

CMG'11

Performance Requirements: An Attempt of a Systematic View

Paper 1107

Session 425

Alexander Podelko

alex.podelko@oracle.com

@apodelko

December 7, 2011

1

Performance requirements are supposed to be tracked from the system inception through the whole system lifecycle including design, development, testing, operations, and maintenance. However different groups of people are involved on each stage using their own vision, terminology, metrics, and tools that makes the subject confusing when going into details. The presentation is an attempt of a systematic view of the subject.

Introduction

- **The topic is more complicated than it looks**
- **Performance requirements are supposed to be tracked through the whole system lifecycle**
- **Each group of stakeholders has its own view and terminology**
- **An overview of existing issues and an attempt to create a holistic view**

Disclaimer: The views expressed here are my personal views only and do not necessarily represent those of my current or previous employers. All brands and trademarks mentioned are the property of their owners.

2

Performance requirements are supposed to be tracked from the system inception through the whole system lifecycle including design, development, testing, operations, and maintenance. However different groups of people are involved on each stage using their own vision, terminology, metrics, and tools that makes the subject confusing when going into details.

For instance, business analysts use business terms. The architects' community uses its own languages and tools (mostly created for documenting functionality so performance doesn't fit them well). Developers often think about performance through the profiler view. The virtual user notion is central for performance testers. Capacity planners use some mathematical terminology when they come up with queuing models. Production people have their own tools and metrics; and executives are more interested in high-level, aggregated metrics. These views are looking into the same subject – system performance – but through different lenses and quite often these views are not synchronized and differ noticeably. All of these views should be synchronized to allow tracing performance through all lifecycle stages and easy information exchange between stakeholders. Many existing approaches to describing performance requirements try to put these multi-dimensional and cross-dependent performance views into a set of simple flat templates designed for functional requirements.

The presentation provides an overview of current issues and is an attempt to create a holistic view of the subject.

Disclaimer: The views expressed here are my personal views only and do not necessarily represent those of my current or previous employers. All brands and trademarks mentioned are the property of their owners.

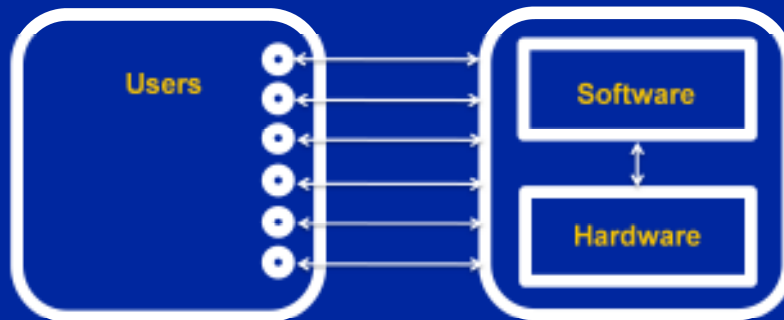
Agenda

- Metrics
- Elicitation
- Analysis and Specification
- Validation and Verification

3

First we discuss different performance metrics and terms, and then we look through the performance requirements process following IEEE Software Engineering Book of Knowledge (SWEBOK) terminology [SWEBOK04].

High-Level View of System



4

Let's take a high-level view of a system. On one side we have users who use the system to satisfy their needs. On another side we have the system, a combination of hardware and software, created (or to be created) to satisfy user's needs.

Users are not interested in what is inside the system and how it functions as soon as their requests get processed in a timely manner (leaving aside personal curiosity and subjective opinions).

Business Performance Requirements

- For today's distributed business systems
- Throughput
- Response / processing times
- All are important

5

So business requirements should state how many requests of each kind go through the system (throughput) and how quickly they need to be processed (response times). Both parts are vital: good throughput with long response times usually is as unacceptable, as are good response times with low throughput. While throughput is definitely comes from business, response times are rather usability requirements when they are good enough, but starting to impact business as soon as they become not so good.

Throughput

- **The rate at which incoming requests are completed**
 - Usually we are interested in a steady mode
- **Straightforward for homogeneous workloads**
 - Not so easy for mixed workloads: mix ratio can change with time
- **Varies with time**
 - Typical hour, peak hour, average, etc.

Throughput is the rate at which incoming requests are completed. Throughput defines the load on the system and is measured in operations per time period. It may be the number of transactions per second or the number of processed orders per hour. In most cases we are interested in a steady mode when the number of incoming requests would be equal to the number of processed requests.

Defining throughput may be pretty straightforward for a system doing the same type of business operations all the time, like processing orders or printing reports when they are homogenous. Clustering requests into a few groups may be needed if requests differ significantly: for example, small, medium, and large reports. It may be more difficult for systems with complex workloads; the ratio of different types of requests can change with the time and season.

Throughput usually varies with time. For example, throughput can be defined for a typical hour, peak hour, and non-peak hour for each particular kind of load. In environments with fixed hardware configuration the system should be able to handle peak load, but in virtualized or cloud environments it may be helpful to further detail what the load is hour-by-hour to ensure better hardware utilization.

Number of Users

- **Number of users by itself doesn't define throughput**
 - Without defining what each user is doing and how intensely
 - 500 users running one short query each minute: throughput 30,000 queries per hour
 - 500 users running one short query each hour: throughput 500 queries per hour
 - Same 500 users, 60X difference between loads

7

Quite often, however, the load on the system is characterized by the number of users. Partially it is coming from the business (in many cases the number of users is easier to find out), partially it is coming from performance tests: unfortunately, quite often performance requirements get defined during performance testing and the number of users is the main lever to manage load in load generation tools.

