

# Continuous Benchmark Generation for Evaluating Enterprise-scale LLM Agents

Divyanshu Saxena<sup>1</sup>, Rishikesh Maurya<sup>2</sup>, Xiaoxuan Ou<sup>2</sup>, Gagan Somashekar<sup>2</sup>, Shachee Mishra Gupta<sup>2</sup>  
Arun Iyer<sup>2</sup>, Yu Kang<sup>2</sup>, Chetan Bansal<sup>2</sup>, Aditya Akella<sup>1</sup>, Saravan Rajmohan<sup>2</sup>

<sup>1</sup>The University of Texas at Austin, <sup>2</sup>Microsoft

## Abstract

The rapid adoption of AI agents across domains has made systematic evaluation crucial for ensuring their usefulness and successful production deployment. Evaluation of AI agents typically involves using a fixed set of benchmarks and computing multiple evaluation metrics for the agent. While sufficient for simple coding tasks, these benchmarks fall short for enterprise-scale agents, where services and requirements evolve continuously and ground-truth examples are sparse. We propose a process of benchmark generation that helps evolve the benchmarks as the requirements change and perform robust evaluation of evolving AI agents. We instantiate this approach for a case study of service migration from one deployment platform to another at a large public enterprise. Our approach relies on semi-structured documents where developers express the high-level intent, and uses state-of-the-art LLMs to generate benchmarks from just a small number of such documents. Overall, this process results in a maintainable evaluation framework, enabling rapid feedback on agent performance and facilitating targeted improvements.

## 1 Introduction

Modern enterprises increasingly rely on large language model (LLM)-based agents to automate software engineering and operational workflows. Examples include *build and CI/CD agents* [11] that adapt build configurations and monitor CI/CD pipelines for issues, *incident management and triaging agents* [8, 12] that route incidents to the appropriate service owners, and *troubleshooting agents* [5] that help reproduce failures and synthesize fixes. These enterprise-scale agents operate across thousands of services and heterogeneous artifacts such as source code, deployment manifests, configuration files, and documentation. Their outputs directly affect production environments—incorrect actions can impact operations, ranging from user dissatisfaction and degraded service quality to complete unavailability of critical services.

While several evaluation mechanisms and benchmarks exist for general-purpose coding agents designed for isolated and generic tasks (e.g., function synthesis or bug fixing), enterprise-scale LLM agents pose fundamentally different evaluation requirements. They must *reliably* produce *high-quality* outputs, even as the task requirements change. Moreover, these agents themselves evolve – operators continually integrate newer model versions and reasoning capabilities—making evaluation an ongoing necessity rather than a one-time exercise.

Existing benchmarks, however, are ill-suited to this setting. Most are static, targeting a fixed and narrow set of tasks. When requirements change, they must be manually revised, often resulting in numerous variants (e.g., SWE-bench [4] and its extensions [1, 6, 10]).

While useful for reproducibility and enabling transparent evaluation, they fail to capture the continuous evolution typical of enterprise environments. Their rigidity makes longitudinal evaluation difficult and imposes a recurring maintenance burden.

We identify two key challenges in evaluating enterprise-scale agents. First, these agents operate under changing dynamics, e.g., deployment platforms, operational policies, or troubleshooting guides may change frequently. An agent that performs well today may degrade as dependencies evolve or organizational requirements shift. Second, the heterogeneity of context and artifacts makes defining ground truth difficult. Conventional benchmarks typically rely on well-defined, static reference outputs, but enterprise-scale tasks span diverse file types, programming languages, and infrastructure components – all of which evolve independently, e.g. for feature additions and bug fixes. Consequently, determining what constitutes the correct or intended outcome for a given task is non-trivial, often requiring significant effort to carefully extract the right ground truths and thus, imposing a recurring maintenance burden.

We propose a new approach to evaluating enterprise-scale agents. Instead of relying on fixed, static benchmarks, we introduce a *continuous benchmark generation pipeline* that evolves alongside changing requirements without imposing significant developer burden on creating these benchmarks. In this paper, we ground the discussion for one such use case from our experience at a major cloud provider Microsoft (anonymized for review), where services were migrated to a new infrastructure platform.

