandra* and *MongoDB*. Extending this approach, [S61] contrasted two cloud application architectures—*Ghost* and *WordPress*. Furthermore, [S62] carried out a performance evaluation of three Function as a Service (FaaS) cloud platforms, including Amazon AWS Lambda, IBM, and Azure Cloud Function.

Step 2: Scenario Analysis and Tuning: This step tunes parameters that specify the input scenario or system resource, which impacts the performance of a system. This step is highly domain-specific. For studies that focus on web applications, [S5] tuned the number of concurrent users; while [S59] controlled the number of requests [S59], [S60], [S62], [S63], [S64] and the queue size of messages [S95]. Studies that focused on data storage systems controlled either the workload [S58], [S62] or the type of stored data [S57]. More specifically, [S58] controlled two use case scenarios with the heavy-reading and heavy-writing. [S57] tuned the complexity of the data structures, e.g., tree-like

textual data or images, and different data replication. Tuning system resource, in terms of hardware configurations of the deployment environment, was another way to impact performance. This included 1) CPU allocation [S59], 2) cache [S59], 3) memory allocation [S56], [S59], and disk fraction [S56].

Step 3: Performance Metrics Collection: Researchers collect performance metrics using dynamic profiling tools and benchmarks [S5], [S56], [S58]. Five studies [S60], [S61], [S62], [S63], [S5] used a *Java* profiling tool, *JMeter*, to collect the mean, maximal and minimal of response time and throughput. In [S58], researchers used *Yahoo Cloud Serving Benchmark* (YCSB) to monitor the throughput of *Cassandra* and *MongoDB*. Similarly, [S64] and [S56] used an open source platform, named *BenchFlow*, to automatically execute the performance tests and providing performance insights. [S57] collected the execution time based on the affiliated graphical user interface of the evaluated databases. [S59] directly collected the performance metrics based on the built-in performance logs from *Microsoft Windows XP Profession* operation system.

Step 4: Comparison: This step entails generating empirical knowledge via comparison. In five out of eleven studies, specific architectural solutions consistently outperformed others [S5], [S57], [S60], [S61]. For instance, [S57] found that the *NoSQL* systems—*MongoDB* and *Neo4j*—delivered superior performance to the traditional *MySQL* + *Hibernate* in both document-oriented and graph-oriented data contexts, with *Neo4j* achieving a performance increase of 1.5 times on average, and *MongoDB* accelerating by up to 33 times. Likewise, studies such as [S60], [S64] compared two deployment architectures for web applications and consistently found the modern architecture to be superior to its older counterpart. In [S5], the *Slim Framework* outperformed the *Laravel Framework* across all evaluated scenarios, although the authors conceded that the *Laravel Framework* may be more suitable for large projects with numerous functions and external libraries. Interestingly, [S61] found the WordPress blogging system to offer not only better performance than the Ghost blogging system, but also to be more stable.

In the other 6 studies, different architectural solutions exhibited superior performance in specific scenarios [S56], [S58], [S59], [S62], [S63]. For instance, [S58] discovered that under a heavy-reading workload, *MongoDB* marginally outperformed *Cassandra*, whereas under a heavy-writing workload, *Cassandra* was twice as efficient as *MongoDB*. Similarly, [S56], [S59], [S62], [S63] found that different architectural solutions achieved optimal performance under different deployment environments and workload scenarios.

4.3.4 Study Template for Self-Adaption

Figure 3d illustrates the general workflow of a self-adaptive framework: 1) Alternative Architectural Solution Acquisition, where the alternative architectural solutions are implemented or generated; 2) Performance Metrics Collections, where performance metrics are acquired by dynamic profiling or model-based performance prediction; 3) Decision Making, which determines whether and which alternative architectural solution should switch to; and 4) Execution, where the system enables the selected architectural solution.

Step 1: Alternative Architectural Solution Acquisition:

In 4 studies [S7], [S66], [S67], [S70], researchers manually implemented alternative architectural solutions for the same application. In the other 2 studies [S68], [S69], researchers firstly recovered the software models from an existing software system, and then proposed algorithms to re-configure the system model as alternative architectural solutions.

Step 2: Performance Metrics Collection:

This step is to collect the performance metrics to support decision making in the next step. A studies [S68] used model-based prediction methods to obtain the performance metrics. The other studies [S7], [S66], [S67], [S69], [S70] continuously monitored the performance metrics of a system.

Step 3: Decision Making:

This step is to analyze the performance metrics from the previous step to determine if an adaption is needed. 4 related studies' decision-making was based on switching policies [S7], [S66], [S67], [S70]. The switching policies are in three categories, including: 1) Threshold based policies, such as mean response time [S66], [S70] and data transmission rate [S7]. Once the performance metric exceeds a threshold, the framework triggers a decision to adapt to suitable architecture solution. 2) Failure based policies [S67]. In [S67], the adaptation framework assumes that the system has a failure when the number of arrival requests stays normal or where the throughput drops to zero. And 3) Fluctuation and deviation based policies [S67]. These policies compare the monitored performance metrics with average metrics measured over a time period. Once it exceeds a specified percentage, the framework triggers a decision to adapt to an alternative architecture to better handle the requests.

The other two studies, [S68] and [S69], used a fitness algorithm to select the optimal solution from a set of candidate solutions, instead of relying on determined policies.

Step 4: Execution:

This step enables adaption to the most efficient architectural solution based on the decision from the previous step. In [S67], the self-adaptive framework activates one of three alternatives—namely the *DTC* architecture, the *HS/HA* architecture, and the *LF* architecture. Similarly, in [S96], the self-adaptive framework activates a more suitable data stream processing mechanism, based on the response time or the sampled delay rate. In [S68] and [S7], the self-adaptive framework switches between alternative components deployment patterns of *IoT* systems. In the other three studies [S66], [S69], the self-adaptive framework re-configured system connections [S69] or their resource allocation [S66]. For example, in [S66], [S70], the self-adaptive framework manages the release of resources for system components to enable adaption.

Some steps like Model Acquisition, Architecture Acquisition or Implementation, and Performance Metrics Collection do appear commonly across templates. These commonalities underscore a shared foundation in the process of integrating architecture and performance analysis. As for practitioners, these commonalities serve as fundamental building blocks, enabling them to leverage existing knowledge and expertise across varied research projects. Moreover, the understanding of these shared steps can streamline the initiation of new projects, helping practitioners

to focus on the unique aspects that align with their specific objectives.

RQ3 Summary: This RQ summarized the typical study templates of the four fundamental purposes from the 109 studies, which offer a guide for reproducing existing studies and advancing the state-of-the-art. Figures 3a to 3d show the templates for *Model-based Analysis*, *Performance Anti-Pattern Analysis*, *Profiling and Comparison*, and *Self-Adaption*. Guidance regarding the key analysis components in the templates, including the architecture models, model annotation, performance models, and anti-pattern detection rules are provided.

4.4 RQ-3: Available Tools

RQ3 identifies tools used in the studies to support the analysis of software architecture, performance, and their integration. Overall, in 13 studies, researchers indicated that the proposed approach or study process was fully automated with tool support [S1], [S2], [S15], [S27], [S30], [S31], [S34], [S35], [S40], [S42], [S43], [S44], [S68]. In the other 12 studies, researchers explicitly mentioned that manual effort was required in the study process [S5], [S6], [S21], [S25], [S45], [S49], [S51], [S52], [S53], [S54], [S66], [S70]. In the remaining 61 studies, researchers did not elaborate on the tools for automation or the required manual effort.

We address RQ3 in two sub-RQs:

- *RQ-3.1: What are the architecture and performance model analysis tools used in the literature?* Here we focus on any tools for constructing and analyzing the architecture and performance modules. We organize the tools' information based on their goals for model acquisition, model annotation, model transformation, and model analytical solution or simulation, which mostly align with the template for model-based prediction.
- *RQ-3.2: What are the profiling tools for collecting performance metrics, and what types of metrics are collected?* Here we focus on tools for collecting run-time performance metrics and mapping the tools to the different types of metrics being collected.

4.4.1 RQ-3.1: What are the architecture and performance model analysis tools used in the literature?

Table ?? provides a summary of the tools used for the acquisition, annotation, transformation, and solution or simulation of the architecture and performance models. We would like to clarify that the table only enumerates the tools that we extracted from the 104 studies that are used in different analytical components of the study templates. However, we cannot assume the inter-changeability or interoperability among these tools.

