

## Web Scaling Frameworks PhD Transfer Event



Thomas Fankhauser Qi Wang **Christos Grecos** Xinheng Wang

Ansgar Gerlicher



### GUEUE Original Idea

-se

### Master's Thesis





WORKER

Thomas Fankhauser 2012





# Project Roadmap

### Work Packages & Publications



### Web Scaling Frameworks: IEEE ICC Conference Paper

### Web Scaling Frameworks: A novel class of frameworks for scalable web services in cloud environments

Thomas Fankhauser<sup>\*†</sup>, Student Member, IEEE, Qi Wang<sup>\*</sup>, Member, IEEE, Ansgar Gerlicher<sup>†</sup>, Member, IEEE, Christos Grecos\*, Senior Member, IEEE and Xinheng Wang\*, Member, IEEE

Abstract—The social web and huge growth of mobile smart devices dramatically increases the performance requirements for web services. State-of-the-art Web Application Frameworks (WAFs) do not offer complete scaling concepts with automatic response times. These functionalities, however, are supported by cloud computing and needed to scale an application to its demands. Components like proxies, load-balancers, distributed caches, queuing and messaging systems have been around for a long time and in each field relevant research exists. Nevertheless, to create a scalable web service it is seldom enough to deploy only one component. In this work we propose to combine those complementary components to a predictable, composed system. The proposed solution introduces a novel class of web frameworks called Web Scaling Frameworks (WSFs) that take over the scaling. The proposed mathematical model allows a universally applicable prediction of performance in the single-machine- and multi-machine scope. A prototypical implementation is created to empirically validate the mathematical model and demonstrates both the feasibility and increase of performance of a WSF. The results show that the application of a WSF can triple the requests handling capability of a single machine and additionally reduce the number of total machines by 44%.

### I. INTRODUCTION

The enormous growth of smart mobile devices in combination with social web services increases the number of requests that need to be processed by modern web platforms in a timely fashion. Whereas cloud computing provides the ability to provision the hardware needed, state-of-the-art Web Application Frameworks (WAFs) do not offer integrated scaling concepts to deal with automatic resource-provisioning and elastic caching or ensure a guaranteed maximum response time

They are rather designed to abstract common functionalities needed for web application development including datamanagement, url-mapping, session-handling and responsegeneration. Today, users progressively access the social web from anywhere using their mobile smart devices, which leads to increased traffic. A single computing resource might not be able to satisfy such an amount of requests - only the junction of multiple computing resources, where each resource gets a small share of the total requests, allows to handle

such huge amounts of requests in aggregation. Handling the exponentially increasing global requests adds the requirement of being able to run multiple instances of an application for highly scalable web services. The major challenges that resource-provisioning, elastic caching or guaranteed maximum are introduced by this requirement are the management of the shared resources, the balancing of the requests among all instances and the decision when to spawn or terminate instances. These challenges are collectively referred to as horizontal scaling [13], [14], [16].

> Our experiments have showed that WAFs have different strengths and weaknesses. A highly abstracted WAF like Ruby on Rails, for example, was slower than the very thin WAF node.js but more powerful regarding data management and interface rendering. If a web service needs to provide both a fast and slim JSON API and a full blown HTML website it is the best solution to combine both WAFs. As both the horizontal scaling and web service composition are very complex matters, it makes sense not to introduce them to WAFs but offload them to another layer - the Web Scaling Framework (WSF) proposed in this paper. Fig. 1 illustrates a WSF that incorporates multiple WAF applications.



Fig. 1. The relationship between the WSF and WAFs

To comply to a proposed class of WSFs, a WSF should:

- take over the responsibilities of scaling and incorporate existing WAFs
- separate the business logic in the web service from the hosting logic
- connect to and combine existing WAFs to a compound web service using standard HTTP requests
- introduce low overhead when added, whilst adding the instant ability to scale
- constantly adapt their infrastructure to fit the required performance at all times

### IEEE International in Sydney, Australia







<sup>\*</sup>School of Computing, University of the West of Scotland, Email: {Thomas.Fankhauser, Qi.Wang, Christos.Grecos, Xinheng.Wang}@uws.ac.uk <sup>†</sup>Mobile Application Development, Stuttgart Media University, Email: {fankhauser, gerlicher}@hdm-stuttgart.de





### Web Scaling Frameworks: IEEE ICC Conference Paper







# Project Roadmap

### Work Packages & Publications



### Web Scaling Frameworks: IEEE TSC Journal Paper

DRAFT 0.5: IEEE TRANSACTIONS ON SERVICES COMPUTING, VOL. X, NO. Y, JANUARY 2015

### Introducing Elastic Scalability to Web Services in the Cloud with Web Scaling Frameworks

Thomas Fankhauser, Student Member, IEEE, Qi Wang, Member, IEEE, Ansgar Gerlicher, Member, IEEE, Christos Grecos, Senior Member, IEEE, and Xinheng Wang, Member, IEEE