But the number of users doesn't, by itself, define throughput. Without defining what each user is doing and how intensely (i.e. throughput for one user), the number of users doesn't make much sense as a measure of load. For example, if 500 users are each running one short query each minute, we have throughput of 30,000 queries per hour. If the same 500 users are running the same queries, but only one query per hour, the throughput is 500 queries per hour. So there may be the same 500 users, but a 60X difference between loads (and at least the same difference in hardware requirements for the application – probably more, considering that not all systems achieve linear scalability).

Concurrency

- **Number of simultaneous users or threads**
 - Number of active users
- **Take resources even if doing nothing**
- **Number of named users**
 - Rather a data-related metric
- **Number of “really concurrent” users**
 - Number of requests in the system
 - Not an end-user performance metric

Homogenous throughput with randomly arriving requests (sometimes assumed in modeling and requirements analysis) is a simplification in most cases. In addition to different kind of requests, most systems use a kind of sessions: some system resources are associated with the user (source of requests). So the number of parallel users (sessions) would be an important requirement further qualifying throughput. In a more generic way this metric may be named concurrency: the number of simultaneous users or threads. It is important: connected but inactive users still hold some resources.

The number of online users (the number of parallel sessions) looks like the best metric for concurrency (complementing throughput and response time requirements). However terminology is somewhat vague here, sometimes “the number of users” may have a completely different meaning:

- Total or named users (all registered or potential users). This is a metric of data the system works with. It also indicates the upper potential limit of concurrency. In some cases may be used as a way to find out concurrency as a percentage of total user population, but definitely is not a concurrency metric.

- “Really concurrent” users: the number of users running requests at the same time. In most cases it is matching the number of requests in the system. While that metric looks appealing, it is not a load metric: the number of “really concurrent” requests depends on the processing time for this request. The shorter is processing time, the fewer concurrent requests we have in the system. For example, let’s assume that we got a requirement to support up to 20 “concurrent” users. If one request takes 10 sec, 20 “concurrent” requests mean throughput of 120 requests per minute. But here we get an absurd situation that if we improve processing time from 10 to one second and keep the same throughput, we miss our requirement because we have only two “concurrent” users. To support 20 “concurrent” users with a one-second response time, we really need to increase throughput 10 times to 1,200 requests per minute.

It is important to understand what users we are discussing: the difference between each of these three “number of users” metrics may be drastic.

Response Times

- **How fast requests are processed**
- **Depends on context**
 - 30 minutes may be excellent for a large batch job
- **Depends on workload**
 - Conditions should be defined
- **Aggregate metrics usually used**
 - Average, percentiles, etc.

9

Response times (in the case of interactive work) or processing times (in the case of batch jobs or scheduled activities) define how fast requests should be processed. Acceptable response times should be defined in each particular case. A time of 30 minutes could be excellent for a big batch job, but absolutely unacceptable for accessing a Web page in a customer portal. Response times depend on workload, so it is necessary to define conditions under which specific response times should be achieved; for example, a single user, average load or peak load.

Response time is the time in the system (the sum of queuing and processing time). Usually there is always some queuing time: just because the server is a complex object with sophisticated collaboration multiple components including processor, memory, disk system, and other connecting parts. That means that response time is larger than service time (to use in modeling) in most cases.

Significant research has been done to define what the response time should be for interactive systems, mainly from two points of view: what response time is necessary to achieve optimal user's performance (for tasks like entering orders) and what response time is necessary to avoid Web site abandonment (for the Internet). Most researchers agreed that for most interactive applications there is no point in making the response time faster than one to two seconds, and it is helpful to provide an indicator (like a progress bar) if it takes more than eight to 10 seconds.

Response times for each individual transaction vary, so we need to use some aggregate values when specifying performance requirements, such as averages or percentiles (for example, 90 percent of response times are less than X). Apdex standard uses a single number to measure user satisfaction [Apdex].

For batch jobs, it is important to specify all schedule-related information, including frequency (how often the job will be run), time window, dependency on other jobs and dependent jobs (and their respective time windows to see how changes in one job may impact others).

Context

- All performance metrics depend on context like:
 - Volume of Data
 - Hardware resources provided
 - Functionality included in the system

10

It is very difficult to consider performance (and, therefore, performance requirements) without full context. It depends, for example, on the volume of data involved, hardware resources provided, and functionality included in the system. So if any of that information is known, it should be specified in the requirements. Not everything may be specified at the same point: while the volume of data is usually determined by the business and should be documented at the beginning, the hardware configuration is usually determined during the design stage.

Internal (Technological) Requirements

- Important for IT
- Derived from business and usability requirements
 - During design and development
- Resources
- Scalability

11

The performance metrics of the system (the right side of the high-level view) are not important from the business (or user) point of view, but are very important for IT (people who create and operate the system). These internal (technological) requirements are derived from business and usability requirements during design and development and are very important for the later stages of the system lifecycle. Traditionally such metrics were mainly used for monitoring and capacity management because they are easier to measure and only recently tools measuring end-user performance get some traction.

Resources

- CPU, I/O, memory, and network
- Resource Utilization
 - Related to a particular configuration
 - Often generic policies like CPU below 70%
- Relative values (in percents) are not useful if configuration is not given
 - Commercial Off-the-Shelf (COTS) software
 - Virtual environments

12

The most wide-spread metric, especially in capacity management and production monitoring, is resource utilization. The main groups of resources are CPU, I/O, memory, and network. However, the available hardware resources are usually a variable in the beginning. It is one of the goals of the design process to specify hardware needed for the system from the business requirements and other inputs like company policies, available expertise, and required interfaces.