**Motivating Use Case: Service Migration Agents.** Enterprise services are typically deployed on shared infrastructure platforms, which themselves evolve over time as underlying components are modified or replaced. We observed one such instance over the past six months in which existing services were transitioned to new infrastructure platform. Such platform-level changes are often intrusive, requiring the dependent services to undergo corresponding migrations. For example, when the underlying operating system changes (e.g., transitioning between Windows and Linux Operating Systems), build and compilation scripts must be updated to reflect the new environment. Similarly, if the service relies on OS-specific libraries, such as the System Drawing package (see example in Table 1) that is only functional for Windows, the source code must be updated to use cross-platform libraries (e.g., ImageSharp). As is the case for enterprise-scale agents, this service migration task spanned several heterogeneous artifacts, including build and compilation scripts, deployment files, and even source code. We discuss examples of these changes in more detail in Section 2.

**Our Proposal: Continuous Benchmark Generation.** Our vision is to develop a methodology that allows benchmarks to be evolved as operational and organizational requirements change,

File Type	Reason	Example
Script files	File path formats differ across operating systems	Change directory paths from forward slash ("/") in Linux to back slash ("\") in Windows
Source and build files	Change in supported libraries	'System.Drawing.Common' is a .NET library that provides access to Windows Graphics Device Interface. SDC does not exist on Linux systems and must be updated to cross-platform libraries like ImageSharp or SkiaSharp
Dependency configuration files	Different installation or deployment setups across OSes	Jaeger, a popular monitoring library needs to be deployed as a Windows service on Windows, but as a container on Linux/Kubernetes
Service orchestration files	Changes in requirements from the platform management teams	Modify or add new Dockerfiles for successful deployment.

**Table 1: Examples of file-level changes required during service migration.**

and enables thorough evaluation of agents over a diverse set of tasks and services. We achieve this via two key components.

The first of these are *developer-authored semi-structured Knowledge Bases (KBs)*. The term KB is generally used to denote resources used by LLM agents [3, 7, 9], but in our context, it describes the specifics of the evaluation task at hand, such as the goal, context, and expected outcomes for the task. We envision that a task (e.g., service migration) can be specified as a suite of *KB documents* – one document per sub-task (e.g., update the logging library in the service code). Crucially, any updates to the requirements can be captured simply by editing an existing KB document or adding a new one. Moreover, by developing distinct documents for heterogeneous artifacts observed in enterprise-scale services, we can get fine-grained metrics about agent performance.

The second component in our proposal are *reference implementations* derived from a small number of real commits that serve as “gold-standard” examples. By linking these examples to the corresponding KB documents, we can construct the ground truths that can be used in the benchmark. As services evolve and new migrations are completed – whether manually or automatically – new commits can be incorporated to expand and refine the benchmark.

Overall, our framework can enable longitudinal, adaptive evaluation of enterprise-scale agents as both the agents and the underlying platforms evolve. In the remainder of this paper, we (1) present our case study of platform migration at Microsoft (2) describe our vision for continuous benchmark generation (Section 3); and (3) discuss broader implications of using our approach (Section 4).

## 2 Case Study: Service Migration

The deployment environment we worked with hosts thousands of services serving over hundreds of millions of users. Performance and efficiency are critical for hosting and developing services at this scale. In an effort to achieve higher performance standards and good cluster efficiency, infrastructure teams continuously upgrade various pieces of the platform. Further, services at large enterprises may need to migrate from one platform to another, in search of better performance and efficiency trade-offs, resulting from improved hardware capabilities and up-to-date software stack support.

**Service Migration Challenges.** Modifications to the infrastructure components may involve changes transparent to the services, e.g., changing the infrastructure to add a new generation of memory devices, but they may often also involve changes that are more

intrusive in nature. For instance, one such platform migration at Microsoft required switching the underlying Operating System (e.g., between Windows OS and Linux). This, in turn, required updates to various build and source code files. We provide a representative list of such changes and examples in Table 1.

These changes are often needed across a heterogeneous set of artifacts because of several reasons. First, there are some fundamental differences between the two Operating Systems, requiring changes to how file paths are formatted, which dependencies are used, and how services are built (see Table 1). Secondly, platform management teams may also specify requirements on which libraries or dependencies they can support, and various best practices to optimally manage services. Thirdly, these changes often result in modifications to deployment files as well, because any changes in the dependencies must be reflected in the associated deployment files, such as Dockerfiles.

Table 1 highlights a representative subset of the several classes of changes needed for this task. In the specific platform migration at Microsoft, such changes were needed across a diverse set of services. Given the scale of service migration needed for our use case, we sought LLM agents for migration, aiming to automate repetitive tasks and accelerate planning. To put into perspective the extent of the complexity, a small set of 11 services that had been migrated manually by developers, took nearly 20-25 weeks per service on average and a total of over 350 PRs.

**Benchmarking Needs.** There are a wide variety of language models and tools that one can use to develop agents. Given the heterogeneous nature of the artifacts requiring migration, benchmarking is necessary to understand how well state-of-the-art agents perform, and which change patterns (such as the ones shown in 1) are best suited for which agents.