1) Model Acquisition Tools: Nine tools were mentioned in the literature for acquiring software architecture models. They are: *UMLet* [S1], *Enterprise Architect* [S1], *Palladio Bench* [S9], [S25], [S40], [S51], *Eclipse Modeling Framework (EMF)* [S20], [S44], *Understand SciTool* [S32], [S33]. *UMLet* and *MagicDraw* allow users to manually draw UML diagrams. *Enterprise Architect* not only supports manually drawing UML diagrams, but also reverse-engineering UML diagrams from system source code. *Palladio Bench* is the default tool for

reverse-engineering PCM models from a code base. *Enterprise Modeling Framework (EMF)* supports creating both PCM models [S20] and DML models [S44]. *Understand SciTool* can recover the method-level call dependencies from a code base.

2) Model Annotation Tools: We identified three common tools for annotating the architecture model in preparation for performance prediction. They are the *UML Profile for Schedulability, Performance, and Time (UML-SPT)*, *UML Profile for Modeling and Analysis of Real-Time and Embedded systems (UML-MARTE)* [S14], [S17], [S19], [S29], [S35], [S41], [S50], [S52], and *Resource Demanding Service Effect Specification (RDSEFF)* [S9], [S11], [S20], [S25], [S26], [S28], [S45], [S51]. Both *UML-SPT* and *UML-MARTE* can annotate UML diagrams. They can annotate performance information such as process duration, message size, branch action probability, network speed, and workload status. In addition, *UML-MARTE* extends the annotation information to hardware aspects. For instance, *UML-MARTE* can annotate a tag named "HwDevice", for distinguishing whether the hardware device is controlled by or interacts with a component in a software system. Similarly, the tag "HwComponent" annotates hardware components that physically contain other devices [S35]. *RDSEFF* is for annotating the PCM models with information including request size (e.g. I/O resource demand), request type—read or write, the probability branch execution, and the file set that the request is operating on.

3) Model Transformation Tools We found 18 tools or techniques for transforming architecture models into either performance models or textual artifacts for *model-based prediction*. Most tools work on transforming UML diagrams and PCM.

The UML models have nine transformation tools. Seven of them, *XML Metadata Interchange (XMI)* [S74], *Query/View/Transformation (QVT) language* [S41], *Janus Transformation Language (JTL)* [S35], *PRIMA-UML* [S52], and *PUMA PRIMA-UML* [S29], [S52] transform UML diagrams into *Queuing Networks*. Four tools are based on declarative transformation languages, which iteratively match the components in the source (i.e., architecture) model to the target (i.e., performance) model. *SPE Editor (SPE-ED)* [S36] not only supports the transformation from annotated UML diagrams into *Queuing Networks*, but also support solving the generated *Queuing Network* models with analytical techniques or with simulation [S8]. Finally, *ArgoSPE* [S17] is a tool to translate UML diagrams into *Stochastic Petri Nets*. There are five commonly used transformation tools for the PCM models. *PCM2LQN* [S28] and *PCM2QPN* [S28] transform PCM models to *Layered Queuing Networks* and *Queueing Petri Nets* respectively. *KLAPER* [S40] is used to transform PCM into *Queueing Petri Nets*. *PCM2SimuCom* [S11], [S20], [S25], [S28], [S51] can automatically transform PCM into simulation code.

The remaining five tools are for transforming other architectural models. *DML2SimQPN* [S44] and *DML2LQNS* [S44] are for transforming DML models into *Layered Queuing Network (LQN)*, *Queueing Petri Net (QPN)*, or simulation code. For transforming *Æmilía*, a specific ADL, [S49] transforms it to simulation code by using a declarative transformation language, *Gra2Mol*, and to a performance model by using the tool, *TwoTower*. For *Calling Dependency*

Graph model.

4) Model Solution or Simulation Tools: We identified 13 tools for solving a performance model or executing simulation code to obtain performance metrics for a system.

There are five different tools for the analytical solution of *Queueing Network* models. *LQNS* is the default solver for a *Layered Queueing Network (LQN)*; it is used in [S14], [S24], [S43]. *DLV Solver* and *OPERA Solver* also solve LQN performance models. *Java Modeling Tools (JMT)* is also a popular tool to approximate solutions for *Queueing Network* models, used in [S10], [S12], [S31], [S50]. Finally, *SHARPE* is the default solver for *Product Form Queueing Network*.

There are three tools for solving or simulating *Petri Nets* models. *QPN Solver* is for basic *Queueing Petri Nets* [S42]. *GreatSPN* solves *Generalized Stochastic Petri Net (GSPN)* models [S17], [S29]. *SimQPN* performs discrete-event simulation for *Queueing Petri Net* models, used in [S15], [S27], [S42], [S44].

The remaining five tools are for executing simulation code transformed from architecture models. *SimuCom* is the most frequently used simulator for *PCM* models [S9], [S11], [S20], [S25], [S28], [S51], [S84]. *SimuLizar* is similar to *SimuCom* but provides higher efficiency, scalability, and reliability. *SimSo* [S23] is for the simulation of the *Architecture and Analysis Design Language (AADL)*. *SMTQA* executes simulation code transformed from a *UML Use Case Diagram*. *SLAstic.SIM* executes simulation code transformed from *Business Process Execution Language (BPEL)*.

4.4.2 RQ-4.2: What are the profiling tools for collecting performance metrics, and what types of metrics are collected?

All studies collected performance metrics. 33 out of 109 studies clearly mentioned the profiling tools that they used. In several studies, researchers did not use a dedicated profiling tool, but built their own profiling programs or collected performance metrics from operating system logs.

Six types of performance metrics were collected including 1) response time; 2) resource utilization; 3) throughput; 4) processing time; 5) latency; and 6) data transmission loss. The three most commonly collected metrics are response time, resource utilization, and throughput. This finding is consistent with results from previous studies [6], [9].

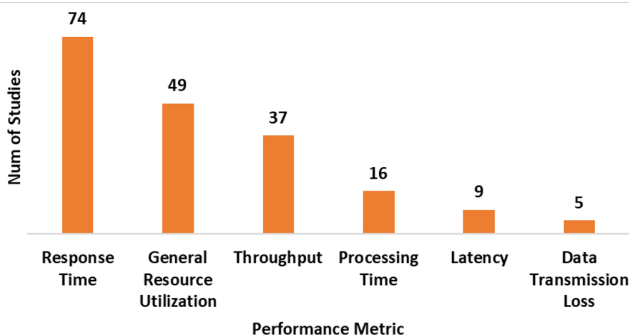


Fig. 4: Performance Metrics vs. Number of Studies

The array of studies reviewed in this paper employed an assortment of tools to gather metrics and oversee the performance of a broad variety of applications. For Java projects, profiling was often carried out using JMeter [S5], [S6], [S27], [S35], [S60], [S61], [S62], [S63], Java Flight

Recorder [S25], and YourKit [S6], [S32]. JMeter, being used in several studies, is capable of generating JUnit test cases automatically for more detailed profiling results. The Java Flight Recorder is notably light-weight, producing almost no performance overhead, whereas YourKit extends support for performance profiling to both Java and .NET projects. Java Management eXtension (JMX) has been employed for monitoring the performance of service-oriented Java applications.

For the performance monitoring of cloud-based applications with distributed clusters like Apache Hadoop, MRPerf [S4] and CloudWatch [S66] have proven to be useful. In [S97], the tools ContiPerf and JUnitPerf were used to turn JUnit 4 tests into performance tests to measure performance within existing tests, facilitating the smooth integration of performance testing into the unit testing process.

In terms of executing performance experiments and overseeing the performance of database systems like Cassandra and MongoDB, open-source frameworks such as Prometheus [S34], BenchFlow [S56], [S64], and Yahoo Cloud Serving Benchmark (YCSB) [S58] have been employed. It's crucial to note that these tools have demonstrated versatility in their applicability, extending to various contexts and purposes beyond those described in this review.

RQ4 Summary: This RQ offers a catalog of tools for architecture analysis and for performance profiling curated from the 109 studies. This can help researchers and practitioners identify resources for their own research, or for reproduction studies.

5 DISCUSSION OF FUTURE DIRECTIONS

5.1 Lack of Reproducibility:

The lack of reproducible research has raised concerns in the software engineering community [55], [56]. Although the literature that integrates software architecture and performance analysis has great potential, reproducibility of related studies remains a concern due to the lack of easy to reuse tools, benchmark datasets, and hands-on evaluations of these resources. Practical, hands-on evaluations of the tools and instruments could augment our study and provide more pragmatic insights into their effectiveness and usability. However, conducting such experiments involves many complexities, including the availability and usability of tools, and could constitute a series of separate studies in itself. This, we suggest, would be a valuable future direction for research in this field.

Therefore, increasing the reproducibility of related studies should be prioritized in future research. The community should emphasize the importance of including replication packages when publishing new studies. This could help to increase the adoption of the proposed techniques. In addition, replication studies should be valued as an important type of scientific study in the community. This enables solid validation of new techniques by allowing study replication in multiple sites and by enforcing cross-project comparisons with the state-of-the-art.

