

Benchmarking as Empirical Standard in Software Engineering Research

Wilhelm Hasselbring
hasselbring@email.uni-kiel.de
Kiel University
Kiel, Germany

ABSTRACT

In empirical software engineering, benchmarks can be used for comparing different methods, techniques and tools. However, the recent ACM SIGSOFT Empirical Standards for Software Engineering Research do not include an explicit checklist for benchmarking. In this paper, we discuss benchmarks for software performance and scalability evaluation as example research areas in software engineering, relate benchmarks to some other empirical research methods, and discuss the requirements on benchmarks that may constitute the basis for a checklist of a benchmarking standard for empirical software engineering research.

CCS CONCEPTS

• **Software and its engineering;**

KEYWORDS

benchmarking, empirical software engineering, empirical standards

ACM Reference Format:

Wilhelm Hasselbring. . Benchmarking as Empirical Standard in Software Engineering Research. In . ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3463274.3463361>

1 INTRODUCTION

The recent ACM SIGSOFT Empirical Standards for Software Engineering Research [23] are meant to give reviewers guidelines to evaluate manuscripts against expectations determined by a scientific community. Empirical standards describe these community expectations. If publication venues adopt these standards, authors know the expectations in advance and can follow the essential criteria laid out by the standard. Reviewers then check submitted papers against the specific criteria in the relevant standards. These standards may also help to educate the next generations of software engineering researchers.

The ACM SIGSOFT Empirical Standards provide a catalog of such standards, starting with the General Standard that applies to all empirical research. Next, there is a set of methodology-specific standards such as experiments, questionnaire surveys and case studies. Several *supplements* for cross-cutting concerns like information visualization and sampling complement the catalog. The standards catalog is modular to reduce duplication, so most research projects will use multiple standards. Each empirical standard is essentially a one-page checklist of specific criteria that can be used by authors to conduct and report research, and by reviewers to evaluate manuscripts.

In empirical software engineering, benchmarks can be used for comparing different methods, techniques and tools. Tichy [27] summarizes the benefits for benchmarks for software research as follows:

“Software research could benefit tremendously from benchmarks. Benchmarks can be tested repeatedly and quickly without requiring human subjects. They help weed out poor techniques quickly and direct attention to the successful ones. There is a certain upfront cost for constructing the benchmark, but that effort could be shared among many researchers.”

Benchmarking should be added to the arsenal of empirical methods in order to speed up progress [26, 27]. The creation and widespread use of a benchmark within a research area is frequently accompanied by rapid technical progress and community building [24]. The Darpa Grand Challenge for self-driving cars represents an example for rapid technical progress inspired by benchmarks [27]. The existence of a benchmark is indicative of the maturity of a scientific discipline [24].

Benchmarks are used to compare different platforms, methods, tools, or techniques. They define standardized measurements to provide repeatable, objective, and comparable results [20]. In computer science, benchmarks are used to compare, for instance, the performance of database management systems [5], information retrieval algorithms [21] and cloud services [3].

For (quantitative) simulations, for instance, there exists a checklist in the Empirical Standards.¹ So far, there exists no such checklist for benchmarks in the Empirical Standards. However, benchmarking – similar to simulation – is relevant for evaluating engineering research, which is research that invents and evaluates technological artifacts [22]. This is already mentioned in the following two quotes of the Engineering Research Standard:²

- “[...] empirically compares the artifact to one or more state-of-the-art benchmarks”
- “Antipatterns: [...] evaluation consists *only* of quantitative performance data that is not compared to established benchmarks or alternative solutions”

Thus, benchmarking is already mentioned to assess new results of engineering research. However, benchmarking is not, as yet, described as an empirical standards on its own.

Section 2 discusses benchmarks for software performance and scalability evaluation as example research areas in software engineering. We relate benchmarks to some other empirical methods

¹<https://github.com/acmsigsoft/EmpiricalStandards/blob/master/docs/QuantitativeSimulation.md>

²<https://github.com/acmsigsoft/EmpiricalStandards/blob/master/docs/EngineeringResearch.md>

in Section 3. Section 4 discusses the requirements on benchmarks before Section 5 concludes the paper.

2 BENCHMARKS FOR SOFTWARE PERFORMANCE AND SCALABILITY EVALUATION

Performance benchmarks are part of the measurement-based approaches in the field of Software Performance Engineering [25]. The employment of performance benchmarks has contributed to improve generations of systems [32].

In the following subsections, we briefly present a few exemplary benchmarks for software performance and scalability evaluation to illustrate benchmarking in software engineering research. Section 2.1 presents the TeaStore benchmark that provides an example microservice-based software application together with synthetic workloads to execute the benchmarks. Section 2.2 presents the MooBench benchmark to measure the performance overhead of monitoring frameworks, with an emphasis on integrating regression benchmarking into continuous integration pipelines. The Theodolite scalability benchmark for distributed stream processing engines is presented in Section 2.3.

2.1 The TeaStore Benchmark for Performance and Scalability Benchmarking

The TeaStore is an online store for tea and tea related utilities [34]. The TeaStore benchmark application has been used as a distributed system for evaluating and extracting software performance models, for testing single and multi-tier auto-scalers, and for software energy-efficiency analysis and management.

The TeaStore software architecture consists of five distinct services and a Registry service as shown in Figure 1. All services communicate with the Registry. Additionally, the WebUI service issues calls to the Image-Provider, Authentication, Persistence and Recommender services. The TeaStore uses a client-side load balancer to allow replication of instances of one service type.