When resource requirements are measured as resource utilization, they are related to a particular hardware configuration. They are meaningful metrics when the hardware configuration is known. But these metrics doesn't make any sense as requirements until the hardware configuration would be decided upon: how we can talk, for example, about processor utilization if we don't know yet how many processors we would have? And such requirements are not useful as requirements for software if it get deployed to different hardware configurations, and, especially, for Commercial Off-the-Shelf (COTS) software.

Only way we can speak about resource utilization on early phases of the system lifecycle is as a generic policy. For example, corporate policy may be that CPU utilization should be below 70 percent.

Resources: Absolute Values

- **Absolute values**
 - # of instructions, I/O per transaction
 - Seen mainly in modeling
 - MIPS in mainframe world
- **Importance increases again with the trends of virtualization, cloud computing, and service-oriented architectures**
 - VMware: CPU usage in MHz
 - Microsoft: Megacycles

13

When required resources are specified in absolute values, like the number of instructions to execute or the number of I/O operations per transaction (as sometimes used, for example, for modeling), it may be considered as a performance metric of the software itself, without binding it to a particular hardware configuration. In the mainframe world, MIPS was often used as such metric for CPU consumption, but there is no such widely used metric in the distributed systems world.

The importance of resource-related requirements is increasing again with the trends of virtualization, cloud computing, and service-oriented architectures. When we depart from the “server(s) per application” model, it becomes difficult to specify requirements as resource utilization, as each application will add only incrementally to resource utilization. There are attempts to introduce such metrics. For example, the ‘CPU usage in MHz’ or ‘usagemhz’ metric used in the VMware world or the ‘Megacycles’ metric sometimes used by Microsoft [Microsoft10]. Another related metric sometimes (but rarely) used is efficiency when it is defined as throughput divided by resources (however the term is often used differently).

In the ideal case (for example, when the system is CPU bound and we can scale the system linearly just adding processors) we can easily find hardware configuration needed if we have an absolute metric of resources required.

For example, if software needs X units of hardware power per request and a processor has Y units of hardware power, we can calculate the number of such processors N needed for processing Z requests as $N=Z*X/Y$. The reality, of course, is more sophisticated. First of all, we have different kinds of hardware resources: processors, memory, I/O, and network. Usually we concentrate on the most critical one keeping in mind others as restrictions.

Scalability

- Ability of the system to meet performance requirements as the demand increases
- Increasing # of users, transaction volumes, data sizes, new workloads, etc.
- Performance requirements as a function, for example, of load or data and configuration
 - No free / ideal scalability

14

Scalability is system's ability to meet the performance requirements as the demand increases (usually by adding hardware). Scalability requirements may include demand projections such as an increasing of the number of users, transaction volumes, data sizes, or adding new workloads. How response times will increase with increasing load or data is important too (load or data sensitivity).

From a performance requirements perspective, scalability means that you should specify performance requirements not only for one configuration point, but as a function of load or data. For example, the requirement may be to support throughput increase from five to 10 transactions per second over next two years with response time degradation not more than 10 percent.

Scalability is also a technological (internal IT) requirement. Or perhaps even a “best practice” of systems design. From the business point of view, it is not important how the system is maintained to support growing demand. If we have growth projections, probably we need to keep the future load in mind during the system design and have a plan for adding hardware as needed.

Agenda

- Metrics
- Elicitation
- Analysis and Specification
- Validation and Verification

IEEE SWEBOK

- IEEE Software Engineering Book of Knowledge defines four stages for requirements:
 - Elicitation
 - Where come from and how to collect them
 - Analysis
 - Classify / Elaborate / Negotiate
 - Specification
 - Production of a document
 - Validation

16

IEEE Software Engineering Book of Knowledge defines four stages for requirements [SWEBOK04]:

- Elicitation: Where requirements come from and how to collect them.
- Analysis: Classify / Elaborate / Negotiate.
- Specification: Production of a document. While documenting requirements is important, the way to do this depends on software development methodology used, corporate standards, and other factors.
- Validation: making sure that requirements are correct.

Where do performance requirements come from?

- Business
- Usability
- Technology

17

If we look at the performance requirements from another point of view, we can classify them into business, usability, and technological requirements.

Business Requirements

- **Comes from the business, may be caught before design starts**
 - **Number of orders per hour**
- **The main trap is to immediately link them to a specific design and technology thus limiting the number of available choices**
 - **For example, it may be one page per order or a sequence of two dozen screens**
 - **Each of the two dozen may be saved separately or all at the end**

18

Business requirements come directly from the business and may be captured very early in the project lifecycle, before design starts. For example, "a customer representative should enter 20 requests per hour and the system should support up to 1000 customer representatives". If translated into more technical terms, requests should be processed in five minutes on average, throughput would be up to 20,000 requests per hour, and there could be up to 1,000 parallel sessions.

The main trap here is to immediately link business requirements to a specific design, technology, or usability requirements, thus limiting the number of available design choices. If we consider a Web system, for example, it is probably possible to squeeze all the information into a single page or have a sequence of two dozen screens. All information can be saved at once or each page of these two-dozen can be saved separately. We have the same business requirements, but response times per page and the number of pages per hour would be different.

Requirements Elicitation

- ***Final* requirements should be quantitative and measurable**
- **Business people know what the system should do and may provide some information**
 - They are not performance experts
- **Document real business requirements in the form they are available**
 - Then elaborate them into quantitative and measurable

19

While the final requirements should be quantitative and measurable, it is not an absolute requirement for initial requirements. Scott Barber, for example, advocates that we need to gather qualitative requirements first [Barber07]. While business people know what the system should do and may provide some numeric information, they are not trained in requirement elaboration and system design. If asked to provide quantitative and measurable requirements, they finally provide them based on whatever assumptions they have about system's design and human-computer interaction, but quite often this results in wrong assumptions being documented as business requirements. We need to document real business requirements in the form they are available, and only then elaborate them into quantitative and measurable.