TABLE 9: Profiling Tools vs. Performance Metrics

Tool Name	Response Time	Resource Utilization	Throughput	Processing Time	Data Transmission Loss	Latency
<i>JMeter</i>	[S6], [S35], [S60], [S61], [S62], [S63], [S83]	[S6], [S83]	[S5], [S6]		[S5]	[S6]
<i>Java Flight Recorder</i>	[S25]	[S25]				
<i>YourKit</i>	[S6]	[S6]	[S6]	[S32]		[S6]
<i>JMX</i>	[S75]	[S75]				
<i>MRPerf</i>	[S4]	[S4]				
<i>CloudWatch</i>	[S66]	[S66]	[S66]			
<i>ContiPerf</i>	[S97]	[S97]				
<i>JUnitPerf</i>	[S97]	[S97]				
<i>WebGUI</i>	[S87]					
<i>Nagios & Cacti</i>		[S95]		[S95]		
<i>SNMP</i>	[S75]	[S75]				
<i>Prometheus</i>	[S34]	[S34]				[S34]
<i>BenchFlow</i>		[S56]	[S56]			
<i>YCSB</i>		[S56]	[S58]			

5.2 Challenges with Emerging Software System Domains:

Most studies focused on software systems in traditional domains as their evaluation subjects, such as E-commerce systems. Only four studies focused on cyber-physical system [S22], real-time embedded system [S23], [S92], and Internet of Things system [S8], which posed challenges to the architecture and performance modeling resulted from the dynamism of the physical environment.

New application domains such as block-chain systems, Web 3.0 systems, and virtual reality systems, may pose new performance challenges that are not sufficiently addressed in existing literature. For example, when building an architecture model of a blockchain system, one has to consider the consensus protocol that manages the transaction order, such as *Proof of Work (PoW)* and *Practical Byzantine Fault Tolerance (PBFT)*. In addition, new performance profiling techniques are needed to monitor performance metrics such as state updating time, consensus-cost time, and contract execution time [57], which cannot be collected using existing profiling instruments. As discussed in RQ4, most existing performance profiling instruments focus on measures of response time, throughput, and resource utilization, which may not meet the analysis needs of projects in specific domains.

Therefore, another future direction is to develop new architecture modeling and performance analysis techniques that can address challenges in emerging software system domains.

5.3 Limitations of the Evaluation Methodology:

A total of 30 studies only conducted a case study to show the feasibility of the proposed approach, without providing any specific evaluation criteria [S1], [S4], [S6], [S9], [S11], [S13], [S14], [S15], [S18], [S19], [S20], [S21], [S24], [S25], [S26], [S27], [S28], [S36], [S42], [S44], [S45], [S47], [S68], [S84], [S85], [S86], [S88], [S94], [S98], [S99]. And for studies that did focus on certain evaluation criteria, they mostly focused on technical effectiveness, such as accuracy of performance prediction, precision and recall of anti-pattern detection.

In addition, there is a lack of involvement from industrial practitioners in the evaluations. As discussed in RQ5, only 29% of the studies evaluated real-world systems [?], [S4], [S6], [S7], [S13], [S15], [S16], [S17], [S19], [S20], [S24], [S27], [S29], [S32], [S33], [S41], [S42], [S46], [S48], [S57], [S58], [S68], [S72], [S84], [S96], [S99], [S100]. In addition, these studies were mostly research-oriented—executing test cases of the system in lab environments to collect performance metrics. Evaluations that involve industrial practitioners can greatly

benefit the research, promoting the adoption of the research outcomes in practical settings.

5.4 Lack of ML/AI Techniques:

Machine learning (ML) and artificial intelligence (AI) techniques have been increasingly adopted in software engineering activities [58], [59], [60]. However, in the literature we reviewed, only three studies [S7], [S68], [S99] leveraged machine learning techniques. These studies applied machine learning for selecting influential performance factors and optimizing architectural solutions for performance self-adaptation.

Despite the rising prominence of ML/AI, its application in the area of software performance engineering is not sufficiently explored, due, in part, to the inherent complexities of these systems. They often lack the detailed information necessary for effective performance modeling, presenting unique challenges that traditional software architectures do not face.

In *Model-based Prediction* studies, intricate architecture modeling, annotation, transformation, and solution processes often demand considerable manual effort. Researchers could explore leveraging ML and AI techniques to alleviate these challenges. In performance anti-pattern detection studies, detection rules are primarily heuristic, relying on practitioners' knowledge and expertise. The adoption of ML and AI techniques could enhance the accuracy of detection in practical settings. Similarly, *Self-Adaption* studies, which largely depend on rules and thresholds for switching policies, could benefit from the integration of ML and AI techniques for more intelligent adaptation.

5.5 Tentative Future Research Road-Map

Based on the above discussion of limitations with existing studies, we would like to propose a tentative road map for researchers and practitioners to refer to and build upon:

Reproducibility: Given the large quantity and high complexity of techniques and tools for integrating software architecture and performance analysis, reproducibility poses unique challenges to this direction. Here, we would like to propose the following research questions, derived from our research results:

- 1) To what level the tools curated from RQ4 is actually available, supported, and usable for future research? It calls for a more in-depth investigation to actually retrieve, try, and test the tools as instructed in the original studies.

- 2) To what level the datasets used in the evaluation of the 104 studies are actually available and useful to the research community? As discussed in RQ5, most of the evaluation subjects are lab implementation and simulated systems. A small portion is real-world software projects. Examination of the availability and quality of the used dataset could serve for the construction of a shared benchmark, which is critical to reproducibility.
- 3) To what extent the claimed results in existing studies can be reproduced, upon the confirmation of tool and data availability and quality? There is no guarantee that the same results could be achieved when repeated by a different research team. In particular, software architecture analysis is doomed subjective to an analyst's expertise; while dynamic performance profiling is volatile to the environment. These factors together add more challenges to reproducing prior results, underscoring its importance.
- 4) To what extent that the various tools and techniques proposed in different studies are inter-operatable with each other? As presented in RQ3, different study templates share the same analysis components. As shown in RQ4, different concrete tools and techniques are employed for the same analytical components of the template in different studies. It remains largely open how the tools summarized in our RQ4 from different studies are actually interoperable with each other when plugged into the templates. This investigation would largely increase the flexibility and potential of architecture and performance analysis.

Emerging Domains: As presented in RQ1 (see Figure 2), researchers have attempted to provide customized solutions that meet the various demands of model-based analysis, resulting from specific architecture styles, specific system features, and certain analytical capabilities in special environments. However, with the emerging software applications domains encompassing all aspects of society, there is a significant gap between existing software architecture and performance analysis techniques and these emerging domains. Therefore, we would like to point to the following research questions for the community:

- 1) How well do existing software architecture and performance analysis techniques apply to emerging domains, such as Web 3.0, Blockchain, Virtual Reality, etc.? This could be investigated based on the structure presented in RQ1 (see Figure 2). First, do the four fundamental purposes apply to emerging domains? Additionally, do the specific focuses of model-based performance analysis techniques cover the challenges of emerging domains?
- 2) What are the unique challenges that are faced with the emerging domains? How existing techniques could be customized and enhanced to overcome them? We envision that emerging domains pose unique challenges that identify the limitations in existing studies. It would require multi-disciplinary research to picture their unique challenges for architecture and performance analysis.

Practical Adoption: With the narrow focused evaluation subjects and factors, the practical adoption of the architecture and performance analysis techniques remains questionable. Therefore, we would like to add the following research question to the future road-map:

What are the most significant challenges and impacting factors on the practical usage of architecture and performance analysis techniques? We assume that usability, learning curve, and cost-and-effectiveness all play critical roles in this realm. However, the most effective approach is to conduct human subject studies to reveal what challenges are actually encountered that will prevent the practical adoption. On top of that, engaging participants from practice to provide direct input through surveys, interviews, and focus groups would also be profoundly valuable.

AI Enabled Analysis: AI has gained widespread application in various domains, as well as in facilitating software engineering practice. Unique challenges arise in the application of AI for architecture and performance analysis. We would like to propose the following research questions that are customized to the unique challenges of this direction:

- 1) How to shift architecture analysis' reliance on experts' expertise to AI? It is widely recognized that architecture analysis tends to be manual, labor-intensive, and error-prone. Therefore its success and effectiveness heavily rely on the experience and expertise of the analyst. It is yet an open question if AI would be able to provide a mitigation to this reliance, and to what extent.
- 2) How to leverage AI for capturing the complexity and uncertainty in performance analysis? From our survey, we learned that performance analysis faces the unique challenges of uncertainty resulting from huge configuration space, composed of parameters and inter-play from the system as well as from the execution environment. Capturing the underlying complexity of such uncertainties often requires statistical analysis of the system's execution profile and environment. AI holds great promise in this area given that it is built upon powerful statistical analysis itself.