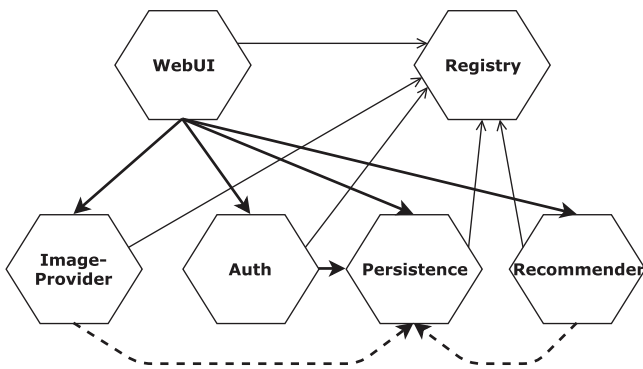


Figure 1: TeaStore Architecture [34].

As the TeaStore is a benchmarking application, it is open source and available to instrumentation using available monitoring solutions. Pre-instrumented Docker images for each service include the Kieker monitoring framework [9, 31].

Besides the TeaStore application, the benchmark provides synthetic user profiles for automated load testing. Figure 2 shows the Browse user profile. Users log in, browse the store for products, add these products to the shopping cart and then log out. The number of users is chosen depending on the maximum load. Besides such artificial user profiles, TeaStore employs synthetic workloads that are derived from two real-world traces (FIFA World Cup 1998 [1] and BibSonomy [2]). The TeaStore application employs an auto-scaler to automatically scale the store at run-time as the load intensity varies. The scaling behavior on both the FIFA and BibSonomy traces are shown in Figure 3 and in Figure 4. Both figures are structured as follows: The horizontal axis shows the experiment time in minutes; the vertical axis represents the current number of scaling units.

In Figure 3, for example, the system is in an under-provisioned state in the entire interval between minute 2 and 5. Overall, the under-provisioning and over-provisioning time-shares indicate good scaling behavior in this experiment [34].

2.2 The MooBench Monitoring Benchmark

The MooBench [37] micro-benchmark has been developed to quantify the overhead for application-level monitoring frameworks under controlled and repeatable conditions. MooBench has also been used by other researchers for replicable performance experiments comparing monitoring frameworks [18]; thus, fostering research on monitoring frameworks.

Waller & Hasselbring [37] employ the MooBench benchmark to evaluate the monitoring overhead of the Kieker [9, 31] monitoring framework and to measure the influence of different configurations for multi-core processors in this context. Fig. 6 shows the linear increase of the overhead with MooBench, when applied to Kieker.

MooBench has been integrated into the continuous integration pipeline of Kieker, allowing for automatic regression benchmarking [36]. An example visualization of a series of regression benchmark results is presented in Figure 7. It shows a performance regression that happened in March 2013 with Kieker release version 1.7, and later diagnosed as a bug in the implementation of adaptive monitoring. The daily results of Kieker with MooBench may be consulted at <http://kieker-monitoring.net/performance-benchmarks/>.

2.3 The Theodolite Scalability Benchmark for Distributed Stream Processing Engines

Scalability is usually defined as the ability of a system to continue processing an increasing workload with additional resources provided [11]. Whereas benchmarking the *performance* of stream processing engines such as throughput or latency is heavily performed by academia and industry [17], approaches for benchmarking their *scalability* are scarce. Theodolite [10] provides a method for benchmarking the scalability of distributed stream processing engines. With Theodolite, individual benchmarks are designed based on typical use cases for stream processing within microservices [7, 8]. Microservice architectures aim, in particular, at scalability [6, 19]. Theodolite supports evaluating scalability independently along different dimensions of increasing workloads. As an example, Fig. 5 compares the results for Kafka Streams and Flink via Theodolite's benchmark for aggregating time attributes.

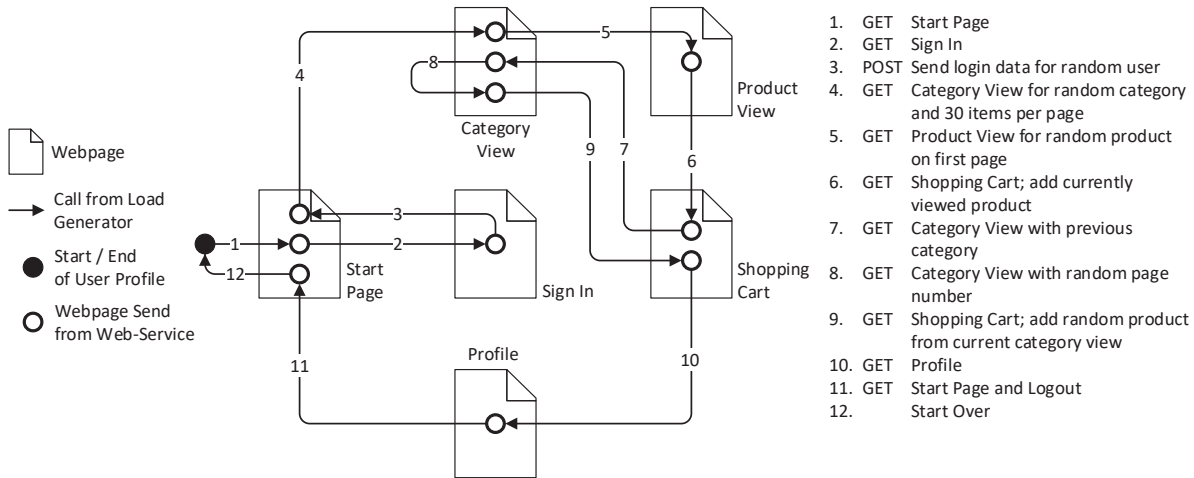


Figure 2: Browse user profile [34].

