

Using Event- and Actor-driven Paradigms to Increase Web Server Performance

FELIX HESSENBERGER

MASTERARBEIT

eingereicht am
Fachhochschul-Masterstudiengang

INTERACTIVE MEDIA

in Hagenberg

im Juni 2014

© Copyright 2014 Felix Hessenberger

This work is published under the conditions of the *Creative Commons License Attribution–NonCommercial–NoDerivatives* (CC BY-NC-ND)—see <http://creativecommons.org/licenses/by-nc-nd/3.0/>.

Declaration

I hereby declare and confirm that this thesis is entirely the result of my own original work. Where other sources of information have been used, they have been indicated as such and properly acknowledged. I further declare that this or similar work has not been submitted for credit elsewhere.

Hagenberg, June 30, 2014

Felix Hessenberger

Contents

Declaration	iii
Abstract	vi
Kurzfassung	vii
1 Introduction	1
1.1 Motivation	1
1.2 Objective	3
1.3 Structure	4
2 Technical Background	5
2.1 Terms and Definitions	5
2.1.1 Network Communication	5
2.1.2 Dynamic Content	5
2.1.3 Asynchronous Requests	6
2.1.4 Request Frequency and Response Time	6
2.1.5 Scalability	7
2.1.6 Development	8
2.2 Program Flow	8
2.2.1 Threads	8
2.2.2 Events	10
2.2.3 SEDA	11
2.2.4 Actors	11
2.2.5 Reactive Architecture	11
3 State of the Art	12
3.1 Event-based	12
3.1.1 Express	12
3.1.2 Watson	14
3.1.3 Twisted	15
3.2 Actor-based	16
3.2.1 Spray	16
3.2.2 Lift	18

3.2.3	Play! Framework	20
3.3	Other	25
3.3.1	Ruby on Rails	25
3.3.2	Node.scala	27
4	Implementation	32
4.1	Considerations	32
4.2	Synchronous Version	33
4.3	Asynchronous Version	36
5	Evaluation	40
5.1	Approach	40
5.2	Performance	43
5.3	Maintenance	47
5.4	Results	49
6	Conclusion and Future Development	52
	References	55
	Literature	55
	Online Sources	55

Abstract

Over the past years websites have advanced from merely displaying content to representing interfaces to dynamic server-side applications of various scales; other environments like mobile platforms tend to use the same HTTP interfaces as well. To limit the cost of server hardware various software-based approaches aim to maximise the number of simultaneous operations by shifting from the classic per-request threading model to more sophisticated concurrency patterns. This thesis presents and compares a number of different approaches to server-side concurrency implementations from the view of a programmer. Typical use-cases for server-side information flow are contrived and evaluated regarding asynchronous processing. Patterns are then reviewed based on their performance in scalable high-throughput networking applications by the example of live applications as well as experimental settings.

Kurzfassung

Franz rast im komplett verwehrlosten Taxi quer durch Bayern.

Chapter 1

Introduction

1.1 Motivation

a

1. Introduction

3

a

1.2 Objective

a

1.3 Structure

Chapter 2

Technical Background

2.1 Terms and Definitions

A Web server can be utilised to handle rather different tasks, from serving entire Web pages to representing an endpoint for raw data retrieval to merely delivering static assets like images. This section aims to give an overview of the basic requirements a modern Web server architecture needs to fulfil. Moreover, important performance factors are elaborated with regard on high-demand setups.

2.1.1 Network Communication

The eponymous task of a Web server is to serve Web-connected clients over the medium of the Internet. This involves receiving and sending messages using different implementations of network protocols. The most widely used protocol of the Web, *HTTP*¹, is a request-response protocol, which means that for every message a client sends to a server, a response is sent back [12]. To minimise networking latency, it is preferable for a Web server to have a high-speed connection to the Internet, fast system I/O² and capable routing hardware. However, these parameters are not directly related to software and are thus neglected during the further course of this thesis.

2.1.2 Dynamic Content

Originally, the Web was intended to be a network of interconnected text files, which later was augmented with images and style sheets; Web servers were basically required to understand requests and respond with static content

¹Hypertext Transfer Protocol

²Input and Output, esp. hardware

accordingly [12]. With the release of *PHP*³, *ASP*⁴ and *Java*⁵ – 1995, 1996 and 1997, respectively – dynamic webpages, i.e. views that are prepared by the server based on dynamic data like database content, became widespread [20]. From that point on, Web servers needed more processing capabilities for script execution and database access; however, the number of requests remained roughly the same, except for occasional form submissions [20].