We acknowledge that the above road map is by no means comprehensive and authoritative. However, we believe that it provides a starting point for the community to build upon.

6 RELATED WORK

There are existing literature review studies that focus on software architecture and performance [6], [7], [8], [9], [10], [11], [12]. Here we discuss how our study differs from existing work.

Aleti *et al.* [6] presented a comprehensive systematic literature review on 188 studies (published between 1992 and 2011) that share the theme of software architecture optimization. The authors derived a taxonomy of architecture optimization concerns from these 188 studies, including: 1) performance (84 papers), 2) financial costs (74 papers), 3) reliability (71 papers), 4) availability (25 papers), 5) energy saving (18 papers), and 6) safety (4 papers). Similarly, the study by Falessi *et al.* [12] summarized 11 empirical studies

(between 2006 to 2008) that evaluate the methods, techniques, and tools that support architecture design, analysis, and review. This work focused on 1) the activities in software development (e.g., documentation design or maintenance evaluation), 2) research questions, and 3) study process of these 11 studies. Our work, compared with [6] and [12], specifically focuses on the integration of architecture and performance analysis. Those papers have much broader scope, considering other quality concerns to motivate architecture optimization. In addition, our work focus on literature in the most recent decade, compared to the range of 1992 to 2011 in Aleti *et al.*'s work [6] and the range of 2006 to 2008 in Falessi *et al.* [12].

A number of secondary studies have already summarized the state-of-the-art in model-based performance prediction [7], [8], [9]. Balsamo *et al.* [9] conducted a literature review that surveyed 15 studies in model-based performance prediction. Their study focused on three aspects: 1) the identified software models and performance models; 2) the application in the software development life-cycle; and 3) the degree of automation. They found that *UML Diagrams* are the most commonly used software architecture model; and *Queuing Network* is the most preferred performance model. Most studies applied model-based prediction in the design phase. And, most of the approaches have a high degree of automation. The study [7] is a survey of performance models used for distributed software system architectures. It revealed that performance prediction approaches are mostly based on *Queuing Networks*, *Petri Nets*, *Queuing Petri Nets*, and *Hierarchical Performance Model*. Similarly, [8] is a survey of 17 performance prediction studies, with a focus on component-based systems. This survey revealed that *UML Diagrams* are the most commonly used software architecture model; and *Queuing Network* is the most preferred performance model, which is consistent with our findings.

Two surveys focused on self-adaptive systems. Arcelli presented a survey of 10 studies that apply performance modeling and assessment for supporting self-adaptive software systems [10]. For each study, this survey analyzed the time of application (e.g., design time or running time), the architecture model (e.g., *PCM*), the performance analysis methodology (e.g., analytical solution or simulated solution), as well as the applicability of the proposed tools. This survey pointed out that the state-of-art does not provide automated, performance-driven decision-making towards convenient alternative self-adaptive system architectures, by design-space exploration. This finding later motivated the author of the survey to propose a self-adaptive approach that is based on a genetic algorithm [S68], which is included in our literature. Becker *et al.* presented another survey focusing on six self-adaptive approaches [11]. This study classified the six approaches based on their use scenarios as in design-time or run-time. Our work differs from these studies [7], [8], [9], [10], [11] in three ways: 1) Our goal is to understand how software architecture and performance analysis can be integrated. Model-based prediction and self-adaptive systems are just two of the four study purposes that emerged from our literature set; 2) The model-based prediction literature surveyed in our work is more recent and comprehensive. We reviewed the literature between 2010

and 2020, while prior work mostly focused on studies before 2010. And our work is also more comprehensive (with 94 studies); 3) For self-adaptive systems, existing surveys only find approaches that rely on a set of pre-defined policies to switch to an alternative architecture. However, we found two recent studies [S7], [S68] where researchers leverage machine learning techniques in making decisions for selecting suitable alternative architecture solutions. And, finally 4) We summarized the reference study templates of four research purposes, the available tools and instruments, and the evaluation methodology and domains for studies that integrate software architecture and performance analysis.

7 CONCLUSION

This paper presented a literature review of 109 studies that integrate software architecture and performance analysis. We addressed five research questions that revealed the study purposes and the facilitated development activities (RQ1), the reference research templates for different research purposes (RQ2), as well as the available tools and instruments (RQ3). The findings could provide reference and guidance for practitioners and researchers who are interested in software architecture, software performance, and the integration of the two.

More specifically, in RQ1, we found that the majority of the surveyed studies focused on leveraging architecture modeling for performance analysis, while the rest addressed performance anti-pattern, performance profiling, and construction of self-adaptive architectures. We also discovered that the majority of the studies focused on addressing performance in the design phase. In RQ2, we summarized typical research templates of the four fundamental objectives of integrating software architecture and performance analysis, which could serve as quick guidance for researchers and practitioners to reproduce or advance existing studies. In RQ3, we summarized a catalog of the tools used in performance and architectural analysis, focusing on architecture modeling and performance profiling separately. This provides support for researchers and practitioners to identify tools that are relevant to their needs.

We also discussed potential future research opportunities based on limitations observed in existing studies. These include: 1) Lack of reproducibility—most studies lack replication package with available tools and dataset to reproduce their studies; 2) Challenges with emerging software domains—most studies focus on projects that are from traditional domains, such as e-commerce; while the focus on emerging domains such as Blockchain still needs more research; 3) Limitations with the evaluation methodology—most studies did not evaluate their techniques based on factors that impact the practical usage, such as usability; and 4) Lack of ML/AI techniques—modern machine learning and artificial intelligence techniques are not sufficiently used in the integration of architecture and performance modeling and analysis. Based on the identified limitations, we proposed a tentative future research road map for researchers and practitioners to refer to and build upon.