To evaluate, we construct an initial evaluation dataset based on the 11 manually migrated services: for each of those, the benchmark consists of the pre-migrated state of the service repo (that will be passed to the agent), and the oracle patch (diff between the pre- and post-migrated states of the repo).

However, just this evaluation set was insufficient for three main reasons. First, the evaluation set reflected idiosyncratic changes, specific to the services rather than systematic ones that can be generalized across services. This is because at the enterprise-level, services still undergo feature updates, bug fixes, etc. *while* platform migration is happening. For instance, a deployment file update

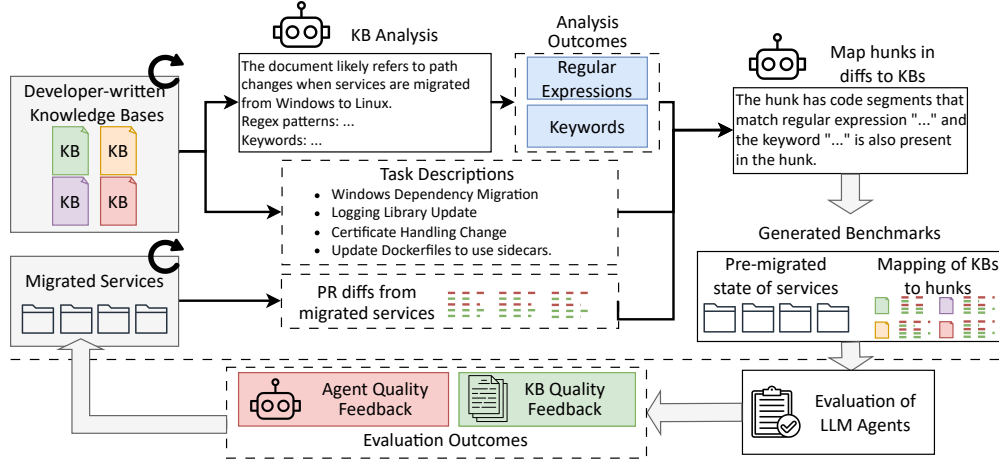


Figure 1: A pipeline to generate benchmarks from manually migrated repos.

required the addition of environment variables to specify the underlying Operating System and the respective compilation toolchain to use, however the commit from a service included additional environment variables that were later found to be service-specific. Second, benchmarks based on fixed set of services fail to capture the breadth of deployment heterogeneity because not all 11 services used all parts of the new platform. For example, some services adopted containerization with Docker, thus using Dockerfile and Kubernetes manifests for deployment, while others continued using legacy PowerShell (.ps1) and shell (.sh) scripts for deployment. This resulted in a mix of containerized and non-containerized services within the same migration wave. Finally, simply evaluating each agent and model on the entire dataset made it hard to extract the performance sliced by the change patterns – the same commit may include changes for more than one change patterns, even the same file may have changes across multiple change patterns! This motivated the design of a new benchmark generation methodology tailored for enterprise-scale, continuous migration tasks.

### 3 Continuous Benchmark Generation Pipelines

To address the limitations of existing static and manually-labor intensive benchmarks, our proposal is to develop *continuously-evolving benchmark generation pipelines* that automatically produces benchmarks, from migrated services given a set of high-level task description documents. The key principle is a clean separation between *what to test* (requirements specification) and *how to test* (i.e., produce concrete instances for evaluation).

In the remainder of this section, we describe our vision – a pipeline that uses natural language documents and migration-related commits. Figure 1 depicts the overall architecture.

**High-level Overview.** At a high level, the pipeline takes the tuple  $I = (\mathcal{S}, \mathcal{T})$  as inputs, where  $\mathcal{S}$  represents a small set of migrated services that act as the foundation for constructing the benchmark, and  $\mathcal{T}$  represents a set of task description documents prepared by developers (more details in Section 3.1). The output is the set  $\mathcal{O} = \{(\mathbf{a}_s, t, \mathcal{H}_{s,t}) \mid \forall s \in \mathcal{S}, t \in \mathcal{T}\}$ , where  $\mathbf{a}_s$  denotes the

pre-migration state of the service  $s$ , and  $\mathcal{H}_{s,t}$  denotes diff hunks from migration-related commits of service  $s$ , relevant for the task  $t$ .

#### 3.1 Knowledge Bases

At the heart of our proposal are developer-authored documents (denoted by the set  $\mathcal{T}$  in the above formulation) describing the various tasks that must be performed by the AI agent to achieve its objectives. We call these task description documents Knowledge Bases (KBs) that can be thought of as external knowledge sources that can be used by LLM agents [3].