2.1.3 Asynchronous Requests

The advent of *AJAX*⁶ and mobile applications in the late 2000's changed requirements drastically. Rather than refreshing the whole view for every piece of information sent and received, data could now be transferred in a more granular fashion. By asynchronously communicating with an API⁷ endpoint in the background, operations like deleting an item from a list could be performed invisibly and ubiquitously without reloading the page context. Especially applications that aim to provide desktop-like capabilities – commonly called Rich Internet Applications – make heavy use of asynchronous requests [10, p. 4]. This inherently also changed users' expectations for websites from anticipating a certain amount of load time to implicating real-time behaviour [13]. To achieve low latency while maintaining client-server information consistency, the server's performance has to meet the combined request frequency of all clients at a given point in time.

2.1.4 Request Frequency and Response Time

Since in many cases the responsiveness of the user interface depends on the duration of the server communication roundtrip, maintaining acceptable response times is often crucial [8, p. 1]. Request frequency and response time correlate in the sense that request frequency represents the demand on a server endpoint while response time – given equally demanding operations per request – can be interpreted as the potential of the server to meet the demand. When the processing limit of the server is met, response times become generally inversely proportional to the request frequency, as illustrated in figure 2.1 [18]. At this point, the server may choose to neglect the request (ideally by returning the status code *503 Service Unavailable* [12]), not respond at all or even stop serving clients altogether (i.e. “crash”).

³Recursive acronym: PHP Hypertext Preprocessor, <http://php.net/>

⁴Active Server Pages, <http://msdn.microsoft.com/en-us/library/aa286483.aspx>

⁵<https://www.java.com/>

⁶Asynchronous JavaScript and XML (Extensible Markup Language)

⁷Application Programming Interfaces

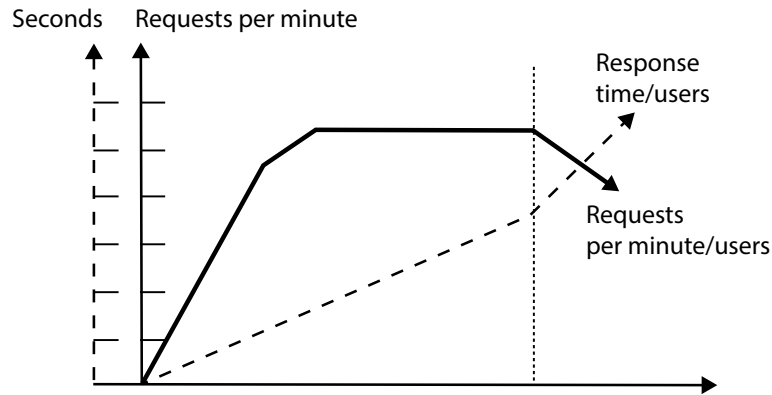


Figure 2.1: Correlation between request frequency and response time in a typical Web server setup. After the server has reached its limit of linearly serving clients (indicated by the dotted line), response times become inversely proportional to the request frequency. Image source: [18]

2.1.5 Scalability

Demands on Web servers typically are lower during the initial phase of a business and grow with the popularity of the service. Since business growth and server load can not be exactly predicted, it is necessary to be able to adjust (i.e. *scale*) the entire server architecture according to current needs in a timely manner. The *Slashdot Effect* describes a sudden spike in service popularity and can, due to the open nature of the Web, lead to a tremendous increase in activity over a relatively short timespan [2, p. 1].

Today's hardware is well suited to meet high demands and can be configured flexibly: If a larger number of physical server units as well as the necessary infrastructure is available, requests can be distributed and the load on a single unit decreases. If single units are outfitted with more memory and faster processors, the number of request operations per unit increases. Since acquiring and maintaining server units and other infrastructure components is expensive, well-designed software can make a significant difference in system efficiency, which in turn can greatly benefit any business [6, p. 11].