REFERENCES

- [1] Len Bass, Paul Clements, and Rick Kazman. *Software architecture in practice*. Addison-Wesley Professional, 2003.
- [2] Xavier Franch and Pere Botella. Putting non-functional requirements into software architecture. In *Proceedings Ninth International Workshop on Software Specification and Design*, pages 60–67. IEEE, 1998.
- [3] Shahed Zaman, Bram Adams, and Ahmed E Hassan. A qualitative study on performance bugs. In *2012 9th IEEE working conference on mining software repositories (MSR)*, pages 199–208. IEEE, 2012.
- [4] Adrian Nistor, Tian Jiang, and Lin Tan. Discovering, reporting, and fixing performance bugs. In *2013 10th working conference on mining software repositories (MSR)*, pages 237–246. IEEE, 2013.
- [5] Steffen Becker. *Coupled model transformations for QoS enabled component-based software design*. PhD thesis, Universität Oldenburg, 2008.
- [6] Aldeida Aleti, Barbora Buhnova, Lars Grunske, Anne Koziulek, and Indika Meedeniya. Software architecture optimization methods: A systematic literature review. *IEEE Transactions on Software Engineering*, 39(5):658–683, 2012.
- [7] Stephen Olatunde Olabiyisi, Elijah Olusayo Omidiora, Faith-Michael E Uzoka, Mbarika Victor, and Boluwaji A Akinuwaesi. A survey of performance evaluation models for distributed software system architecture. In *World Congress on Engineering 2012. July 4-6, 2012. London, UK.*, volume 2186, pages 35–43. International Association of Engineers, 2010.
- [8] Steffen Becker, Lars Grunske, Raffaella Mirandola, and Sven Overhage. Performance prediction of component-based systems. In *Architecting Systems with Trustworthy Components*, pages 169–192. Springer, 2006.
- [9] Simonetta Balsamo, Antinisca Di Marco, Paola Inverardi, and Marta Simeoni. Model-based performance prediction in software development: A survey. *IEEE Transactions on Software Engineering*, 30(5):295–310, 2004.
- [10] Davide Arcelli. Exploiting queuing networks to model and assess the performance of self-adaptive software systems: a survey. *Procedia Computer Science*, 170:498–505, 2020.
- [11] Matthias Becker, Markus Luckey, and Steffen Becker. Model-driven performance engineering of self-adaptive systems: a survey. In *Proceedings of the 8th international ACM SIGSOFT conference on Quality of Software Architectures*, pages 117–122, 2012.
- [12] Davide Falesti, Muhammad Ali Babar, Giovanni Cantone, and Philippe Kruchten. Applying empirical software engineering to software architecture: challenges and lessons learned. *Empirical Software Engineering*, 15(3):250–276, 2010.
- [13] Barbara Kitchenham, Pearl Brereton, Zhi Li, David Budgen, and Andrew Burn. Repeatability of systematic literature reviews. In *15th Annual Conference on Evaluation & Assessment in Software Engineering (EASE 2011)*, pages 46–55. IET, 2011.
- [14] Barbara Ann Kitchenham, David Budgen, and Pearl Brereton. *Evidence-based software engineering and systematic reviews*, volume 4. CRC press, 2015.
- [15] David Garlan. Software architecture. *Wiley Encyclopedia of Computer Science and Engineering*, 2007.
- [16] Philippe B Kruchten. The 4+ 1 view model of architecture. *IEEE software*, 12(6):42–50, 1995.
- [17] Christine Hofmeister, Robert L Nord, and Dilip Soni. Describing software architecture with uml. In *Working Conference on Software Architecture*, pages 145–159. Springer, 1999.
- [18] Claudio Riva, Petri Selonen, Tarja Systa, and Jianli Xu. Uml-based reverse engineering and model analysis approaches for software architecture maintenance. In *20th IEEE International Conference on Software Maintenance, 2004. Proceedings.*, pages 50–59. IEEE, 2004.
- [19] Connie U Smith and Lloyd G Williams. Software performance engineering. In *UML for Real*, pages 343–365. Springer, 2003.
- [20] Connie U Smith and Lloyd G Williams. *Performance solutions: a practical guide to creating responsive, scalable software*, volume 1. Addison-Wesley Reading, 2002.
- [21] Günter Haring, Christoph Lindemann, and Martin Reiser. *Performance evaluation: Origins and directions*. Springer, 2003.
- [22] Gianfranco Balbo. *Approximate solutions of queueing network models of computer systems*. Purdue University, 1979.
- [23] Shahed Zaman, Bram Adams, and Ahmed E Hassan. Security versus performance bugs: a case study on firefox. In *Proceedings of the 8th working conference on mining software repositories*, pages 93–102. ACM, 2011.
- [24] Vittorio Cortellessa. How far are we from the definition of a common software performance ontology? In *Proceedings of the 5th international workshop on Software and performance*, pages 195–204, 2005.
- [25] Murray Woodside, Greg Franks, and Dorina C Petriu. The future of software performance engineering. In *Future of Software Engineering (FOSE’07)*, pages 171–187. IEEE, 2007.
- [26] Connie U Smith. Introduction to software performance engineering: Origins and outstanding problems. In *International School on Formal Methods for the Design of Computer, Communication and Software Systems*, pages 395–428. Springer, 2007.
- [27] Elaine J Weyuker and Filippos I Vokolos. Experience with performance testing of software systems: issues, an approach, and case study. *IEEE transactions on software engineering*, 26(12):1147–1156, 2000.
- [28] Alberto Avritzer and ER Weyuker. The automatic generation of load test suites and the assessment of the resulting software. *IEEE Transactions on Software Engineering*, 21(9):705–716, 1995.
- [29] Mrs Charmy Patel and Ravi Gulati. Software performance testing tools—a comparative analysis. *Int. J. Eng. Res. Dev.*, 3(9):58–61, 2012.
- [30] Alberto Avritzer and Elaine J Weyuker. Investigating metrics for architectural assessment. In *Proceedings Fifth International Software Metrics Symposium. Metrics (Cat. No. 98TB100262)*, pages 4–10. IEEE, 1998.
- [31] Khaled M Mustafa, Rafa E Al-Qutaish, and Mohammad I Muhairat. Classification of software testing tools based on the software testing methods. In *2009 Second International Conference on Computer and Electrical Engineering*, volume 1, pages 229–233. IEEE, 2009.
- [32] 7 Best Performance Testing Tools in 2020 | Load Testing Tools, August 2020.
- [33] Will Sobel, Shanti Subramanyam, Akara Sucharitakul, Jimmy Nguyen, Hubert Wong, Arthur Klepchkov, Sheetal Patil, Armando Fox, and David Patterson. Cloudstone: Multi-platform, multi-language benchmark and measurement tools for web 2.0. In *Proc. of CCA*, volume 8, page 228, 2008.
- [34] Nishi Srivastava, Ujjwal Kumar, and Pawan Singh. Software and performance testing tools. 2021.
- [35] Luis San Andrés, Stephen Phillips, and Dara Childs. A water-lubricated hybrid thrust bearing: measurements and predictions of static load performance. *Journal of Engineering for Gas Turbines and Power*, 139(2), 2017.
- [36] Philip Markey and Gary Clync. A performance analysis of ws-*(soap) and restful web services for implementing service and resource orientated architectures. *The IT&T*, page 93, 2013.
- [37] Emily H Halili. *Apache JMeter: A practical beginner's guide to automated testing and performance measurement for your websites*. Packt Publishing Ltd, 2008.
- [38] Srikrishna Prasad and SB Avinash. Smart meter data analytics using opentsdb and hadoop. In *2013 IEEE Innovative Smart Grid Technologies-Asia (ISGT Asia)*, pages 1–6. IEEE, 2013.
- [39] RPT Applying. Using rational performance tester version 7.
- [40] M Raveendra Kumar and R Hari Kumar. Architectural refactoring of a mission critical integration application: A case study. In *Proceedings of the 4th India Software Engineering Conference*, pages 77–83, 2011.
- [41] Omid Bushehrian. Automatic actor-based program partitioning. *Journal of Zhejiang University SCIENCE C*, 11(1):45–55, 2010.
- [42] Sara Mahdavi-Hezavehi, Vinicius HS Durelli, Danny Weyns, and Paris Avgeriou. A systematic literature review on methods that handle multiple quality attributes in architecture-based self-adaptive systems. *Information and Software Technology*, 90:1–26, 2017.
- [43] SI Mal’kovskii, Aleksei A Sorokin, Sergey P Korolev, AA Zatsarinnyi, and GI Tsoi. Performance evaluation of a hybrid computer cluster built on ibm power8 microprocessors. *Programming and Computer Software*, 45(6):324–332, 2019.
- [44] Lucia Kapova, Babora Buhnova, Anne Martens, Jens Happe, and Ralf Reussner. State dependence in performance evaluation of component-based software systems. In *Proceedings of the first joint WOSP/SIPEW international conference on Performance engineering*, pages 37–48, 2010.
- [45] Holger Eichelberger, Cui Qin, Klaus Schmid, and Claudia Niederée. Adaptive application performance management for big data stream processing. *Softwaretechnik Trends*, 35(3):35–37, 2015.
- [46] Michele Guerriero, Saeed Tajfar, Damian A Tamburri, and Elisabetta Di Nitto. Towards a model-driven design tool for big data architectures. In *Proceedings of the 2nd international workshop on BIG data software engineering*, pages 37–43, 2016.