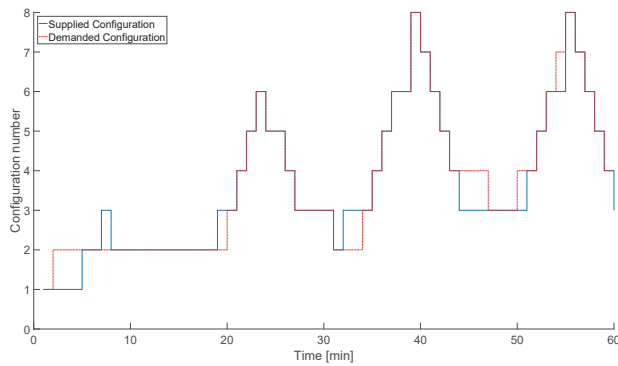


Figure 3: Scaling behavior for the FIFA trace [34].

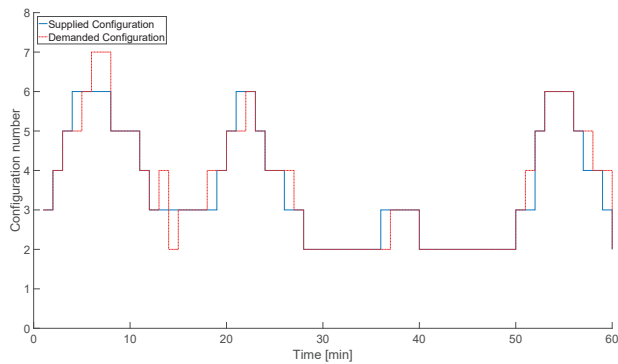


Figure 4: Scaling behavior for the BibSonomy trace [34].

A benchmark such as Theodolite does not only consist of benchmarking data, but it provides a software framework for executing the benchmarks. Figure 8 depicts the Theodolite framework architecture for executing scalability benchmarks. It consists of the following components:

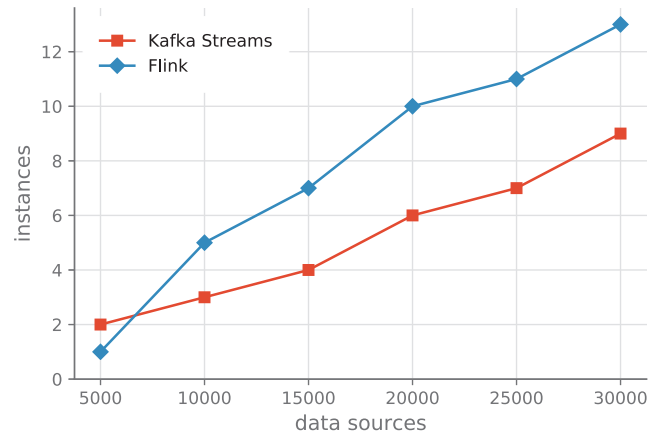


Figure 5: Comparison of Kafka Streams and Flink via Theodolite's benchmark for aggregating time attributes [10].

Experiment Control The central experiment control is started at the beginning of each scalability benchmark and runs throughout its entire execution. For each experiment, it starts and configures the workload generator component to generate the current workload of the tested dimension. Further, it starts and replicates the SUT (system under test) according to the evaluated number of instances. After each experiment, this component resets the messaging system, ensuring no queued data can be accessed by the following subexperiment.

Workload Generator This component generates a configurable constant workload of a configurable workload dimension. It fulfills the function of a data source in a big data streaming system, such as an IoT device or another microservice. Since different use cases require different data input formats, Theodolite allows for individual workload generators per

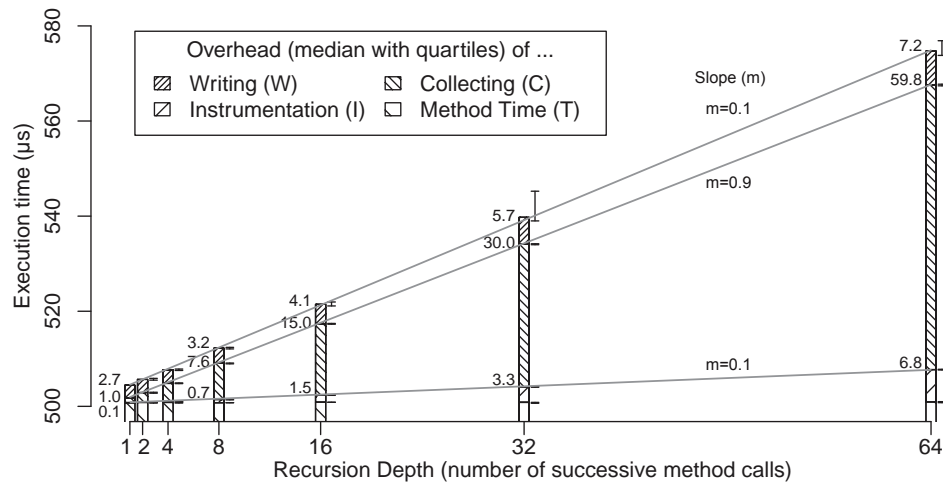


Figure 6: Benchmarking the performance overhead of monitoring frameworks with the MooBench benchmark [37].