Ideally, the server software should be hardware agnostic, i.e. should behave consistently independent of the hardware it runs on. For instance, if the software depends heavily on sharing application state via RAM⁸, scaling out on more than one machine will be unsuccessful [11]. Scalability can be measured by the relationship between hardware resources and the increase of performance. If this relationship is nearly linear, the system can be considered to scale well.

⁸Random Access Memory

2.1.6 Development

Not a part of the production system itself, but nonetheless a vital part of all Web server applications is their development. A structured, idiomatic way of writing application logic doubtlessly contributes to every software product. Modularisation of components facilitate the use of third-party software like libraries and frameworks. In return, using existing solutions can greatly reduce development time and effort, while simultaneously providing proven solutions. Web server applications particularly benefit from frameworks since they often handle standard tasks like network I/O², database access and caching[9, Foreword]. Integrating and maintaining these frameworks is a major part in implementing a Web server application; thus, not only the performance, but also the ease of use of selected frameworks and their language environments by the developer are treated as criteria in this thesis.

2.2 Program Flow

Since a Web application in production is usually publicly accessible, serving multiple clients simultaneously is the rule, rather than the exception. Depending on the popularity of the service, the number of concurrent requests can range anywhere from dozens to several thousands, e.g. for social media sites [2, p. 1]. A server process with a single flow of control would only be able to serve one client at once, with all requests received while the server is busy being neglected. Therefore, networking applications always have to be implemented with multiple program flows that can be executed concurrently [22]. This section lists various paradigms associated with designing an application capable of maintaining multiple flows of control.

2.2.1 Threads

A thread is a sequence of instructions within a program. Allocating processing time to threads is handled by an operating system scheduler. To have a program execute multiple logic structures concurrently, they have to be explicitly abstracted in the form of threads. Physical concurrency occurs, when threads are executed simultaneously – i.e. at the exactly same time – on different processor cores; in contrast, logical concurrency describes that multiple threads are executed sequentially at roughly the same time, thus giving the impression of simultaneous execution. Physical concurrency is inherently more efficient [17].

A great advantage of threads in the context of Web server applications lies in the natural abstraction level regarding multiple parallel requests: Client communication is commonly treated as a set of mutually independent connections; this approach of abstraction facilitates a clear program flow structure

[11]. Accordingly, every request can be treated as an isolated flow of control [1, p. 2]. However, since threads are not isolated from each other and share state via a common memory address space, this only holds true as long as resources like queues or database components are accessed sequentially. Thus, close attention has to be paid by the developer to avoid race-conditions, deadlocks and access violations – complications that result from improper thread coordination [3, p. 1]. Generally, the implementation of large-scale systems heavily relying on threads – as an evolutionary improvement from sequential computing – always introduces additional complexity; Edward Lee even goes as far as declaring it to be nondeterministic [7, p. 1].

Traditionally, Web server applications process each request on a dedicated thread throughout its whole lifespan, from accepting it to responding to it [4, p. 162]. This behaviour can be observed for instance in implementations of the popular LAMP⁹ server stack configuration [4, p. 48]. A less experienced programmer might find this ideal, since concurrency stays mostly hidden and the application logic is orientated on the flow of a single request – smaller projects might not experience any drawbacks of this setup at all. However, it is obvious that to scale up a thread-based system the number of threads has to be increased. The number of threads engaging in simultaneous processing, i.e. physical concurrency, is limited by the number of processing cores. This means that on a computer equipped with a quad-core processor, four threads can be executed – and thus, four requests can be served – in parallel¹⁰.

Problems arise when a thread has to wait for another requirement to be fulfilled. The process of meeting a requirement that renders the executing thread unable to proceed is called a *blocking* operation. Such actions include for instance reading or writing a file on mass storage, handling network traffic or file uploads, querying a database, accessing another Web service or doing intensive computations [4, p. 196]. When a thread encounters a blocking operation, it cannot advance in the program flow until the operation completes; the resulting delay can account to anywhere from a few milliseconds to several seconds, for instance when accessing a slow or unresponsive Web service. The only way to counteract the temporary occupation of threads and to continue processing incoming requests is the creation of new threads [6, p. 36]. However, every newly created thread counts towards certain limitations in scalability. On the one hand every thread receives a predefined share of process address space memory – also known as *stack* – upon creation to temporarily store data [21]; since memory is reserved in advance without knowing the exact requirements of the thread, a certain amount of memory overhead is likely. On the other hand, the entirety of all threads has

⁹Linux, Apache, MySQL, PHP

¹⁰Certain implementations of simultaneous multithreading allow for increasing this number at the cost of reduced performance per thread, for instance Intel's Hyper-Threading Technology (<http://www.intel.com/>).

to be orchestrated by an operating system module called *scheduler*, which requires processing time relative to the number of threads [21]. Moreover, a computationally expensive procedure called *context switching* must also be followed upon changing the actively processed thread [16].

