

MODEL PERFORMANCE METRICS

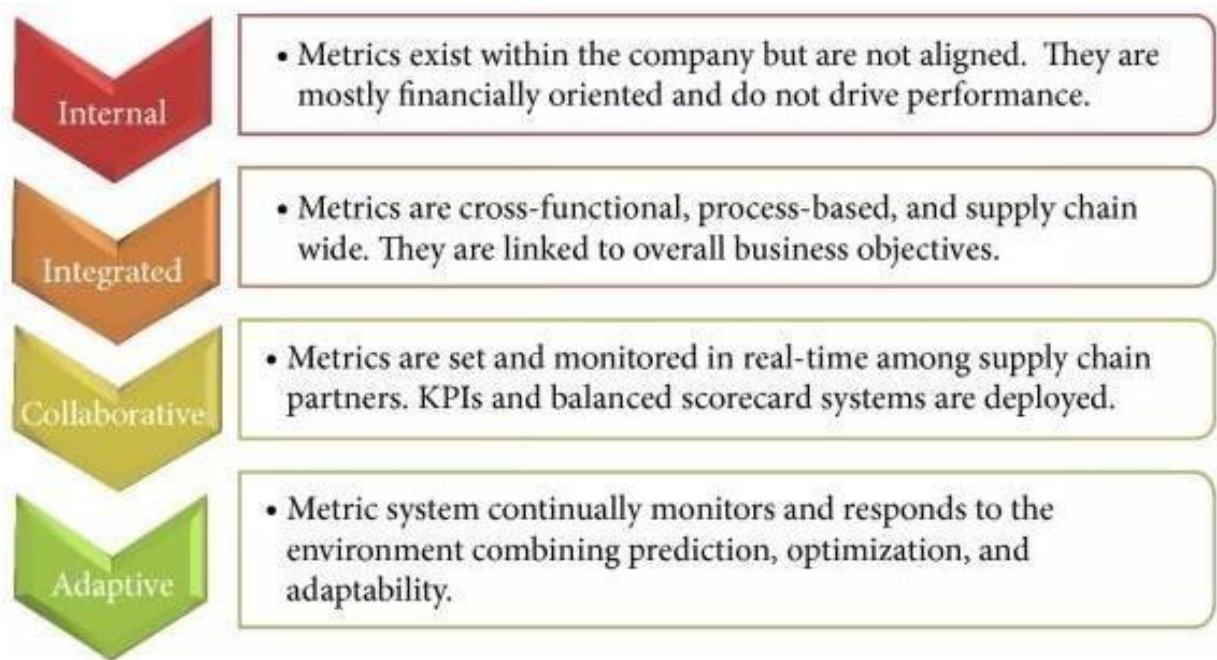
INTRODUCTION

Nowadays, the success of most companies is determined by the quality of their IT services. In addition to internal application systems, e.g., ERP systems, the application systems including the digital services provided to end- users have become business-critical assets of primarily every company. As a part of the digital transformation, customers access digital services as end- users through web browsers or dedicated smart phone apps. Examples include buying products via online retail stores and reserving concert tickets via respective booking sites. Studies have shown that the way customers experience the digital services of the companies has a direct impact on the business KPIs: satisfied customers buy products; customers who experience services that are not available or slow will move to competitors. For example, Google loses 20 % of traffic if their websites get 500 ms slower; Amazon loses 1 % of revenue for every 100 ms in latency¹. Hence, application systems and the resulting end-user experience have to be measured and optimized continuously.

COLLECTION OF PERFORMANCE MEASUREMENTS

Managing application performance requires the continuous collection of data about all relevant parts of the system starting from the end user all the way through the system. This collected data is the basis for getting a holistic end-to-end and up-to-date view of the application state including the end-user experience. In this chapter, we will discuss what data to collect, and from where and how to collect the data in order to achieve this view. Most application systems are implemented in a way that, in addition to the application logic executed at the provider's site parts of the application are executed at client's site.