Building these Knowledge Bases requires developers to break the high-level objective into smaller tasks, each describing a specific and related set of changes. For instance, for the high-level objective of service migration that we studied in Section 2, sub-tasks could be ‘update the logging dependencies based on the Operating System’ or ‘update the Kubernetes manifest files to utilize the latest recommended sidecar frameworks’. In our approach, a Knowledge Base is a set of documents, where each document maps to a narrow sub-task like the ones described above.

Capturing requirements in the form of KBs imposes little developer burden – these requirements are described even today in the form of migration manuals or troubleshooting guides. These KBs can then evolve naturally as platforms deprecate features, introduce new APIs, or modify resource specifications. This makes them a durable anchor point for continuous benchmark generation, as we discuss below. This intent-based system design has been a common insight in several other domains as well [2].

#### 3.2 Benchmarks from KBs and Migrated Repos

The key challenge in using migrated enterprise-level services as benchmarks is in isolating migration-related changes from idiosyncratic service updates and ensuring that outdated changes do not compromise benchmark validity. Our insight is that all the needed changes for migration, are captured in the set of task descriptions  $\mathcal{T}$ , and hence, only those changes from the migrated services must be used in the benchmark that map to some task in the set  $\mathcal{T}$ .

Repository	Precision	Recall	F1 Score
Repo1	1	0.5	0.66
Repo2	1	0.4	0.57
Repo3	1	0.667	0.8
Repo4	1	0.25	0.4

Table 2: Benchmark Generation Results

To perform this mapping, we rely on two techniques: First, if the KB document explicitly provides keywords (for instance, in a separate section in the document) that must be present in a relevant code segment, we can directly search for these keywords and map each hunk in a diff that has those keywords to this document. These keywords are used by the agents to identify "where" the changes are to be made. Secondly, KB documents may not have precise keywords, but rather descriptions of what 'lines of code' to look for. For instance, for the file path format migration (see Table 1), the KB may specify 'Windows drive names' implying that we must look for strings matching drive name patterns (such as 'C:' or 'D:'). In this case, we can leverage the reasoning and text-generation capabilities of LLMs to generate regular expressions for these types of descriptions.

**Evolving the Benchmarks.** Any updates to the requirements, can be simply incorporated into the KB documents. Further, this process of benchmark generation can also act as feedback for KB quality itself – say a new edit was made to the KB, but the benchmarks do not show the update despite having a commit for it, then one can infer that the KB content needs to be updated. Similarly, the benchmark can be expanded as more services are migrated, and more migration-related commits are available.

### 3.3 Evaluation Metrics

Our benchmarks are constructed of code diffs for each provided task. We can use this benchmark to evaluate agents by computing the following metrics:

1. *Line-edit Precision and Recall.* We can compare the agent's line edits with the ground-truth edits produced by our benchmark pipeline. Each predicted edit is matched to a ground-truth edit within a fixed edit distance. Line-edit precision measures the fraction of correctly predicted edits, while line-edit recall measures the fraction of required edits successfully performed by the agent.

2. *Per-KB Precision and Recall.* While line-edit metrics capture syntactic fidelity, they do not reveal which tasks the agent handled well or poorly. To complement them, we evaluate correctness at the KB granularity. An agent is said to cover a KB if at least one of its predicted changes matches the ground truth for that KB. Per-KB precision measures the fraction of covered KBs that were actually relevant to the service, and per-KB recall measures the fraction of required KBs that the agent successfully covered.

We believe more metrics can provide richer feedback to the pipeline, and is an interesting direction of research.

### 3.4 Evaluating the Benchmarks

We implemented the proposed pipeline and evaluated it using 137 KB documents across 4 services. For comparison, a group of developers manually constructed benchmarks based on migration-related

commits from the same services. We then assessed how well the automatically generated benchmarks aligned with these manually designed ones. As shown in Table 2, the pipeline-generated benchmarks achieved high precision and strong recall.

Interestingly, the automated benchmarks also proved to be more reliable. For instance, when evaluated against the first version of the manually-designed benchmarks, we saw low recall scores. A deeper inspection revealed that the manually written benchmarks included obsolete files that were still present in the repositories. Since our pipeline derives ground truth directly from reference commits, such files are automatically excluded, yielding cleaner and more accurate benchmarks.