- [47] Aniello Castiglione, Marco Gribaudo, Mauro Iacono, and Francesco Palmieri. Modeling performances of concurrent big data applications. *Software: Practice and Experience*, 45(8):1127–1144, 2015.
- [48] Simon Spinner, Jürgen Walter, and Samuel Kounev. A reference architecture for online performance model extraction in virtualized environments. In *Companion Publication for ACM/SPEC on International Conference on Performance Engineering*, pages 57–62, 2016.
- [49] Lorenzo Pagliari, Mirko D’Angelo, Mauro Caporuscio, Raffaella Mirandola, and Catia Trubiani. To what extent formal methods are applicable for performance analysis of smart cyber-physical systems? In *Proceedings of the 13th European Conference on Software Architecture-Volume 2*, pages 139–146, 2019.
- [50] Xin Zhou, Yuqin Jin, He Zhang, Shanshan Li, and Xin Huang. A map of threats to validity of systematic literature reviews in software engineering. In *2016 23rd Asia-Pacific Software Engineering Conference (APSEC)*, pages 153–160. IEEE, 2016.
- [51] Connie U Smith and Lloyd G Williams. Software performance antipatterns. In *Proceedings of the 2nd international workshop on Software and performance*, pages 127–136, 2000.
- [52] Connie U Smith and Lloyd G Williams. New software performance antipatterns: More ways to shoot yourself in the foot. In *Int. CMG Conference*, pages 667–674, 2002.
- [53] Connie U Smith and Lloyd G Williams. More new software performance antipatterns: Even more ways to shoot yourself in the foot. In *Computer Measurement Group Conference*, pages 717–725. Citeseer, 2003.
- [54] Simonetta Balsamo, Moreno Marzolla, Antinisa Di Marco, and Paola Inverardi. Experimenting different software architectures performance techniques: a case study. In *Proceedings of the 4th international workshop on Software and performance*, pages 115–119, 2004.
- [55] Carl Boettiger. An introduction to docker for reproducible research. *ACM SIGOPS Operating Systems Review*, 49(1):71–79, 2015.
- [56] Lech Madeyski and Barbara Kitchenham. Would wider adoption of reproducible research be beneficial for empirical software engineering research? *Journal of Intelligent & Fuzzy Systems*, 32(2):1509–1521, 2017.
- [57] Peilin Zheng, Zibin Zheng, Xiapu Luo, Xiangping Chen, and Xuanzhe Liu. A detailed and real-time performance monitoring framework for blockchain systems. In *2018 IEEE/ACM 40th International Conference on Software Engineering: Software Engineering in Practice Track (ICSE-SEIP)*, pages 134–143. IEEE, 2018.
- [58] Du Zhang and Jeffrey JP Tsai. Machine learning and software engineering. *Software Quality Journal*, 11(2):87–119, 2003.
- [59] Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. Software engineering for machine learning: A case study. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*, pages 291–300. IEEE, 2019.
- [60] Hironori Washizaki, Hiromu Uchida, Foutse Khomh, and Yann-Gaël Guéhéneuc. Studying software engineering patterns for designing machine learning systems. In *2019 10th International Workshop on Empirical Software Engineering in Practice (IWESEP)*, pages 49–495. IEEE, 2019.
- [S5] Andri Sunardi et al. Mvc architecture: A comparative study between laravel framework and slim framework in freelancer project monitoring system web based. *Procedia Computer Science*, 157:134–141, 2019.
- [S6] Catia Trubiani, Alexander Bran, André van Hoorn, Alberto Avritzer, and Holger Knoche. Exploiting load testing and profiling for performance antipattern detection. *Information and Software Technology*, 95:329–345, 2018.
- [S7] Javier Cámara, Henry Muccini, and Karthik Vaidhyanathan. Quantitative verification-aided machine learning: A tandem approach for architecting self-adaptive iot systems. In *Proceedings of the 17th International Conference on Software Architecture*, pages 11–22. IEEE, 2020.
- [S8] Connie U Smith and Catalina M Lladó. Spe for the internet of things and other real-time embedded systems. In *Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering Companion*, pages 227–232, 2017.
- [S9] Lucia Happe, Barbora Buhnova, and Ralf Reussner. Stateful component-based performance models. *Software & Systems Modeling*, 13(4):1319–1343, 2014.
- [S10] Davide Arcelli, Vittorio Cortellessa, and Catia Trubiani. Performance-based software model refactoring in fuzzy contexts. In *International Conference on Fundamental Approaches to Software Engineering*, pages 149–164. Springer, 2015.
- [S11] Andreas Brunnert and Helmut Krcmar. Continuous performance evaluation and capacity planning using resource profiles for enterprise applications. *Journal of Systems and Software*, 123:239–262, 2017.
- [S12] Leire Etxeberria, Catia Trubiani, Vittorio Cortellessa, and Goiuria Sagardui. Performance-based selection of software and hardware features under parameter uncertainty. In *Proceedings of the 10th international ACM Sigsoft Conference on Quality of Software Architectures*, pages 23–32, 2014.
- [S13] Fabian Brosig, Nikolaus Huber, and Samuel Kounev. Architecture-level software performance abstractions for online performance prediction. *Science of Computer Programming*, 90:71–92, 2014.
- [S14] Catia Trubiani, Antinisa Di Marco, Vittorio Cortellessa, Nariman Mani, and Dorina Petriu. Exploring synergies between bottleneck analysis and performance antipatterns. In *Proceedings of the 5th ACM/SPEC International Conference on Performance engineering*, pages 75–86, 2014.
- [S15] Simon Eismann, Jürgen Walter, Jóakim von Kistowski, and Samuel Kounev. Modeling of parametric dependencies for performance prediction of component-based software systems at run-time. In *Proceedings of the 15th IEEE International Conference on Software Architecture*, pages 135–13509. IEEE, 2018.
- [S16] Marco Gribaudo, Mauro Iacono, and M Kiran. A performance modeling framework for lambda architecture based applications. *Future Generation Computer Systems*, 86:1032–1041, 2018.
- [S17] Elena Gómez-Martínez, Rafael Gonzalez-Cabero, and José Merseguer. Performance assessment of an architecture with adaptative interfaces for people with special needs. *Empirical Software Engineering*, 19(6):1967–2018, 2014.
- [S18] Dominik Werle, Stephan Seifermann, and Anne Koziulek. Data stream operations as first-class entities in component-based performance models. In *European Conference on Software Architecture*, pages 148–164. Springer, 2020.
- [S19] Jose-Ignacio Requeno, José Merseguer, and Simona Bernardi. Performance analysis of apache storm applications using stochastic petri nets. In *2017 IEEE International Conference on Information Reuse and Integration (IRI)*, pages 411–418. IEEE, 2017.
- [S20] Johannes Kroß and Helmut Krcmar. Model-based performance evaluation of batch and stream applications for big data. In *Proceedings of the IEEE 25th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems*, pages 80–86. IEEE, 2017.
- [S21] Filip Nalepa, Michal Batko, and Pavel Zezula. Model for performance analysis of distributed stream processing applications. In *Database and Expert Systems Applications*, pages 520–533. Springer, 2015.
- [S22] Riccardo Pincirolì and Catia Trubiani. Model-based performance analysis for architecting cyber-physical dynamic spaces. In *2021 IEEE 18th International Conference on Software Architecture (ICSA)*, pages 104–114. IEEE, 2021.
- [S23] Thanh Dat Nguyen, Yassine Ouhammou, and Emmanuel Grolleau. Parad repository: On the capitalization of the performance

LITERATURE DATASET

- [S1] Bixin Li, Li Liao, and Yi Cheng. Evaluating software architecture evolution using performance simulation. In *Proceedings of the 4th Intl Conf on Applied Computing and Information Technology/3rd Intl Conf on Computational Science/Intelligence and Applied Informatics/1st Intl Conf on Big Data, Cloud Computing, Data Science & Engineering*, pages 7–13. IEEE, 2016.
- [S2] Xin Du, Xin Yao, Youcong Ni, Leandro L Minku, Peng Ye, and Ruliang Xiao. An evolutionary algorithm for performance optimization at software architecture level. In *2015 IEEE Congress on Evolutionary Computation*, pages 2129–2136. IEEE, 2015.
- [S3] Davide Arcelli, Vittorio Cortellessa, Antonio Filieri, and Alberto Leva. Control theory for model-based performance-driven software adaptation. In *Proceedings of the 11th International ACM SIGSOFT Conference on Quality of Software Architectures (QoSA)*, pages 11–20. IEEE, 2015.
- [S4] Xing Wu, Yan Liu, and Ian Gorton. Exploring performance models of hadoop applications on cloud architecture. In *Proceedings of the 11th International ACM SIGSOFT Conference on Quality of Software Architectures*, pages 93–101, 2015.