Goals vs. Requirements

- **Most response times "requirements" are goals**
 - Missing them won't prevent deploying the system
- **For response times, the difference between goals and requirements may be large**
 - For many web applications goals are two-five seconds and requirements somewhere between eight seconds and one minute

20

One often missed issue, as Scott Barber notes, is goals versus requirements [Barber07]. Most of response time “requirements” (and sometimes other kinds of performance requirements,) are goals, not requirements: something that we want to achieve, but missing them won’t necessarily prevent deploying the system.

In many cases, especially for response times, there is a big difference between goals and requirements (the point when stakeholders agree that the system can’t go into production with such performance). For many interactive web applications, response time goals are two to five seconds and requirements may be somewhere between eight seconds and one minute.

One approach may be to define both goals and requirements. The problem is that, except when coming from legal or contractual obligation, requirements are very difficult to get. Even if stakeholders define performance requirements, quite often, when it comes to go/no go decisions, it becomes clear that it was not the real requirements, but rather second-tier goals.

See The Whole Picture

- For example, the requirement is 10 seconds
- We got 15 seconds for peak load
- But what if
 - Only on busiest day of the year
 - All other days it will be below 10 seconds
 - It is CPU-constrained and may be fixed by additional hardware

21

In addition, multiple performance metrics only together provide the full picture. For example, you may have a 10-second requirement and you get 15-second response time under the full load. But what if you know that this full load is the high load on the busiest day of year, that response times for the maximal load for other days are below 10 second, and you see that it is CPU-constrained and may be fixed by a hardware upgrade? Real response time requirements are so environment- and business-dependent that for many applications it may be problematic to force people to make hard decisions in advance for each possible combination of circumstances. One approach may be to specify goals (making sure that they make sense) and only then, if they are not met, make the decision what to do with all the information available.

Determining Specific Requirements

- It depends
- Approach the subject from different points of view
- Just to illustrate here are 10 methods suggested by Peter Sevcik to find T in APDEX
 - T is threshold between satisfied and tolerating users; should be strongly correlated with the response time goal

22

Determining what specific performance requirements is another large topic that is difficult to formalize. Consider the approach suggested by Peter Sevcik for finding T, the threshold between satisfied and tolerating users. T is the main parameter of the Apdex (Application Performance Index) methodology, providing a single metric of user satisfaction with the performance of enterprise applications. Peter Sevcik defined ten different methods [Sevcik08].

Methods 1-5 to Find T (by Peter Sevcik)

- Default value (4 sec)
- Empirical data
- User behavior model (# of elements/task repetitiveness)
- Outside references
- Observing users

23

- Default value (the Apdex methodology suggest 4 sec)
- Empirical data
- User behavior model (number of elements viewed / task repetitiveness)
- Outside references
- Observing the user

Methods 6-10 to Find T (by Peter Sevcik)

- Controlled performance experiment
- Best time multiple
- Find frustration threshold F first and calculate T from F ($F=4T$ in APDEX)
- Interview stakeholders
- Mathematical inflection point

24

- Controlled performance experiment
- Best time multiple
- Find frustration threshold F first and calculate T from F (the Apdex methodology assumes that $F = 4T$)
- Interview stakeholders
- Mathematical inflection point.

Each method is discussed in details in [Sevcik08].

Suggested Approach

- So Peter Sevcik suggests to use several of these methods: if all come approximately to the same number it will be T
- A similar approach can be used for performance requirements: use several methods to get the numbers – you get goal/requirement if they are close
 - Investigate / sort out if they differ significantly

25

The idea is use of several (say, three) of these methods for the same system. If all come to approximately the same number, they give us T. While the approach was developed for production monitoring, there is definitely a strong correlation between T and the response time goal (having all users satisfied sounds a pretty good goal), and between F and the response time requirement. So the approach probably can be used for getting response time requirements with minimal modifications. While some specific assumptions like four seconds for default or the $F=4T$ relationship may be up for argument, the approach itself conveys the important message that there are many ways to determine a specific performance requirement and it would be better to get it from several sources for validation purposes. Depending on your system, you can determine which methods from the above list (or maybe some others) are applicable, calculate the metrics and determine your requirements.

Usability Requirements

- Many researchers agree that response times should be 2-5 seconds for maximum performance and users lose focus if response times are more than 8-10 seconds
- Sometimes linked closely to business requirements
 - Make sure that response times are not worse than competitor's

26

Usability requirements, mainly related to response times, are based on the basic principles of human-computer interaction. Many researchers agree that users lose focus if response times are more than eight to 10 seconds and that response times should be two to five seconds for maximum productivity. These usability considerations may influence design choices (such as using several Web pages instead of one). In some cases, usability requirements are linked closely to business requirements; for example, make sure that your system's response times are not worse than response times of similar or competitor systems.