The client site usually constitutes a system tier accessing the back-end



EXTRACTION OF PERFORMANCE-RELEVANT SYSTEM INFORMATION

The previous chapter focused on the collection of performance measurements from the relevant locations of the application system. This chapter focuses on the representation of higher the application system. While time series represent summary statistics (e.g., counts, percentile, etc.) over time, execution traces provide a detailed representation of the application-internal control flow that results from individual system requests.

From this data, architectural information, including logical and physical deployments and interactions (topology), can be extracted. For all cases, we will highlight examples and use cases in the context of APM level performance-relevant information about the system and their end-users that can be extracted from this data and that is used for APM visualization and reasoning, as detailed in the next chapters. Notably, we will focus on three commonly used representations, namely time series, execution traces, and augmented information about the architecture.

When depicting the number of users accessing a system, time series usually show a periodic pattern, e.g., based on the weekdays and the hours of the day. Other interesting patterns are spikes, for instance, indicating peaks in workload or hiccups.



EXECUTION TRACES

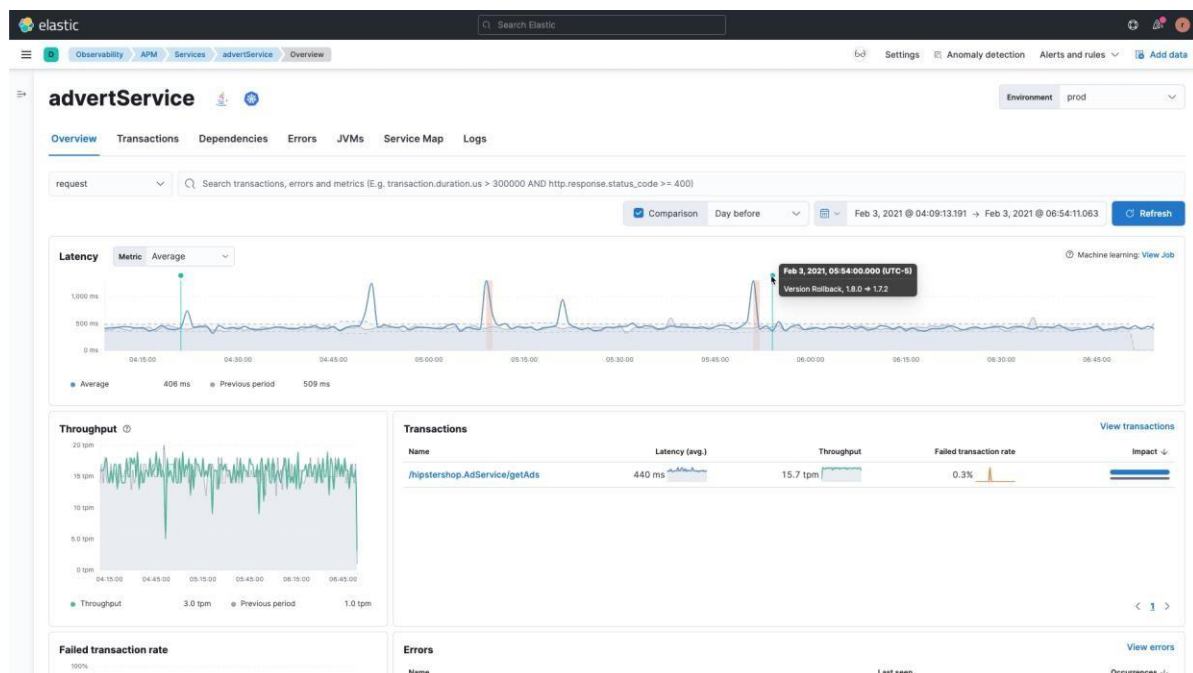
We concluded the previous section with the statement that time series are not suitable for analyzing individual requests. A data structure commonly used in APM for this purpose is an execution trace. Informally, an execution trace is a representation of the execution flow of a request through the system—ideally starting from the end user. As an example, Figure 3 depicts a schematic execution trace.

The execution trace starts with an operation called `do Filter` that is commonly found as an entry point in web-based applications. It can be observed that the execution of the `do Filter` operation includes a sequence of additional nested operation executions, until the `list` operation performs a sequence of calls to a database.