Some of the problems of threads can be addressed by using a *thread pool*: Instead of spawning new threads upon each request, a fixed number of threads is spawned in advanced and workload is distributed among them. However, this procedure is not without problems and introduces the delicate step of setting the thread pool size [14]. It can be concluded that a lower thread overhead can benefit the overall performance of a process. Furthermore, when scaling an application, the maximum number of simultaneously processed threads can at best only increase linearly in relation to the number of processing cores in a system [15].

2.2.2 Events

While threads present a natural abstraction for handling Web requests, recent years have seen an incline towards event-driven flow of control. Seen from a different perspective, events are at least equally idiomatic: The Web server has no control over the arriving requests, yet it has to respond by executing application logic. Instead of forcefully maintaining control over the execution context, the Web server may relinquish control and let itself be controlled by events. This strategy follows the principle of *inversion of control* [5]. At its simplest, an event-driven application consists of two major components: On the one hand an *event loop* containing an *event listener* and on the other hand an *event handler*. The event loop is a lightweight structure passing incoming events from a queue to event listeners that have subscribed to a certain kind of event, e.g. an incoming network request. The targeted event listener then passes the event on to a handler function, which executes necessary application logic and may create another event upon completion. In the case of an event-driven HTTP server this returned event may contain an HTTP response to be sent to the initial client. Larger applications typically have one event loop per process and a number of listeners and handlers [6, p. 33]. For an illustration of a basic event-driven application setup, see figure 2.2.

Event-driven programming does not preclude the existence of threads; neither is it the opposite or an evolutionary step. All major operating systems use threads as a means of managing process execution, thus even a purely event-driven program needs at least one thread. However, this is not a favourable scenario, especially for a Web server: If the event loop and the handler would run on the same thread, the event loop would block if the handler blocks and thus would not be able to accept further events. Therefore the event loop (and with it the application I/O²) commonly runs on a dedicated thread with the handlers run on other threads or – in more

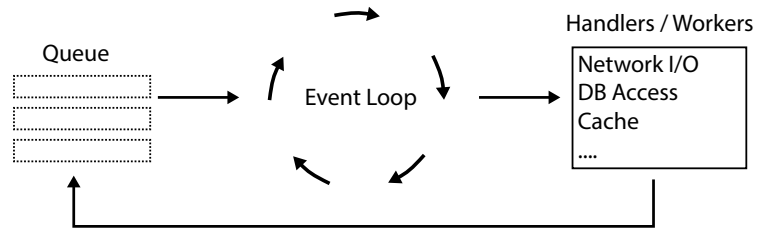


Figure 2.2: Basic flow of control in an event-driven application. Operations that would normally block the event loop are executed separately and create further events upon completion.

sophisticated implementations – in a thread pool (see section 2.2.1).

2.2.3 SEDA

In this thesis, the term *event-driven* always describes intra-system event architecture, also known as *staged event-driven architecture (SEDA)* [19].

When incoming demand exceeds system processing capacity, it can simply queue up requests in its internal buffer queues.

2.2.4 Actors

2.2.5 Reactive Architecture

Chapter 3

State of the Art

3.1 Event-based

3.1.1 Express

a

a

3.1.2 Watson

a

3.1.3 Twisted

a

3.2 Actor-based

3.2.1 Spray

a

a

3.2.2 Lift

a

a

3.2.3 Play! Framework

a

a

a

a

a

3.3 Other

3.3.1 Ruby on Rails

a

a

3.3.2 Node.scala

a

a

a

a

Chapter 4

Implementation

4.1 Considerations

a Just an example, no Real appPlay Setup, Threadpools, Local Testing,
Cloud Deployment

4.2 Synchronous Version

a

a

a Prerequisites, ...