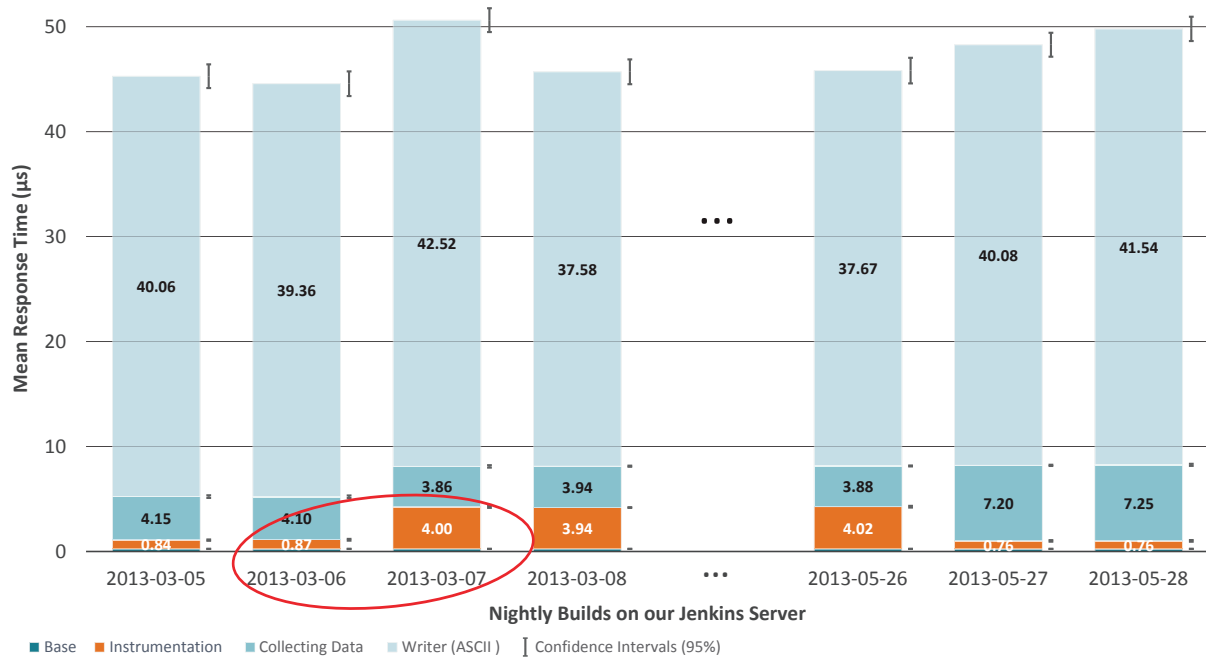


Figure 7: Scenario for detecting performance anomalies between releases via benchmarks in continuous integration [36].

use case. However, individual workload generators can share large parts of their implementations.

Messaging System In event-driven, microservice-based architectures, individual services usually communicate with each other via a dedicated messaging system. The Theodolite benchmarking architecture therefore contains such a system, serving as a message queue between workload generator and stream processing engine and as a sink for processed data. State-of-the-art messaging systems already partition the data for the stream processing engine and are, thus, likely

to have high impact on the engine's scalability. They provide plenty of configuration options, making it reasonable to benchmark different configurations against each other.

Microservice (SUT) This component acts as a microservice that applies stream processing and, thus, is the actual SUT. This microservice fulfills a specific use case. An implementation of this microservice uses a certain stream processing engine along with a certain configuration, which should be benchmarked. The stream processing engine receives all data to be processed from the messaging system and, optionally, writes processing results back to it.

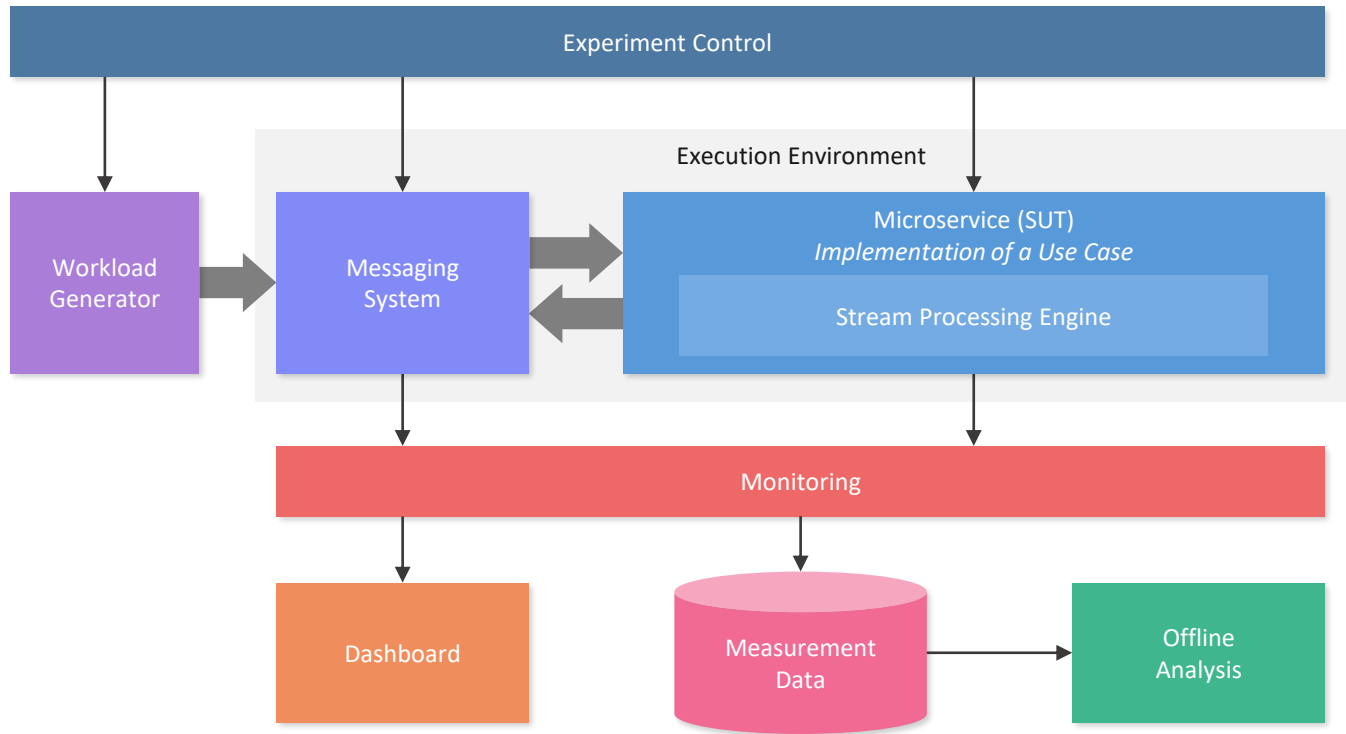


Figure 8: The Theodolite framework architecture for executing scalability benchmarks [10].