- analysis process for aadl designs. In *European Conference on Software Architecture*, pages 22–39. Springer, 2017.
- [S24] Gururaj Maddodi, Slinger Jansen, and Michiel Overeem. Aggregate architecture simulation in event-sourcing applications using layered queuing networks. In *Proceedings of the ACM/SPEC International Conference on Performance Engineering*, pages 238–245, 2020.
- [S25] Rajitha Yasaweerasinghelage, Mark Staples, Ingo Weber, and Hye-Young Paik. Predicting the performance of privacy-preserving data analytics using architecture modelling and simulation. In *Proceedings of the 15th International Conference on Software Architecture (ICSA)*, pages 166–16609. IEEE, 2018.
- [S26] Qais Noorshams, Roland Reeb, Andreas Rentschler, Samuel Kounev, and Ralf Reussner. Enriching software architecture models with statistical models for performance prediction in modern storage environments. In *Proceedings of the 17th international ACM Sigsoft symposium on Component-based software engineering*, pages 45–54, 2014.
- [S27] Jürgen Walter, Christian Stier, Heiko Kozirolek, and Samuel Kounev. An expandable extraction framework for architectural performance models. In *Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering Companion*, pages 165–170. IEEE, 2017.
- [S28] Fabian Brosig, Philipp Meier, Steffen Becker, Anne Kozirolek, Heiko Kozirolek, and Samuel Kounev. Quantitative evaluation of model-driven performance analysis and simulation of component-based architectures. *IEEE transactions on software engineering*, 41(2):157–175, 2014.
- [S29] Murray Woodside, Dorina C Petriu, José Merseguer, Dorin B Petriu, and Mohammad Alhaj. Transformation challenges: from software models to performance models. *Software & Systems Modeling*, 13(4):1529–1552, 2014.
- [S30] Elham Eshraghian and Vahid Rafe. Performance measurement of models specified through component-based software architectural styles. *Measurement*, 73:372–383, 2015.
- [S31] Catia Trubiani, Achraf Ghabi, and Alexander Egyed. Exploiting traceability uncertainty between software architectural models and performance analysis results. In *Proceedings of the 2015 European Conference on Software Architecture*, pages 305–321. Springer, 2015.
- [S32] Yutong Zhao, Lu Xiao, Xiao Wang, Zhifei Chen, Bihuan Chen, and Yang Liu. Butterfly space: An architectural approach for investigating performance issues. In *Proceedings of the 17th International Conference on Software Architecture*, pages 202–213. IEEE, 2020.
- [S33] Yutong Zhao, Lu Xiao, Xiao Wang, Lei Sun, Bihuan Chen, Yang Liu, and Andre B Bondi. How are performance issues caused and resolved?—an empirical study from a design perspective. In *Proceedings of the ACM/SPEC International Conference on Performance Engineering*, pages 181–192, 2020.
- [S34] Li Wu, Johan Tordsson, Erik Elmroth, and Odej Kao. Microrca: Root cause localization of performance issues in microservices. In *NOMS 2020-2020 IEEE/IFIP Network Operations and Management Symposium*, pages 1–9. IEEE, 2020.
- [S35] Davide Arcelli, Vittorio Cortellessa, Daniele Di Pompeo, Romina Eramo, and Michele Tucci. Exploiting architecture/runtime model-driven traceability for performance improvement. In *Proceedings of the 16th IEEE International Conference on Software Architecture*, pages 81–90. IEEE, 2019.
- [S36] Anne Martens, Heiko Kozirolek, Lutz Prechelt, and Ralf Reussner. From monolithic to component-based performance evaluation of software architectures. *Empirical Software Engineering*, 16(5):587–622, 2011.
- [S37] Diego Didona, Francesco Quaglia, Paolo Romano, and Ennio Torre. Enhancing performance prediction robustness by combining analytical modeling and machine learning. In *Proceedings of the 6th ACM/SPEC international conference on performance engineering*, pages 145–156, 2015.
- [S38] Anshul Jindal, Vladimir Podolskiy, and Michael Gerndt. Performance modeling for cloud microservice applications. In *Proceedings of the 2019 ACM/SPEC International Conference on Performance Engineering*, pages 25–32, 2019.
- [S39] Floriment Klinaku, Alireza Hakamian, and Steffen Becker. Architecture-based evaluation of scaling policies for cloud applications. In *2021 IEEE International Conference on Autonomic Computing and Self-Organizing Systems (ACSOS)*, pages 151–157. IEEE, 2021.
- [S40] Fabian Gorsler, Fabian Brosig, and Samuel Kounev. Performance queries for architecture-level performance models. In *Proceedings of the 5th ACM/SPEC international conference on Performance engineering*, pages 99–110, 2014.
- [S41] Abel Gómez, Connie U Smith, Amy Spellmann, and Jordi Cabot. Enabling performance modeling for the masses: Initial experiences. In *International Conference on System Analysis and Modeling*, pages 105–126. Springer, 2018.
- [S42] Simon Eismann, Johannes Grohmann, Jürgen Walter, Jóakim Von Kistowski, and Samuel Kounev. Integrating statistical response time models in architectural performance models. In *Proceedings of the 16th IEEE International Conference on Software Architecture*, pages 71–80. IEEE, 2019.
- [S43] Catia Trubiani, Aldeida Aleti, Sarah Goodwin, Pooyan Jamshidi, Andre van Hoorn, and Samuel Gratzl. Visarch: Visualisation of performance-based architectural refactorings. In *European Conference on Software Architecture*, pages 182–190. Springer, 2020.
- [S44] Nikolaus Huber, Fabian Brosig, Simon Spinner, Samuel Kounev, and Manuel Bähr. Model-based self-aware performance and resource management using the descartes modeling language. *IEEE Transactions on Software Engineering*, 43(5):432–452, 2016.
- [S45] Manar Mazkatli, David Monschein, Johannes Grohmann, and Anne Kozirolek. Incremental calibration of architectural performance models with parametric dependencies. In *Proceedings of the 17th IEEE International Conference on Software Architecture*, pages 23–34. IEEE, 2020.
- [S46] Tse-Hsun Chen, Weiyei Shang, Zhen Ming Jiang, Ahmed E Hassan, Mohamed Nasser, and Parminder Flora. Detecting performance anti-patterns for applications developed using object-relational mapping. In *Proceedings of the 36th International Conference on Software Engineering*, pages 1001–1012, 2014.
- [S47] Davide Arcelli and Daniele Di Pompeo. Applying design patterns to remove software performance antipatterns: a preliminary approach. *Procedia Computer Science*, 109:521–528, 2017.
- [S48] Sara Fioravanti, Fulvio Patara, and Enrico Vicario. Engineering the performance of a meta-modeling architecture. In *Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering Companion*, pages 203–208, 2017.
- [S49] Martina De Sanctis, Catia Trubiani, Vittorio Cortellessa, Antiniscia Di Marco, and Mirko Flammijn. A model-driven approach to catch performance antipatterns in adl specifications. *Information and Software Technology*, 83:35–54, 2017.
- [S50] Davide Arcelli, Vittorio Cortellessa, and Daniele Di Pompeo. Performance-driven software model refactoring. *Information and Software Technology*, 95:366–397, 2018.
- [S51] Catia Trubiani, Anne Kozirolek, Vittorio Cortellessa, and Ralf Reussner. Guilt-based handling of software performance antipatterns in palladio architectural models. *Journal of Systems and Software*, 95:141–165, 2014.
- [S52] Vittorio Cortellessa, Antiniscia Di Marco, and Catia Trubiani. An approach for modeling and detecting software performance antipatterns based on first-order logics. *Software & Systems Modeling*, 13(1):391–432, 2014.
- [S53] Alexander Wert, Marius Oehler, Christoph Heger, and Roozbeh Farahbod. Automatic detection of performance anti-patterns in inter-component communications. In *Proceedings of the 10th international ACM Sigsoft Conference on Quality of Software Architectures*, pages 3–12, 2014.
- [S54] Alberto Avritzer, Ricardo Britto, Catia Trubiani, Barbara Russo, Andrea Janes, Matteo Camilli, André van Hoorn, Robert Heinrich, Martina Rapp, and Jörg Henß. A multivariate characterization and detection of software performance antipatterns. In *Proceedings of the ACM/SPEC International Conference on Performance Engineering*, pages 61–72, 2021.
- [S55] Catia Trubiani and Anne Kozirolek. Detection and solution of software performance antipatterns in palladio architectural models. In *Proceedings of the 2nd ACM/SPEC International Conference on Performance engineering*, pages 19–30, 2011.
- [S56] Davide Arcelli, Vittorio Cortellessa, and Catia Trubiani. Performance-based software model refactoring in fuzzy contexts. In *International Conference on Fundamental Approaches to Software Engineering*, pages 149–164. Springer, 2015.
- [S57] Sara Fioravanti, Simone Mattolini, Fulvio Patara, and Enrico Vicario. Experimental performance evaluation of different data models for a reflection software architecture over nosql persistence layers. In *Proceedings of the 7th International Conference on Performance Engineering*, pages 297–308. ACM, 2016.