Response Times: Review of Research

- In 1968 Robert Miller defined three threshold levels of human attention
- Instantaneous 0.1-0.2 seconds
- Free interaction 1-5 seconds
- Focus on dialog 5-10 seconds

27

As long ago as 1968, Robert Miller's paper "Response Time in Man-Computer Conversational Transactions" described three threshold levels of human attention [Miller68]. Jakob Nielsen believes that Miller's guidelines are fundamental for human-computer interaction, so they are still valid and not likely to change with whatever technology comes next [Nielsen94]. These three thresholds are:

- Users view response time as instantaneous (0.1-0.2 second)
- Users feel they are interacting freely with the information (1-5 seconds)
- Users are focused on the dialog (5-10 seconds)

Instantaneous Response Time

- Users feel that they directly manipulate User Interface (UI)
- For example, between typing a symbol and its appearance on the screen
- 0.1-0.2 seconds
- Often beyond the reach of application developers
 - System/UI libraries, client-side

28

Users view response time as instantaneous (0.1-0.2 second): Users feel that they directly manipulate objects in the user interface. For example, the time from the moment the user selects a column in a table until that column highlights or the time between typing a symbol and its appearance on the screen. Robert Miller reported that threshold as 0.1 seconds. According to Peter Bickford 0.2 second forms the mental boundary between events that seem to happen together and those that appear as echoes of each other [Bickford97].

Although it is a quite important threshold, it is often beyond the reach of application developers. That kind of interaction is provided by operating system, browser, or interface libraries, and usually happens on the client side, without interaction with servers (except for dumb terminals, that is rather an exception for business systems today).

Free Interaction

- Notice delay, but "feel" the computer is "working"
- Earlier researchers reported 1-2 sec
 - Simple terminal interface
- For problem solving tasks no performance degradation up to 5 sec
 - Depends on the number of elements and repetitiveness of the task

29

Users feel they are interacting freely with the information (1-5 seconds): They notice the delay, but feel the computer is "working" on the command. The user's flow of thought stays uninterrupted. Robert Miller reported this threshold as one-two seconds [Miller68].

Peter Sevcik identified two key factors impacting this threshold [Sevcik03]: the number of elements viewed and the repetitiveness of the task. The number of elements viewed is, for example, the number of items, fields, or paragraphs the user looks at. The amount of time the user is willing to wait appears to be a function of the perceived complexity of the request. The complexity of the user interface and the number of elements on the screen both impact thresholds. Back in 1960s through 1980s the terminal interface was rather simple and a typical task was data entry, often one element at a time. So most earlier researchers reported that one to two seconds was the threshold to keep maximal productivity. Modern complex user interfaces with many elements may have higher response times without adversely impacting user productivity. Users also interact with applications at a certain pace depending on how repetitive each task is. Some are highly repetitive; others require the user to think and make choices before proceeding to the next screen. The more repetitive the task, the better expected response time.

That is the threshold that gives us response time usability goals for most user-interactive applications. Response times above this threshold degrade productivity. Exact numbers depend on many difficult-to-formalize factors, such as the number and types of elements viewed or repetitiveness of the task, but a goal of three to five seconds is reasonable for most typical business applications.

Does It Change with Time?

- Do expectations increase with time?
 - 2009 Forrester research suggests 2 second response time, in 2006 similar research suggested 4 seconds
 - The approach is often questioned: they just ask. It is known that user perception of time may be misleading
 - What page are we talking about?

30

There are researches suggests that response time expectations increase with time. Forrester research [Forrester09] of 2009 suggests 2 second response time, in 2006 similar research suggested 4 seconds (both researches were sponsored by Akamai, a provider of Web accelerating solutions). While the trend probably exists, the approach of this research was often questioned: they just ask. It is known that user perception of time may be misleading. Also, as mentioned earlier, response time expectations depends on the number of elements viewed, the repetitiveness of the task, user assumptions of what the system is doing, and UI showing the status. Stating standard without specification of what page we are talking about may be overgeneralization.

Focus on Dialog

- **Users are focused on the task: 5-10 sec**
- **Half of users abandon Web pages after 8.5 sec - Peter Bickford, 1997**
 - 2 min delay after 27 quick interactions
 - Watch cursor kept users 20 sec, animated cursor 1 min, progress bar until the end
- **Users should reorient themselves after a delay above the threshold**

31

Users are focused on the dialog (5-10 seconds): They keep their attention on the task. Robert Miller reported that threshold as 10 seconds [Miller68]. Users will probably need to reorient themselves when they return to the task after a delay above this threshold, so productivity suffers.

Peter Bickford investigated user reactions when, after 27 almost instantaneous responses, there was a two-minute wait loop for the 28th time for the same operation. It took only 8.5 seconds for half the subjects to either walk out or hit the reboot [Bickford97]. Switching to a watch cursor during the wait delayed the subject's departure for about 20 seconds. An animated watch cursor was good for more than a minute, and a progress bar kept users waiting until the end. Bickford's results were widely used for setting response times requirements for Web applications.

That is the threshold that gives us response time usability requirements for most user-interactive applications. Response times above this threshold cause users to lose focus and lead to frustration. Exact numbers vary significantly depending on the interface used, but it looks like response times should not be more than eight to 10 seconds in most cases. Still, the threshold shouldn't be applied blindly; in many cases, significantly higher response times may be acceptable when appropriate user interface is implemented to alleviate the problem.

Agenda

- Metrics
- Elicitation
- Analysis and Specification
- Validation and Verification

Technological Requirements

- Comes from the chosen design and used technology
 - We call ten web services sequentially to show a page within 3 sec. It translates into requirements of 200-250 ms for each web service
 - resource utilization requirements

33

The third category, technological requirements, comes from chosen design and used technology. Some of technological requirements may be known from the beginning if some design elements are given, but others are derived from business and usability requirements throughout the design process and depend on the chosen design.