In addition to the execution flow, capturing components (e.g., Java classes or microservices) and operations, and locations (e.g., application

server, IP address), execution traces usually include further measurements. One type of performance measurement commonly found in execution traces is the response time (or duration) of each operation execution.

In the example, the response time for each operation execution is included in the second column. Moreover, execution traces may include information such as the parameters of the operation executions.

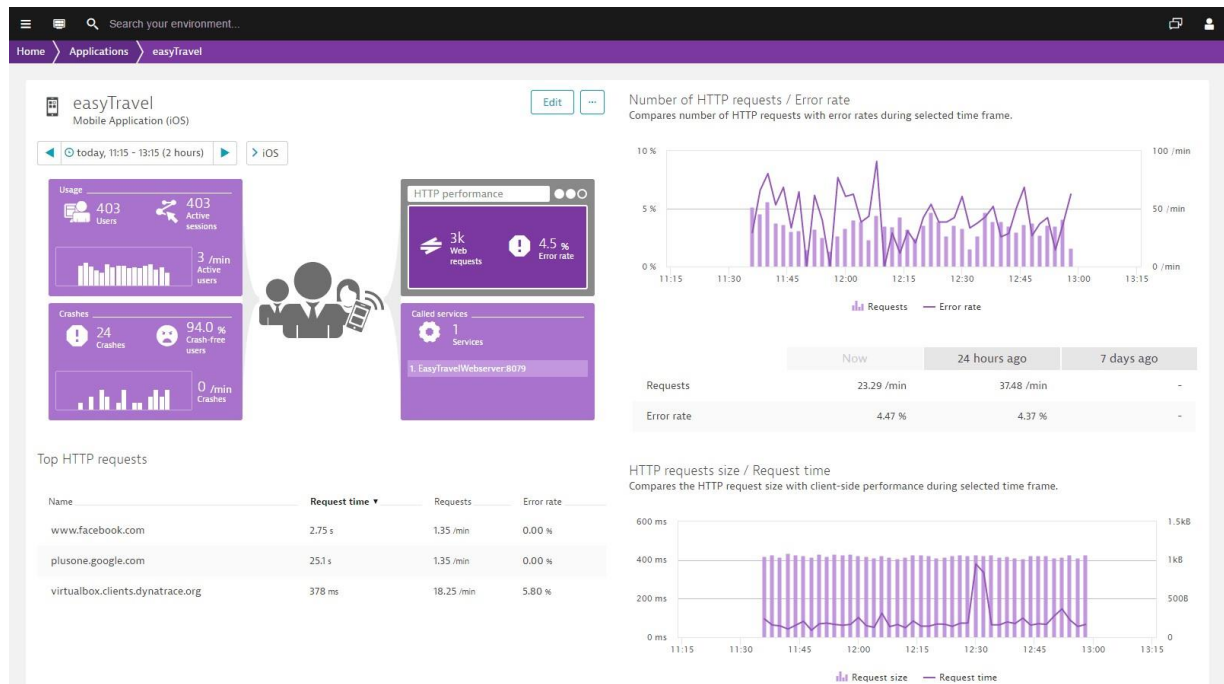


ARCHITECTURAL INFORMATION

Time series and execution traces allow to analyze the chronological order of performance measurements and of individual requests respectively. This information is commonly used to derive and represent performance-relevant architectural information of a system. The architecture of a system includes structural and dynamic information. Examples for structural information are the existence and deployment of software and hardware components.

The dynamic information includes interactions (e.g., number of calls, average response times) between components and associated information about the runtime behavior, e.g., a health state or time series. In Chapter 4 we include example of performance-augmented architectural

information. This representation is useful to have an overall state of the system and it provides a basis for a detailed manual or automated.



CONCLUSION

APM allows to provide deep insights into the run-time of application systems and its user experience, and supports the detection, diagnosis, and resolution of incidents. In the previous chapters, we gave an overview about typical APM activities and tooling support, and current challenges. We would like to emphasize the following key takeaways:

- APM is not a purely technical topic any more. Instead, it is a cross-cutting concern involving all organizational layers and units of a company.
- APM needs to be thought of from the perspective of the user experience instead of focusing purely on the back-end systems.