Abstract—In the social web, web services have to accommodate a significant number of requests due to the high interactivity of current applications. Applications have to be built in a scalable fashion so the number of machines can be adapted to highly dynamic traffic situations. In the deployment context, web service providers often need to focus initially on the business logic which prevents detailed scalability considerations. If however a critical mass of customers is reached, providers need to be able to scale-up immediately their web services to stay in business. State-of-the-art Web Application Frameworks (WAFs) focus on the creation of application logic including data validation, view composition and session handling. However, they don't offer integrated cloud scaling concepts that handle the automatic provisioning of resources. Web service providers have to create custom-built systems that consider scalability issues manually. As the creation of such scaling-systems is a very complex, we proposed in our recent work the concept of Web Scaling Frameworks (WSFs) in order to offload scaling to another layer of abstraction. WSFs are composed of traditional WAFs with multiple other components to provide scalability right from the launch of the deployment cycle. In this work, a detailed design for WSFs including necessary modules, interfaces and components is presented. A mathematical model used for performance rating is evaluated and enhanced on a computing cluster of 42 machines. Traffic traces from over 25 million real-world applications are analysed and evaluated on the cluster to compare the WSF performance with a traditional scaling approach using WAFs and caches. The results show that the application of WSFs can reduce the number of total machines needed for three representative real-world applications a social network, a trip planner and the FIFA World Cup 98 website - by 32%, 63% and 92% respectively.

Index Terms—scalability, web service, cloud computing, web scaling frameworks, performance evaluation, software architecture

### **1** INTRODUCTION

deal with a high level of interactivity in applications, which which often prevents detailed scalability considerations. in turn introduces enormous amounts of requests. Static State-of-the-art Web Application Frameworks (WAFs) and find intelligent routes based on live data.

The above scenarios introduce new challenges to web scalability issues manually. service providers and developers. As single-server systems

- T. Fankhauser, Q. Wang and X. Wang are with the School of Computing, University of the West of Scotland
- E-mail: {Thomas.Fankhauser, Qi.Wang, Xinheng.Wang}@uws.ac.uk • C. Grecos is an Independent Imaging Consultant.
- E-mail: grecoschristos@gmail.com • *A. Gerlicher is with the Media University Stuttgart.* E-mail: gerlicher@hdm-stuttgart.de

THE demands for modern web services increase due to mass is reached, the web services need to be able to scale-**I** the soaring social nature of the web and the upsurge up immediately to stay in business. Before this threshold of the total number of mobile devices. Web services have to is reached, providers need to focus on the business logic,

websites are replaced by dynamic and highly interactive are designed to abstract common functionalities needed for applications. For instance, TV shows deploy apps that allow the efficient implementation of web services. They focus users to influence the course of the show, advertisements are on the creation of application logic, data structures, data brought to customers only if they remain in the vicinity of validation, view layer presentation and session handling. advertised target locations, smart sensors deliver data for They don't offer integrated scaling concepts that handle the all kinds of metrics which users see on their mobile devices, provisioning of resources, manage elastic caching or ensure and cars communicate traffic situations, report traffic jams guaranteed response times. Today, web service providers have to create custom-built systems that consider these

The creation of such scaling systems is a very complex are not able to handle the increased load, applications matter. Hence, we proposed the introduction of Web Scaling need to be built in a scalable fashion. Requests have to be Frameworks (WSFs) in our recent work [1]. WSFs offload balanced over all available machines, resources need to be scaling to another layer of abstraction. They take over the shared without conflicting versions, distributed transactions responsibilities of scaling by embedding existing WAFs in have to be processed in a fault tolerant manner, and the a larger system. The application logic stays on the side of number of machines has to be adapted to highly dynamic the WAFs while the scaling logic is provided by the WSF. traffic situations. Typically, web service providers need to Fig. 1 illustrates the interplay between a WSF and multireach a critical mass of users to be successful. If the critical ple WAFs. To utilise an existing interface, the frameworks communicate with each other using HTTP. WSFs provide instant scalability to common WAF applications whilst only introducing a low overhead. To meet the performance requirements at all times, the infrastructure is adapted automatically. Resources are provisioned on a pay-per-use basis to benefit from the concept of cloud computing. WSFs are able to transparently use Software-as-a-Service (SaaS) or machine-cluster components.

### Submitted to:





























### Web Scaling Frameworks: IEEE TSC Journal Paper

Normal Version vs. Scaled Version

Evaluation





# Project Roadmap

### Work Packages & Publications



Reactive/Proactive Post-Processing: Conference Paper







Write Up

### Reactive/Proactive Post-Processing: Conference Paper





Dependency Analysis Data Structures Declaration Generation Visualisation Link Analysis



Reactive/Proactive Post-Processing: Conference Paper

Pre/Post-Processing Synchronous vs. Asynchronous Optimisation Parallelisation Fragmentation **Eventual Cache** 





Web Scaling Frameworks: Thesis

Thesis focus: Create and evaluate a full-stack Web Scaling Framework with dependency declaration and optimised post-processing algorithms







2016



## Web Scaling Frameworks PhD Transfer Event



Thomas Fankhauser Qi Wang **Christos Grecos** Xinheng Wang

Ansgar Gerlicher