For example, if we need to call ten Web services sequentially to show the Web page with a three-second response time, the sum of response times of each Web service, the time to create the Web page, transfer it through the network and render it in a browser should be below 3 second. That may be translated into response time requirements of 200-250 milliseconds for each Web service. The more we know, the more accurately we can apportion overall response time to Web services.

Another example of technological requirements is resource consumption requirements. For example, CPU and memory utilization should be below 70% for the chosen hardware configuration.

Analysis and Modeling

- Final requirements are elaborated from business requirements by applying usability and technological requirements
- Requirements traceability
 - Where it came from
- Input for Software Performance Engineering
 - For example, defining service / stored procedure response times by its share in the end-to-end performance budget

34

Business requirements should be elaborated during design and development, and merge together with usability and technological requirements into the final performance requirements, which can be verified during testing and monitored in production. The main reason why we separate these categories is to understand where the requirement comes from: is it a fundamental business requirement and the system fails if we miss it or a result of a design decision that may be changed if necessary.

Performance requirements are important input for Software Performance Engineering [Smith02]. During design and development the requirements are further elaborated. For example, the service / stored procedure response time requirements should be determined by its share in the end-to-end performance budget. In this way, the worst possible combination of all required services, middleware and presentation layer overheads will provide the required time. For example, if there is a Web page with 10 drop-down boxes calling 10 separate services, the response time objective for each service may be 0.2 seconds to get three seconds average response time (leaving one second for network, presentation, and rendering).

Documenting Requirements

- Requirements / Architect's vocabulary
- Quality Attributes
 - Part of Nonfunctional Requirements
- Approaches
 - Text
 - Quality Attribute Scenarios (SEI)
 - Planguage

35

Requirement Engineering / Architect's vocabulary is very different from what is used in performance testing / capacity planning. Performance and scalability are examples of Quality Attributes (QA), part of Nonfunctional Requirements (NFR).

In addition to specifying requirements in plain text, there are multiple approaches to formalize documenting of requirements. For example, Quality Attribute Scenarios by The Carnegie Mellon Software Engineering Institute (SEI) or Planguage (Planning Language) introduced by Tom Gilb.

Quality Attribute Scenarios

- QA scenario defines:
 - Source
 - Stimulus
 - Environment
 - Artifact
 - Response
 - Response Measure

36

QA scenario defines source, stimulus, environment, artifact, response, and response measure [Bass03]. For example, the scenario may be “Users initiate 1,000 transactions per minute stochastically under normal operations, and these transactions are processed with an average latency of two seconds.” For this example:

- Source is a collection of users.
- Stimulus is the stochastic initiation of 1,000 transactions per minute.
- Artifact is always the system's services.
- Environment is the system state, normal mode in our example.
- Response is processing the transactions.
- Response measure is the time it takes to process the arriving events (an average latency of two seconds in our example).

Planguage

- **Tag: unique identifier**
- **Gist: brief description**
- **Scale: unit of measure**
- **Meter: how to measure**
- **Minimum / Plan / Stretch/ Wish : levels to attain**
- **Past / Record / Trend**

37

Planguage (Planning language) was suggested by Tom Gilb and may work better for quantifying quality requirements [Simmons01]. Planguage keywords include:

- Tag: a unique identifier
- Gist: a short description
- Stakeholder: a party materially affected by the requirement
- Scale: the scale of measure used to quantify the statement
- Meter: the process or device used to establish location on a Scale
- Must: the minimum level required to avoid failure
- Plan: the level at which good success can be claimed
- Stretch: a stretch goal if everything goes perfectly
- Wish: a desirable level of achievement that may not be attainable through available means
- Past: an expression of previous results for comparison
- Trend: an historical range or extrapolation of data
- Record: the best-known achievement

It is very interesting that Planguage defines four levels for each requirement: minimum, plan, stretch, and wish.

What Metrics to Use?

- Average
- Max
- Percentiles (X% below Y sec)
- Median
- Typical
- etc.

38

Another question is how to specify response time requirements or goals. Individual transaction response times vary, so aggregate values should be used. For example, such metrics as average, maximum, different kinds of percentiles, or median. The problem is that whatever aggregate value you use, you lose some information.

The Issue

- **SLA (Service Level Agreement)**
 - "99.5% of all transactions should have a response time less than five seconds"
- **What happens with the rest 0.5%?**
 - All 6-7 seconds
 - All failed/timeout
- **Add different types of transactions, different input data, different user locations, etc.**

39

Percentiles are more typical in SLAs (Service Level Agreements). For example, "99.5 percent of all transactions should have a response time less than five seconds". While that may be sufficient for most systems, it doesn't answer all questions. What happens with the remaining 0.5 percent? Do these 0.5 percent of transactions finish in six to seven seconds or do all of them timeout? You may need to specify a combination of requirements: for example, 80 percent below four seconds, 99.5 percent below six seconds, and 99.9 percent below 15 seconds (especially if we know that the difference in performance is defined by distribution of underlying data). Other examples may be average four seconds and maximal 12 seconds, or average four seconds and 99 percent below 10 seconds.

Observability

- **Four different viewpoints**
 - Management
 - Engineering
 - QA Testing
 - Operations
- **Ideal would be different views of the same performance database**
- **Reality is a mess of disjoint tools**

40

As Adrian Cockcroft noted about observability [Cockcroft00], in addition to losing information when you aggregate, there are different viewpoints for performance data that need to be provided for different audiences. You need different metrics for management, engineering, operations, and quality assurance. For operations and management percentiles may work best. If you do performance tuning and want to compare two different runs, average may be a better metric to see the trend. For design and development you may need to provide more detailed metrics; for example, if the order processing time depends on the number of items in the order, it may be separate response time metrics for one to two, three to 10, 10 to 50, and more than 50 items.