Monitoring The monitoring component collects runtime information from both the messaging system and the stream processing engine. This includes data to be displayed by the dashboard and data required to actually measure the scalability of the SUT.

Dashboard Our Theodolite architecture contains a dashboard for observation of benchmark executions. It visualizes monitored runtime information of the execution environment, the messaging system, and the SUT. Thus, it allows to verify the experimental setup (e.g., number of deployed instances and number of generated messages per seconds).

Offline Analysis Based on the raw monitoring data, a dedicated component evaluates the scalability of the SUT by computing the required metrics. This component is executed offline after completing all experiments. Since Theodolite stores monitoring data persistently, you can repeat all computations at any time without re-executing the underlying experiments.

The Theodolite benchmarking framework can be configured by the following parameters:

- (1) An implementation of the use case that should be benchmarked
- (2) Configurations for the SUT including messaging system and execution environment
- (3) The workload dimension for scalability benchmarking
- (4) A workload generator generating workloads along the configured dimensions

- (5) A list of workloads for the configured dimension to be tested

- (6) A list of numbers of instances to be tested

For details refer to Henning & Hasselbring [10].

2.4 Summary

TeaStore, MooBench and Theodolite are example benchmarks for software performance and scalability evaluation. So far, the Empirical Standards do not include review guidelines for such benchmarking experiments.

3 BENCHMARKS RELATED TO OTHER EMPIRICAL RESEARCH METHODS

Empirical research methods related to benchmarking are simulations, case studies and experiments.

3.1 Benchmarks vs. Simulations

A simulation study involves developing and using a mathematical model that imitates a real-world system's behavior. In computational science [15], such as ocean science [14, 16], we have a clear separation between the model and the real world. In computer science and software engineering this is not always the case: both, the implemented model and the real-world system are software systems.

Let us take a look at Peer-to-Peer (P2P) systems as an example domain. In the process of developing P2P systems, simulation has proved to be an essential method for the evaluation of new P2P algorithms and system architectures. The separation between a simulation model of a P2P system and a real P2P system that operates on a real physical network hinders the transition of simulation models to real systems. To bridge this gap, RealPeer [12, 13] introduces the idea of *simulation-based development* of P2P systems [4]. With RealPeer, an initial simulation model of a P2P system is iteratively transformed into the intended real P2P system. The RealPeer framework supports a developer in modeling, simulating and ultimately developing P2P systems.

Fig. 9 illustrates the use of layered models for simulation-based development of P2P systems with RealPeer. Initially, for each layer of the model a developer creates a model of the corresponding aspect (left hand side). These models are incrementally refined until they correspond to the intended real P2P system at the end of the development process (right hand side). The last step (separated by a dashed line) is a special case. In this step, each element of the model is replaced by its real counterpart. For a full introduction to RealPeer, refer to Hildebrandt & Hasselbring [13].

RealPeer supports the modeling and simulation as well as the development of P2P systems. The resulting real P2P system – not the simulation – could eventually become a candidate for a P2P benchmark.

Simulations and benchmarks have a lot in common, for instance, we can *measure* the performance of software systems via benchmarks or *assess* the performance via simulation [35]. However, the essential difference is that simulations execute a *model* of a system while benchmarks execute the *actual* system (within some experiment setup).

3.2 Benchmarks vs. Case Studies

Benchmarks require specified, synthetic workloads [29], while case study workloads should originate from real-live usage. Case studies are an empirical inquiry that investigates a contemporary phenomenon (the ‘case’) in depth and within its real-world context. Unfortunately, the label ‘case study’ is often misused in the software engineering literature [38]. For instance, illustrative examples are often called case study. But neither examples nor benchmarks are case studies according to the Empirical Standards.³

³<https://github.com/acmsigsoft/EmpiricalStandards/blob/master/docs/CaseStudy.md>

The synthetic workloads for benchmarks may be derived from real-live usage [30, 33], but for comparative and reproducible evaluation of the systems under test we need defined synthetic workloads.

An advantage of benchmarking over case studies is that replication is built into the method.

3.3 Benchmarks vs. Controlled Experiments

Experiments require a high level of control over all variables affecting the outcome but also provide reproducibility and easy comparability. Similar to experiments, benchmarks aim for a high control of the influencing variables and for reproducibility. On the other hand, the actual platform, tool, or technique evaluated by the benchmark can be highly variable, thus each benchmark run is similar to an experiment. However, the Empirical Standards only cover experiments with human participants.⁴ Benchmarking is a form of controlled experimentation [27], but not yet included in the Empirical Standards. Benchmarks require and provide particularly high levels of control.

4 REQUIREMENTS ON BENCHMARKS

Benchmarking is not only relevant for research, but has also a long history in industry. Kistowski et al. [28] introduce the primary concerns of benchmark development from the perspectives of the SPEC and TPC committees, thus industrial consortia. Benchmark candidates must undergo a process of several steps, including the definition of measurement methodologies, workload selection, and a number of rigorous benchmark acceptance tests. Kistowski et al. [28] provide a definition of the term *benchmark* in the context of performance evaluation, and define a benchmark as a *standard* tool for the competitive evaluation and comparison of competing systems or components according to specific characteristics, such as performance, dependability, or security.

Key characteristics of benchmarks are [28]:

- **Relevance** How closely the benchmark behavior correlates to behaviors that are of interest to consumers of the results.
- **Reproducibility** The ability to consistently produce similar results when the benchmark is run with the same test configuration.
- **Fairness** Allowing different test configurations to compete on their merits without artificial limitations.
- **Verifiability** Providing confidence that a benchmark result is accurate.
- **Usability** Avoiding roadblocks for users to run the benchmark in their test environments.