## 4 Discussion and Future Directions

In this paper, we presented a case study for the evaluation of LLM agents at enterprise-scale. Using the example of platform migration induced service migration, we showcase why evaluation of enterprise-scale agents is inherently complex and continuous. Our proposal is a pipeline that can generate benchmarks, as requirements evolve or more services are available as ground-truth.

Our approach addresses two key challenges. First, it leverages long-term commit histories to build realistic benchmarks without manually filtering feature or bug-fix changes. Second, it minimizes developer effort by reusing existing natural-language migration documents to automatically construct benchmarks.

Our proposal opens several directions for future work. A key question is how to assess the accuracy of generated benchmarks, requiring new mechanisms and metrics for correctness evaluation. Another is how to use benchmarking feedback to better prompt agents – for instance, one can provide targeted feedback if the LLM agent performs poorly on specific tasks. At enterprise scale, one could even deploy multiple agents specialized for different task types. We believe this evaluation methodology may extend beyond service-migration agents studied in this paper – for instance, to troubleshooting agents leveraging natural-language guides – though this remains a direction for future exploration.

## References

- [1] 2024. SWE-bench Lite: A Lightweight Benchmark for Evaluating LLMs on Software Engineering Tasks. <https://www.swebench.com/lite.html>.
- [2] Vaastav Anand, Yichen Li, Alok Gautam Kumbhare, Celine Irvine, Chetan Bansal, Gagan Somashekar, Jonathan Mace, Pedro Las-Casas, Ricardo Bianchini, and Rodrigo Fonseca. 2025. Intent-based System Design and Operation. In *Proceedings of the 4th Workshop on Practical Adoption Challenges of ML for Systems*.
- [3] Naman Gupta, Shashank Kirtania, Priyanshu Gupta, Krishna Kariya, Sumit Gulwani, Arun Iyer, Suresh Parthasarathy, Arjun Radhakrishna, Sriram K. Rajamani, and Gustavo Soares. 2024. STACKFEED: Structured Textual Actor-Critic Knowledge Base Editing with FeedBack. [arXiv:2410.10584](https://arxiv.org/abs/2410.10584) [cs.AI]
- [4] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can Language Models Resolve Real-world Github Issues?. In *The Twelfth International Conference on Learning Representations*.
- [5] Chang Lou, Dimas Shidqi Parikesit, Yujin Huang, Zhewen Yang, Senapati Diwangkara, Yuzhuo Jing, Achmad Imam Kistijantoro, Ding Yuan, Suman Nath, and Peng Huang. 2025. Deriving Semantic Checkers from Tests to Detect Silent Failures in Production Distributed Systems. In *OSDI'25: 19th USENIX Symposium on Operating Systems Design and Implementation*.
- [6] OpenAI. 2024. Introducing SWE-bench Verified. <https://openai.com/index/introducing-swe-bench-verified/>.
- [7] Xi Wang, Taketomo Isazawa, Liana Mikaelyan, and James Hensman. 2025. KBLaM: Knowledge Base augmented Language Model. [arXiv:2410.10450](https://arxiv.org/abs/2410.10450) [cs.AI]
- [8] Zexin Wang, Jianhui Li, Minghua Ma, Ze Li, Yu Kang, Chaoyun Zhang, Chetan Bansal, Murali Chintalapati, Saravan Rajmohan, Qingwei Lin, Dongmei Zhang,

- Changhua Pei, and Gaogang Xie. 2024. Large Language Models Can Provide Accurate and Interpretable Incident Triage. In *2024 IEEE 35th International Symposium on Software Reliability Engineering (ISSRE)*.
- [9] Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin Huang, Kaidi Cao, Qian Huang, Vassilis N. Ioannidis, Karthik Subbian, James Zou, and Jure Leskovec. 2024. STaRK: Benchmarking LLM Retrieval on Textual and Relational Knowledge Bases. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [10] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2025. Demystifying LLM-Based Software Engineering Agents. *Proc. ACM Softw. Eng.* 2, FSE (June 2025).
- [11] Weiyuan Xu, Juntao Luo, Tao Huang, Kaixin Sui, Jie Geng, Qijun Ma, Isami Akasaka, Xiaoxue Shi, Jing Tang, and Peng Cai. 2025. LogSage: An LLM-Based Framework for CI/CD Failure Detection and Remediation with Industrial Validation. *arXiv:2506.03691 [cs.SE]*
- [12] Zhaoyang Yu, Aoyang Fang, Minghua Ma, Chaoyun Zhang, Ze Li, Murali Chintalapati, Xuchao Zhang, Rujia Wang, Chetan Bansal, Saravan Rajmohan, Qingwei Lin, and et.al. 2025. Triangle: Empowering Incident Triage with Multi-Agent.