Moreover, the tools providing performance information for different audiences are usually different, present information in a different way, and may measure different things. For example, load testing tools and active monitoring tools provide metrics for the used synthetic workload that may differ significantly from the actual production load. This becomes a real issue if you want to implement some kind of process, such as Six Sigma, to keep performance under control throughout the whole system lifecycle.

Metrics to Use

- **Combination of percentile and availability metric works in many cases**
 - 97% below 5 sec, less than 1% failed/timeout
- **An example of another approach:**
 - **Apdex (Application Performance Index)**
 - **Objective user satisfaction metric**
 - **A number between 0 and 1**
 - **0 no users satisfied, 1 all users satisfied**

Things get more complicated when there are many different types of transactions, but a combination of percentile-based performance and availability metrics usually works in production for most interactive systems. While more sophisticated metrics may be necessary for some systems, in most cases they make the process overcomplicated and results difficult to analyze.

There are efforts to make an objective user satisfaction metric. For example, Apdex - Application Performance Index [Apdex]. Apdex is a single metric of user satisfaction with the performance of enterprise applications. The Apdex metric is a number between 0 and 1, where 0 means that no users were satisfied, and 1 means all users were satisfied. The approach introduces three groups of users: satisfied, tolerating, and frustrated. Two major parameters are introduced: threshold response times between satisfied and tolerating users T , and between tolerating and frustrated users F [Apdex, Sevcik08]. There probably is a relationship between T and the response time goal, and between F and the response time requirement. However, while Apdex is a good metric for management and operations, it may be too high-level for engineering.

Agenda

- Metrics
- Elicitation
- Analysis and Specification
- Validation and Verification

Requirements Validation

- **Making sure that requirements are valid**
 - Quite often used to mean checking against test results (instead of verification)
- **Checking against different sources**
- **Reviews, modeling, prototyping, etc.**
- **Iterative process**
- **Tracing**
 - Tracing back to the original requirement

Requirements validation is making sure that requirements are valid. Unfortunately term 'validation' is quite often used to mean checking against test results instead of verification.

A good way to validate a requirement is to get it from different independent sources: if all numbers are about the same, it is a good indication that the requirement is probably valid. Validation may include, for example, reviews, modeling, and prototyping. Requirements process is iterative by nature and requirements may change with time, so it is important to trace requirements back to their source.

Requirements Verification

- **Checking if the system performs according to the requirements**
- **Both requirements and results should use the same aggregates to be compared**
- **Many tools measure only server time (or server and network)**
 - End user time may differ significantly, especially for rich web clients or thick clients
- **Both in load testing and production !**

Requirements verification is checking if the system performs according to the requirements. To make meaningful comparison, both the requirements and results should use the same aggregates. One consideration here is that load testing tools and many monitoring tools measure only server and network time. While end user response times, which business is interested in and usually assumed in performance requirements, may differ significantly, especially for rich web clients or thick clients due to client-side processing and browser rendering. Verification should be done using load testing results as well as during ongoing production monitoring. Checking production monitoring results against requirements and load testing results is also a way to validate that load testing was done properly.

Verification Issue

- Let's consider the following example
- Response time requirement is 99% below 5 sec
- 99% 3-5 sec, 1% 5-8 sec
 - Looks like a minor performance issue
- 99% 3-5 sec, 1% failed or had strangely high response times (more than 30 sec)
 - Looks like a bug or serious performance issue

Requirement verification presents another subtle issue: how to differentiate performance issues from functional bugs exposed under load. Often, additional investigation is required before you can determine the cause of your observed results. Small anomalies from expected behavior are often signs of bigger problems, and you should at least figure out *why* you get them.

When 99 percent of your response times are three to five seconds (with the requirement of five seconds) and 1 percent of your response times are five to eight seconds it usually is not a problem. But it probably should be investigated if this 1 percent fail or have strangely high response times (for example, more than 30 sec) in an unrestricted, isolated test environment. This is not due to some kind of artificial requirement, but is an indication of an anomaly in system behavior or test configuration. This situation often is analyzed from a requirements point of view, but it shouldn't be, at least until the reasons for that behavior become clear.

Requirements Verification: Performance vs. Bug

- Two completely different cases
 - Performance issue: business decision, cost vs. response time trade off.
 - Bug exposed under load: should be traced down first to make decision

46

These two situations look similar, but are completely different in nature:

1) The system is missing a requirement, but results are consistent: this is a business decision, such as a cost vs. response time trade off.

2) Results are not consistent (while requirements can even be met): that may indicate a problem, but its scale isn't clear until investigated.

Unfortunately, this view is rarely shared by development teams too eager to finish the project, move it into production, and move on to the next project. Most developers are not very excited by the prospect of debugging code for small memory leaks or hunting for a rare error that is difficult to reproduce. So the development team becomes very creative in finding "explanations". For example, growing memory and periodic long-running transactions in Java are often explained as a garbage collection issue. That is false in most cases. Even in the few cases when it is true, it makes sense to tune garbage collection and prove that the problem went away.

Another typical situation is getting some transactions failed during performance testing. It may still satisfy performance requirements, which, for example, state that 99% of transactions should be below X seconds – and the share of failed transaction is less than 1 percent. While this requirement definitely makes sense in production where we may have network and hardware failures, it is not clear why we get failed transactions during the performance test if it was run in a controlled environment and no system failures were observed. It may be a bug exposed under load or a functional problem for some combination of data.