Similar and complementary to these characteristics, Gray [5] postulates relevance, portability, scalability and simplicity as basic benchmark requirements. Essential components of benchmarks are [24]:

- **Motivating Comparison** The purpose of a benchmark is to compare, and the motivation aspect refers to the need for the research area, and in turn the benchmark itself.

⁴<https://github.com/acmsigsoft/EmpiricalStandards/blob/master/docs/Experiments.md>

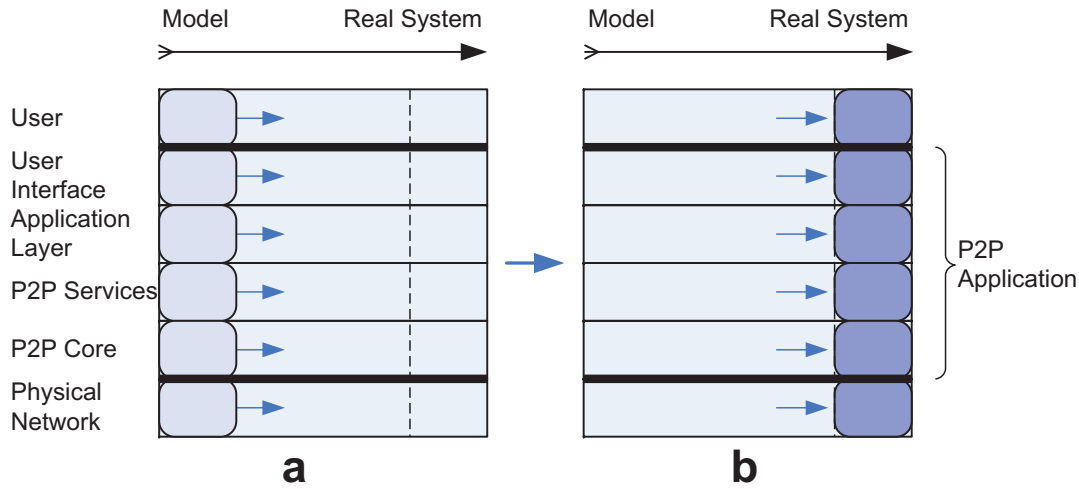


Figure 9: Simulation-based development of P2P systems with RealPeer [13].

- **Task Sample** The tests in the benchmark should be a representative sample of the tasks that the method, tool or technique is expected to solve in actual practice. For performance evaluation this is a representative workload as part of the benchmark.
- **Performance Measures** The measurements can be made by a computer or by a human, and can be quantitative or qualitative.

Tichy [26, 27] explicates that constructing benchmarks is hard work, best shared within a *community*. Furthermore, benchmarks need to evolve from narrowly targeted tests to broader, generalizable tests in order to prevent overfitting for a specific goal. Sim et al. [24] further pursue the community idea and state that benchmarks must always be developed and used in the community, rather than by a single researcher.

However, it should be sufficient to start the development process of a new benchmark with a small group of researches as an offer to a larger scientific community. Such a *proto-benchmark* [24] can act as a template to further the discussion of the topic and to initialize the consensus process. A proto-benchmark is a set of tests that is missing some of the above-mentioned requirements. The most common proto-benchmarks lack a performance measure [24]. Defining appropriate performance measures will be difficult in some areas of software engineering that involve human activities.

These requirements may constitute the basis for a *checklist* of a benchmarking standard for empirical software engineering research. However, it is necessary to continuously scrutinize and adapt benchmarks to avoid over-optimization of research towards the benchmarks.

5 SUMMARY AND OUTLOOK

This proposal paper does not report on an empirical research project; thus, it cannot be evaluated on the basis of the Empirical Standards. Instead, based on previous experience and literature on benchmarking in software engineering research we *argue* for including

benchmarks into the Empirical Standards for Software Engineering Research.

If the community (for this paper the PROPSE Workshop on Properties of Software Engineering Research) confirms my assessment that we should include benchmarking as an empirical standard, I will prepare an appropriate pull request for a benchmarking checklist at <https://github.com/acmsigsoft/EmpiricalStandards/>. This benchmarking checklist could constitute a separate document, or it could be an addition to the Engineering Research and/or Experiments standards. Some specific evaluation criteria for benchmarks, which are not yet included in these standards, are the following:

Essential Attributes

- Justifies the relevance of the benchmark.
- Describes the experimental setup for the benchmark with sufficient detail.
- Specifies the synthetic workload with sufficient detail.
- Allows different test configurations to compete on their merits without artificial limitations.
- Provides confidence that a benchmark result is accurate.
- Avoids roadblocks for users to run the benchmark in their test environments.
- Provides a replication package including datasets and analytical scripts (for Engineering Research this a desirable attribute, for benchmarks this is an essential attribute).

Desirable Attributes

- Reports on independent replication of the benchmark.
- Reports on a large community that uses the benchmark.
- Reports on an independent organization that maintains the benchmark.

ACKNOWLEDGMENTS

I would like to thank Sören Henning, Lutz Prechelt, and the anonymous reviewers for their valuable feedback on earlier versions of this paper.