- [S58] Gerard Haughian, Rasha Osman, and William J Knottenbelt. Benchmarking replication in cassandra and mongodb nosql datastores. In *Proceedings of the 2016 International Conference on Database and Expert Systems Applications*, pages 152–166. Springer, 2016.
- [S59] Chung-Horng Lung, Samuel A Ajila, and Wen-Lin Liu. Measuring the performance of aspect oriented software: A case study of leader/followers and half-sync/half-async architectures. *Information Systems Frontiers*, 16(5):853–866, 2014.
- [S60] Nuno Gonçalves, Diogo Faustino, António Rito Silva, and Manuel Portela. Monolith modularization towards microservices: Refactoring and performance trade-offs. In *2021 IEEE 18th International Conference on Software Architecture Companion (ICSA-C)*, pages 1–8. IEEE, 2021.
- [S61] Ghufuran Alzboon and Amro Al-Said Ahmad. A performance evaluation approach for n-tier cloud-based software services. In *Proceedings of the 2022 6th International Conference on Cloud and Big Data Computing*, pages 31–36, 2022.
- [S62] Kim Long Ngo, Joydeep Mukherjee, Zhen Ming Jiang, and Marin Litoiu. Evaluating the scalability and elasticity of function as a service platform. In *Proceedings of the 2022 ACM/SPEC on International Conference on Performance Engineering*, pages 117–124, 2022.
- [S63] Grzegorz Blinowski, Anna Ojdowska, and Adam Przybyłek. Monolithic vs. microservice architecture: A performance and scalability evaluation. *IEEE Access*, 10:20357–20374, 2022.
- [S64] Alberto Avritzer, Vincenzo Ferme, Andrea Janes, Barbara Russo, André van Hoorn, Henning Schulz, Daniel Menasché, and Vilc Rufino. Scalability assessment of microservice architecture deployment configurations: A domain-based approach leveraging operational profiles and load tests. *Journal of Systems and Software*, 165:110564, 2020.
- [S65] Wilhelm Hasselbring and André van Hoorn. Kieker: A monitoring framework for software engineering research. *Software Impacts*, 5:100019, 2020.
- [S66] Ian Gergin, Bradley Simmons, and Marin Litoiu. A decentralized autonomic architecture for performance control in the cloud. In *2014 IEEE International Conference on Cloud Engineering*, pages 574–579. IEEE, 2014.
- [S67] Chung-Horng Lung, Xu Zhang, and Pragash Rajeswaran. Improving software performance and reliability in a distributed and concurrent environment with an architecture-based self-adaptive framework. *Journal of Systems and Software*, 121:311–328, 2016.
- [S68] Davide Arcelli. A multi-objective performance optimization approach for self-adaptive architectures. In *European Conference on Software Architecture*, pages 139–147. Springer, 2020.
- [S69] Mazen Ezzeddine, Sébastien Tauvel, Françoise Baude, and Fabrice Huer. On the design of sla-aware and cost-efficient event driven microservices. In *Proceedings of the Seventh International Workshop on Container Technologies and Container Clouds*, pages 25–30, 2021.
- [S70] Yunfeng Peng, Congming Shi, Guowei Gao, Jianan Wang, and Hai Liu. Parallel component composition and performance optimization based on agent technology. In *Proceedings of the 4th International Conference on Computer Science and Software Engineering*, pages 243–252, 2021.
- [S71] Mirco Tribastone, Philip Mayer, and Martin Wirsing. Performance prediction of service-oriented systems with layered queueing networks. In *International Symposium On Leveraging Applications of Formal Methods, Verification and Validation*, pages 51–65. Springer, 2010.
- [S72] Arkadiusz Wrzosek. Applying naf for performance analysis: Performance analysis of soa systems using lqn models. In *2012 Military Communications and Information Systems Conference (MCC)*, pages 1–8. IEEE, 2012.
- [S73] Diego Perez-Palacin, José Merseguer, and Simona Bernardi. Performance aware open-world software in a 3-layer architecture. In *Proceedings of the 1st Joint WOSP/SIPEW International Conference on Performance Engineering*, pages 49–56. ACM, 2010.
- [S74] D Evangelin Geetha, TV Suresh Kumar, and K Rajani Kanth. Framework for hybrid performance prediction process model: use case performance engineering approach. *ACM SIGSOFT Software Engineering Notes*, 36(3):1–15, 2011.
- [S75] Marin Litoiu and Cornel Barna. A performance evaluation framework for web applications. *Journal of Software: Evolution and Process*, 25(8):871–890, 2013.
- [S76] Luca Berardinelli, Vittorio Cortellessa, and Antiniscia Di Marco. Performance modeling and analysis of context-aware mobile software systems. In *International Conference on Fundamental Approaches to Software Engineering*, pages 353–367. Springer, 2010.
- [S77] Catia Trubiani, Indika Meedeniya, Vittorio Cortellessa, Aldeida Aleti, and Lars Grunske. Model-based performance analysis of software architectures under uncertainty. In *Proceedings of the 9th international ACM Sigsoft conference on Quality of software architectures*, pages 69–78, 2013.
- [S78] Mohammad Alhaj and Dorina C Petriu. Traceability links in model transformations between software and performance models. In *International SDL Forum*, pages 203–221. Springer, 2013.
- [S79] Davide Arcelli and Vittorio Cortellessa. Software model refactoring based on performance analysis: better working on software or performance side? *Formal Engineering approaches to Software Components and Architectures*, pages 33–47.
- [S80] Romina Eramo, Vittorio Cortellessa, Alfonso Pierantonio, and Michele Tucci. Performance-driven architectural refactoring through bidirectional model transformations. In *Proceedings of the 8th international ACM SIGSOFT conference on Quality of Software Architectures*, pages 55–60. ACM, 2012.
- [S81] Razib Hayat Khan and Poul E Heegaard. A performance modeling framework incorporating cost efficient deployment of collaborating components. In *Proceedings of the 2nd International Conference on Software Technology and Engineering*, pages 331–340. IEEE, 2010.
- [S82] Rasha Tawhid and Dorina C Petriu. Automatic derivation of a product performance model from a software product line model. In *2011 15th international software product line conference*, pages 80–89. IEEE, 2011.
- [S83] Robert Von Massow, André Van Hoorn, and Wilhelm Hasselbring. Performance simulation of runtime reconfigurable component-based software architectures. In *European Conference on Software Architecture*, pages 43–58. Springer, 2011.
- [S84] Jens Happe, Steffen Becker, Christoph Rathfelder, Holger Friedrich, and Ralf H Reussner. Parametric performance completions for model-driven performance prediction. *Performance Evaluation*, 67(8):694–716, 2010.
- [S85] Michael Hauck, Michael Kuperberg, Nikolaus Huber, and Ralf Reussner. Ginpex: deriving performance-relevant infrastructure properties through goal-oriented experiments. In *Proceedings of the 7th International ACM SIGSOFT Conference on Quality of Software Architectures*, pages 53–62. ACM, 2011.
- [S86] Michael Hauck, Jens Happe, and Ralf H Reussner. Automatic derivation of performance prediction models for load-balancing properties based on goal-oriented measurements. In *2010 IEEE International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems*, pages 361–369. IEEE, 2010.
- [S87] Henning Groenda. Improving performance predictions by accounting for the accuracy of composed performance models. In *Proceedings of the 8th international ACM SIGSOFT conference on Quality of Software Architectures*, pages 111–116. ACM, 2012.
- [S88] Fabian Brosig, Nikolaus Huber, and Samuel Kounev. Automated extraction of architecture-level performance models of distributed component-based systems. In *Proceedings of the 26th IEEE/ACM International Conference on Automated Software Engineering*, pages 183–192. IEEE, 2011.
- [S89] Seyed Saber Jalali, Hassan Rashidi, and Eslam Nazemi. A new approach to evaluate performance of component-based software architecture. In *2011 UKSim 5th European Symposium on Computer Modeling and Simulation*, pages 451–456. IEEE, 2011.
- [S90] Alessio Gambi, Giovanni Toffetti, and Sara Comai. Model-driven web engineering performance prediction with layered queue networks. In *International Conference on Web Engineering*, pages 25–36. Springer, 2010.
- [S91] Nariman Mani, Dorina C Petriu, and Murray Woodside. Propagation of incremental changes to performance model due to soa design pattern application. In *Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering*, pages 89–100. ACM, 2013.
- [S92] Vincent Gaudel, Frank Singhoff, Alain Plantec, Stéphane Rubini, Pierre Dissaux, and Jérôme Legrand. An ada design pattern recognition tool for aadl performance analysis. *ACM SIGAda Ada Letters*, 31(3):61–68, 2011.
- [S93] Enrico Barbierato, Mauro Iacono, and Stefano Marrone. Perfbpel: A graph-based approach for the performance analysis of bpel soa

- applications. In *6th International ICST Conference on Performance Evaluation Methodologies and Tools*, pages 64–73. IEEE, 2012.
- [S94] Davide Arcelli, Vittorio Cortellessa, and Catia Trubiani. Antipattern-based model refactoring for software performance improvement. In *Proceedings of the 8th international ACM SIGSOFT conference on Quality of Software Architectures*, pages 33–42, 2012.
- [S95] Nurul Huda Nik Zulkipli and Norazlan Idris. An empirical study on performance evaluation of parallel architecture for web application services. In *2013 IEEE Symposium on Computers & Informatics (ISCI)*, pages 216–221. IEEE, 2013.
- [S96] Javier Cámara, Pedro Correia, Rogério De Lemos, David Garlan, Pedro Gomes, Bradley Schmerl, and Rafael Ventura. Evolving an adaptive industrial software system to use architecture-based self-adaptation. In *2013 8th International Symposium on Software Engineering for Adaptive and Self-Managing Systems (SEAMS)*, pages 13–22. IEEE, 2013.
- [S97] Christoph Heger, Jens Happe, and Roozbeh Farahbod. Automated root cause isolation of performance regressions during software development. In *Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering*, pages 27–38, 2013.
- [S98] Klaus Krogmann, Michael Kuperberg, and Ralf Reussner. Using genetic search for reverse engineering of parametric behavior models for performance prediction. *IEEE Transactions on Software Engineering*, 36(6):865–877, 2010.
- [S99] Yongin Kwon, Sangmin Lee, Hayoon Yi, Donghyun Kwon, Seungjun Yang, Byung-Gon Chun, Ling Huang, Petros Maniatis, Mayur Naik, and Yunheung Paek. Mantis: Automatic performance prediction for smartphone applications. In *2013 USENIX Annual Technical Conference (USENIX ATC 13)*, pages 297–308, 2013.
- [S100] Anne Kozirolek, Heiko Kozirolek, and Ralf Reussner. Peropteryx: automated application of tactics in multi-objective software architecture optimization. In *Proceedings of the 7th international ACM Sigsoft Conference on Quality of Software Architectures*, pages 33–42. ACM, 2011.