The equipment is not operating as expected, and therefore there is a danger that it can operate with even wider deviation in this unexpected and not thoroughly understood way. The fact that this danger did not lead to a catastrophe before is no guarantee that it will not the next time, unless it is completely understood.

**Dr. Richard Feynman
Roger Commission Report on the
Challenger space shuttle accident**

47

When some transactions fail under load or have very long response times in the controlled environment and we don't know why, we have one or more problems. When we have unknown problems, why not trace it down and fix in the controlled environment? It would be much more difficult in production. What if these few failed transactions are a view page for your largest customer, and you won't be able to create any order for this customer until the problem is fixed? In functional testing, as soon as you find a problem, you usually can figure out how serious it is. This is not the case for performance testing: usually you have no idea what caused the observed symptoms and how serious it is, and quite often the original explanations turn out to be wrong.

Michael Bolton described this situation concisely [Bolton06]:

As Richard Feynman said in his appendix to the Rogers Commission Report on the Challenger space shuttle accident, when something is not what the design expected, it's a warning that something is wrong. "The equipment is not operating as expected, and therefore there is a danger that it can operate with even wider deviations in this unexpected and not thoroughly understood way. The fact that this danger did not lead to a catastrophe before is no guarantee that it will not the next time, unless it is completely understood." When a system is in an unpredicted state, it's also in an unpredictable state.

Summary

- **Specify performance requirements at the beginning of any project**
- **What to specify depends on the system**
 - Quantitative and measurable in the end
- **Elaborate and verify requirements throughout Development – Testing – Production**

48

We need to specify performance requirements at the beginning of any project for design and development (and, of course, reuse them during performance testing and production monitoring). While performance requirements are often not perfect, forcing stakeholders just to think about performance increases the chances of project success.

What exactly should be specified – goal vs. requirements (or both), average vs. percentile vs. APDEX, etc. – depends on the system and environment. Whatever it is, it should be something quantitative and measurable in the end. Making requirements too complicated may hurt. We need to find meaningful goals / requirements, not invent something just to satisfy a bureaucratic process.

If we define a performance goal as a point of reference, we can use it throughout the whole development cycle and testing process and track our progress from the performance engineering point of view. Tracing this metric in production will give us valuable feedback that can be used for future system releases.

Questions ?

Alexander Podelko

alex.podelko@oracle.com

apodelko@yahoo.com

@apodelko

*Links and references may be found in
the slide notes and at
www.alexanderpodelko.com*

49

References

[Apdex] Apdex Web site
<http://www.apdex.org/>

[Barber07] Barber, S. Get performance requirements right - think like a user. Compuware white paper, 2007.
http://www.perftestplus.com/resources/requirements_with_compuware.pdf

[Bass03] Bass L., Clements P., Kazman R. Software Architecture in Practice, Addison-Wesley, 2003.
<http://etutorials.org/Programming/Software+architecture+in+practice,+second+edition>

[Bickford97] Bickford P. Worth the Wait? Human Interface Online, View Source, 10/1997.
http://web.archive.org/web/20040913083444/http://developer.netscape.com/viewsource/bickford_wait.htm

[Bolton06] Bolton M. More Stress, Less Distress, Better Software, November 2006.
<http://www.stickyminds.com/sitewide.asp?ObjectId=11536&Function=edetail&ObjectType=ART>

[Cockcroft00] Cockcroft, A., Walker B. Capacity Planning for Internet Services. Quick planning techniques for high growth rates. Sun Microsystems, 2000.
<http://java.coe.psu.ac.th/SunDocuments/SunBluePrints/caphi.pdf>

[Forrester09] eCommerce Web Site Performance Today. Forrester Consulting on behalf of Akamai Technologies, 2009.
http://www.akamai.com/html/about/press/releases/2009/press_091409.html

[Miler68] Miller, R. B. Response time in user-system conversational transactions, In Proceedings of the AFIPS Fall Joint Computer Conference, 33, 1968, 267-277.

[Microsoft10] Mailbox Server Processor Capacity Planning. Microsoft, 2010.
<http://technet.microsoft.com/en-us/library/ee712771.aspx>

[Martin86] Martin, G.L. and Corl, K.G. System response time effects on user productivity, Behavior and Information Technology, 5(1), 1986, 3-13.

[Nielsen94] Nielsen J. Response Times: The Three Important Limits, Excerpt from Chapter 5 of Usability Engineering, 1994.
<http://www.useit.com/papers/responsetime.html>

[Performance07] Performance Testing Guidance for Web Applications, 2007.
<http://perftestingguide.codeplex.com/releases/view/6690>

[Podelko07] Podelko A. Multiple Dimensions of Performance Requirements, CMG, 2007.
http://www.alexanderpodelko.com/docs/Performance_Requirements_CMG07.pdf

[Sevcik03] Sevcik, P. How Fast Is Fast Enough, Business Communications Review, March 2003, 8-9.
http://www.bcr.com/architecture/network_forecasts%10sevcik/how_fast_is_fast_enough?_20030315225.htm

[Sevcik08] Sevcik, P. Using Apdex to Manage Performance, CMG, 2008.
<http://www.apdex.org/documents/Session318.0Sevcik.pdf>

[Simmons01] Simmons E. Quantifying Quality Requirements Using Planguage, Quality Week, 2001.
http://www.clearspecs.com/downloads/ClearSpecs20V01_Quantifying%20Quality%20Requirements.pdf

[Smith02] Smith C., Williams L. "Performance Solutions", Addison-Wesley, 2002.

[SWEBOK04] Guide to the Software Engineering Body of Knowledge (SWEBOK). IEEE, 2004.
<http://www.computer.org/portal/web/swebok>