REFERENCES

- [1] M. Arlitt and T. Jin. 2000. A workload characterization study of the 1998 World Cup Web site. *IEEE Network* 14, 3 (2000), 30–37. <https://doi.org/10.1109/65.844498>
- [2] Dominik Benz, Andreas Hotho, Robert Jäschke, Beate Krause, Folke Mitzlaff, Christoph Schmitz, and Gerd Stumme. 2010. The social bookmark and publication management system bibsonomy. *The VLDB Journal* 19, 6 (Dec. 2010), 849–875. <https://doi.org/10.1007/s00778-010-0208-4>
- [3] David Bermbach, Erik Wittern, and Stefan Tai. 2017. *Cloud service benchmarking*. Springer.
- [4] L. Bischofs, S. Giesecke, M. Gottschalk, W. Hasselbring, T. Warns, and S. Willer. 2006. Comparative evaluation of dependability characteristics for peer-to-peer architectural styles by simulation. *Journal of Systems and Software* 79, 10 (Oct. 2006), 1419–1432. <https://doi.org/10.1016/j.jss.2006.02.063>
- [5] Jim Gray (Ed.). 1993. *The Benchmark Handbook for Database and Transaction Systems* (2nd ed.). Morgan Kaufmann.
- [6] Wilhelm Hasselbring. 2016. Microservices for Scalability: Keynote Talk Abstract. In *Proceedings of the 7th ACM/SPEC International Conference on Performance Engineering (ICPE 2016)* (Delft, The Netherlands). ACM, New York, NY, USA, 133–134. <https://doi.org/10.1145/2851553.2858659>
- [7] Wilhelm Hasselbring. 2018. Software Architecture: Past, Present, Future. In *The Essence of Software Engineering*, Volker Gruhn and Rüdiger Striemer (Eds.). Springer International Publishing, Cham, 169–184. https://doi.org/10.1007/978-3-319-73897-0_10
- [8] Wilhelm Hasselbring and Guido Steinacker. 2017. Microservice Architectures for Scalability, Agility and Reliability in E-Commerce. In *Proc. IEEE International Conference on Software Architecture Workshops*. <https://doi.org/10.1109/ICSAW.2017.11>
- [9] Wilhelm Hasselbring and André van Hoorn. 2020. Kieker: A monitoring framework for software engineering research. *Software Impacts* 5 (Aug. 2020). <https://doi.org/10.1016/j.simp.2020.100019>
- [10] Sören Henning and Wilhelm Hasselbring. 2021. Theodolite: Scalability Benchmarking of Distributed Stream Processing Engines in Microservice Architectures. *Big Data Research* 25, 100209 (July 2021), 1–17. <https://doi.org/10.1016/j.bdr.2021.100209>
- [11] Nikolas Roman Herbst, Samuel Kounev, and Ralf Reussner. 2013. Elasticity in Cloud Computing: What It Is, and What It Is Not. In *10th International Conference on Autonomic Computing (ICAC 13)*. USENIX Association, San Jose, CA, 23–27.
- [12] Dieter Hildebrandt, Ludger Bischofs, and Wilhelm Hasselbring. 2007. RealPeer – A Framework for Simulation-based Development of Peer-to-Peer Systems. In *Proceedings of the 15th Euromicro Conference on Parallel, Distributed and Network-based Processing (PDP 2007)*. IEEE Computer Society Press, Los Alamitos, CA, USA, 490–497. <https://doi.org/10.1109/PDP.2007.70>
- [13] Dieter Hildebrandt and Wilhelm Hasselbring. 2008. Simulation-based Development of Peer-to-Peer Systems with the RealPeer Methodology and Framework. *Journal of Systems Architecture* 54, 9 (Sept. 2008), 849–860. <https://doi.org/10.1016/j.sysarc.2008.01.010>
- [14] Arne Johanson and Wilhelm Hasselbring. 2017. Effectiveness and efficiency of a domain-specific language for high-performance marine ecosystem simulation: a controlled experiment. *Empirical Software Engineering* 22, 4 (Aug. 2017), 2206–2236. <https://doi.org/10.1007/s10664-016-9483-z>
- [15] Arne Johanson and Wilhelm Hasselbring. 2018. Software Engineering for Computational Science: Past, Present, Future. *Computing in Science & Engineering* 20, 2 (March 2018), 90–109. <https://doi.org/10.1109/MCSE.2018.021651343>
- [16] Arne Johanson, Andreas Oschlies, Wilhelm Hasselbring, and Boris Worm. 2017. SPRAT: A spatially-explicit marine ecosystem model based on population balance equations. *Ecological Modelling* 349 (April 2017), 11–25. <https://doi.org/10.1016/j.ecolmodel.2017.01.020>
- [17] Jeyhun Karimov, Tilmann Rabl, Asterios Katsifodimos, Roman Samarev, Henri Heiskanen, and Volker Markl. 2018. Benchmarking Distributed Stream Data Processing Systems. In *IEEE 34th International Conference on Data Engineering (ICDE)*. 1507–1518. <https://doi.org/10.1109/ICDE.2018.00169>
- [18] Holger Knoche and Holger Eichelberger. 2018. Using the Raspberry Pi and Docker for Replicable Performance Experiments. In *Proceedings of the 2018 ACM/SPEC International Conference on Performance Engineering*. ACM. <https://doi.org/10.1145/3184407.3184431>
- [19] Holger Knoche and Wilhelm Hasselbring. 2019. Drivers and Barriers for Microservice Adoption – A Survey among Professionals in Germany. *Enterprise Modelling and Information Systems Architectures (EMISA)* – *International Journal of Conceptual Modeling* 14, 1 (2019), 1–35. <https://doi.org/10.18417/emisa.14.1>
- [20] Samuel Kounev, Klaus-Dieter Lange, and Joakim von Kistowski. 2020. *Systems Benchmarking for Scientists and Engineers*. Springer.
- [21] Federico Nanni, Bhaskar Mitra, Matt Magnusson, and Laura Dietz. 2017. Benchmark for Complex Answer Retrieval. In *Proceedings of the ACM SIGIR International Conference on Theory of Information Retrieval*. ACM. <https://doi.org/10.1145/3121050.3121099>