4.3 Asynchronous Version

a

a

a

Chapter 5

Evaluation

5.1 Approach

a

a

a
Methodology

5.2 Performance

a

a

a

a JMeter loader.io

5.3 Maintenance

a

a

5.4 Results

a

a

Chapter 6

Conclusion and Future Development

a

a

References

Literature

- [1] Rob Von Behren, Jeremy Condit, and Eric Brewer. “Why Events Are A Bad Idea (for high-concurrency servers)”. In: *HotOS IX : The 9th Workshop on Hot Topics in Operating Systems* (2003).
- [2] Vaarnan Drolia, Cedric Ansley, and Chin Shen. “Threads vs Events for Server Architectures”. 2010.
- [3] Jeffrey Fischer, R Majumdar, and Todd Millstein. “Tasks: language support for event-driven programming”. In: *PEPM '07 Proceedings of the 2007 ACM SIGPLAN symposium on Partial evaluation and semantics-based program manipulation* (2007), pp. 134–143.
- [4] Cal Henderson. *Building Scalable Web Sites*. O'Reilly Series May. O'Reilly Media, Inc, 2006, p. 330.
- [5] Gregor Hohpe. “Programming Without a Call Stack â Event-driven Architectures”. In: *Enterprise Integration Patterns* (2006).
- [6] Tom Hughes-Croucher and Mike Wilson. *Node - Up and Running*. 2012.
- [7] Edward A. Lee. “The problem with threads”. In: *Computer* 39.5 (May 2006), pp. 33–42.
- [8] S Nadimpalli and S Majumdar. “Techniques for Achieving High Performance Web Servers”. In: *Parallel Processing, 2000. . . .* (2000).
- [9] Alexander Reelsen. *Play Framework Cookbook*. Packt Publishing, 2011.
- [10] Sencha Inc. “Web Applications Come of Age”. 2011.
- [11] Bryan Veal and Annie Foong. “Performance Scalability of a Multi-Core Web Server”. In: *Proceedings of the 3rd ACM/IEEE Symposium on Architecture for networking and communications systems - ANCS '07* (2007), p. 57.

Online Sources

- [12] R Fielding, U C Irvine, and J Gettys. *Hypertext Transfer Protocol – HTTP/1.1*. 1999. URL: <http://www.ietf.org/rfc/rfc2616.txt>.

- [13] Jesse James Garrett. *AJAX : A New Approach to Web Applications*. 2005. URL: <http://www.adaptivepath.com/ideas/ajax-new-approach-web-applications/>.
- [14] Brian Goetz. *Thread Pools and Work Queues*. 2002. URL: <http://www.ibm.com/developerworks/java/library/j-jtp0730/index.html>.
- [15] Michael McCool. *The Serious Drawbacks of Explicit Multi- Threading*. 2008. URL: <http://software.intel.com/en-us/blogs/2008/06/05/nitrogen-narcosis-part-ii-the-serious-drawbacks-of-explicit-multi-threading/>.
- [16] Mark McGranahan. *Threaded vs Evented Servers*. 2010. URL: <http://mmcgrana.github.io/2010/07/threaded-vs-evented-servers.html>.
- [17] Svetlin Nakov. *Internet Programming with Java*. 2004. URL: <http://www.nakov.com/inetjava/lectures/part-1-sockets/InetJava-1.3-Multithreading.html>.
- [18] Oracle. *Establishing Performance Goals*. 2010. URL: <http://docs.oracle.com/cd/E19900-01/819-4741/fygaj/index.html>.
- [19] Michael Peterson. *Events and Event-Driven Architecture : Part 1*. 2012. URL: <http://thornydev.blogspot.co.at/2012/01/events-and-event-driven-architecture.html>.
- [20] Pingdom. *A History of the Dynamic Web*. 2007. URL: <http://royal.pingdom.com/2007/12/07/a-history-of-the-dynamic-web/>.
- [21] Mark Russinovich. *Pushing the Limits of Windows : Processes and Threads*. 2009. URL: <http://blogs.technet.com/b/markrussinovich/archive/2009/07/08/3261309.aspx>.
- [22] Webopedia. *Multi-core Technology*. 2007. URL: http://www.webopedia.com/TERM/M/multi_core_technology.html.