- [22] Paul Ralph. 2021. ACM SIGSOFT Empirical Standards Released. *ACM SIGSOFT Software Engineering Notes* 46, 1 (Feb. 2021), 19–19. <https://doi.org/10.1145/3437479.3437483>
- [23] Paul Ralph, Nauman bin Ali, Sebastian Baltes, Domenico Bianculli, Jessica Diaz, Yvonne Dittrich, Neil Ernst, Michael Felderer, Robert Feldt, Antonio Filieri, Breno Bernard Nicolau de França, Carlo Alberto Furia, Greg Gay, Nicolas Gold, Daniel Graziotin, Pinjia He, Rashina Hoda, Natalia Juristo, Barbara Kitchenham, Valentina Lenarduzzi, Jorge Martinez, Jorge Melegati, Daniel Mendez, Tim Menzies, Jefferson Moller, Dietmar Pfahl, Romain Robbes, Daniel Russo, Nyyti Saarimäki, Federica Sarro, Davide Taibi, Janet Siegmund, Diomidis Spinellis, Mirosław Staron, Klaas Stol, Margaret-Anne Storey, Davide Taibi, Damian Tamburri, Marco Torchiano, Christoph Treude, Burak Turhan, Xiaofeng Wang, and Sira Vegas. 2021. Empirical Standards for Software Engineering Research. <http://arxiv.org/abs/2010.03525> Version 0.2.0.
- [24] Susan Elliott Sim, Steve Easterbrook, and Richard C. Holt. 2003. Using benchmarking to advance research: a challenge to software engineering. In *25th International Conference on Software Engineering*. IEEE. <https://doi.org/10.1109/icse.2003.1201189>
- [25] Connie U Smith and Lloyd G Williams. 2002. *Performance solutions: a practical guide to creating responsive, scalable software*. Addison-Wesley.
- [26] Walter F. Tichy. 1998. Should computer scientists experiment more? *Computer* 31, 5 (May 1998), 32–40. <https://doi.org/10.1109/2.675631>
- [27] Walter F. Tichy. 2014. Where’s the Science in Software Engineering? Ubiquity Symposium: The Science in Computer Science. *Ubiquity* 2014 (March 2014), 1–6. <https://doi.org/10.1145/2590528.2590529>
- [28] Joakim v. Kistowski, Jeremy A. Arnold, Karl Huppler, Klaus-Dieter Lange, John L. Henning, and Paul Cao. 2015. How to Build a Benchmark. In *Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering*. ACM. <https://doi.org/10.1145/2668930.2688819>
- [29] André van Hoorn, Matthias Rohr, and Wilhelm Hasselbring. 2008. Generating Probabilistic and Intensity-Varying Workload for Web-Based Software Systems. In *Performance Evaluation: Metrics, Models and Benchmarks*. Springer, 124–143. https://doi.org/10.1007/978-3-540-69814-2_9
- [30] Andre van Hoorn, Christian Vögele, Eike Schulz, Wilhelm Hasselbring, and Helmut Krcmar. 2014. Automatic Extraction of Probabilistic Workload Specifications for Load Testing Session-Based Application Systems. In *Proceedings of the 8th International Conference on Performance Evaluation Methodologies and Tools (ValueTools 2014)* (Bratislava, Slovakia). ICST, 139–146. <https://doi.org/10.4108/icst.valuetools.2014.258171>
- [31] André van Hoorn, Jan Waller, and Wilhelm Hasselbring. 2012. Kieker: A Framework for Application Performance Monitoring and Dynamic Software Analysis. In *Proceedings of the 3rd ACM/SPEC International Conference on Performance Engineering (ICPE 2012)* (Boston, Massachusetts, USA, April 22–25, 2012). ACM, 247–248. <https://doi.org/10.1145/2188286.2188326>
- [32] Marco Vieira, Henrique Madeira, Kai Sachs, and Samuel Kounev. 2012. Resilience Benchmarking. In *Resilience Assessment and Evaluation of Computing Systems*. Springer, 283–301. https://doi.org/10.1007/978-3-642-29032-9_14
- [33] Christian Vögele, André van Hoorn, Eike Schulz, Wilhelm Hasselbring, and Helmut Krcmar. 2018. WESSBAS: extraction of probabilistic workload specifications for load testing and performance prediction—a model-driven approach for session-based application systems. *Software & Systems Modeling* 17, 2 (May 2018), 443–477. <https://doi.org/10.1007/s10270-016-0566-5>
- [34] Joakim von Kistowski, Simon Eismann, Norbert Schmitt, Andre Bauer, Johannes Grohmann, and Samuel Kounev. 2018. TeaStore: A Micro-Service Reference Application for Benchmarking, Modeling and Resource Management Research. In *2018 IEEE 26th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*. IEEE, 223–236. <https://doi.org/10.1109/mascots.2018.00030>
- [35] Robert von Massow, André van Hoorn, and Wilhelm Hasselbring. 2011. Performance Simulation of Runtime Reconfigurable Component-Based Software Architectures. In *Software Architecture (Proceedings ECSA 2011) (Lecture Notes in Computer Science, Vol. 6903)*, Ivica Crnkovic, Volker Gruhn, and Matthias Book (Eds.). Springer-Verlag, 43–58. https://doi.org/10.1007/978-3-642-23798-0_5
- [36] Jan Waller, Nils C. Ehmke, and Wilhelm Hasselbring. 2015. Including Performance Benchmarks into Continuous Integration to Enable DevOps. *SIGSOFT Softw. Eng. Notes* 40, 2 (March 2015), 1–4. <https://doi.org/10.1145/2735399.2735416>
- [37] Jan Waller and Wilhelm Hasselbring. 2012. A Comparison of the Influence of Different Multi-Core Processors on the Runtime Overhead for Application-Level Monitoring. In *Proceedings of the International Conference on Multicore Software Engineering, Performance, and Tools (MSEPT 2012) (Lecture Notes in Computer Science, Vol. 7303)*. Springer Berlin / Heidelberg, 42–53. https://doi.org/10.1007/978-3-642-31202-1_5
- [38] Claes Wohlin. 2021. Case Study Research in Software Engineering—It is a Case, and it is a Study, but is it a Case Study? *Information and Software Technology* 133 (May 2021). <https://doi.org/10.1016/j.infsof.2021.106514>