

Implementation Of A Scalable Architecture

David Alejandro Vasquez Carreño
Ingeniería de sistemas
Escuela Colombiana de Ingeniería Julio Garavito
Bogotá, Colombia
david.vasquez@mail.escuelaing.edu.co

Abstract—In today's world, most companies support their products and services in IT systems, like web servers, bank functions or cloud gaming, but there is a more complex architecture behind these services, that allow companies to provide continuous and scalable architectures to the growing set of customers.

Index Terms—Server, CPU, Load.

I. INTRODUCTION

Complex systems have mechanisms for allowing more users, for being able to handle all requests for the company. In data centers, there are scalable mechanisms like load balancers or backup servers that are there for failings in the system. This allows the company to support services 24/7 without major interruptions, except for maintenance of servers, that is previously negotiated with the provider.

Concurrent requests from users are handled by a capable data center, that provides multiple servers in a paid plan, but these systems are often limited or can be sensible to an excessive number of requests, that are not specially planned by the company, like special offers in black Friday, that can take websites to have more users than expected, taking servers to the limit, in networking and CPU load. Data centers try to make the server tasks easier, by distributing the requests between more servers that are not handling that load, so the users can't notice so much concurrence in the system.

Another mechanism, that is the one we are going to be discussing, is the one to launch more servers depending on the CPU load or network load of a system, so more servers are available in case of the increment of users' requests, like in our example of black Friday customers going to buy at the same time at a digital store.

II. WHAT IS A SCALABLE ARCHITECTURE?

Scalability is a property of the systems to adapt itself and react without losing quality. In our case of study, we can look at network, processes or systems, that must change based on the environment and parameters, and make changes to maintain quality, [1]

The image [1] shows an example of horizontal scaling, that consists in adding more servers to the layout, keeping the properties of those, to be equal to the already existing ones.

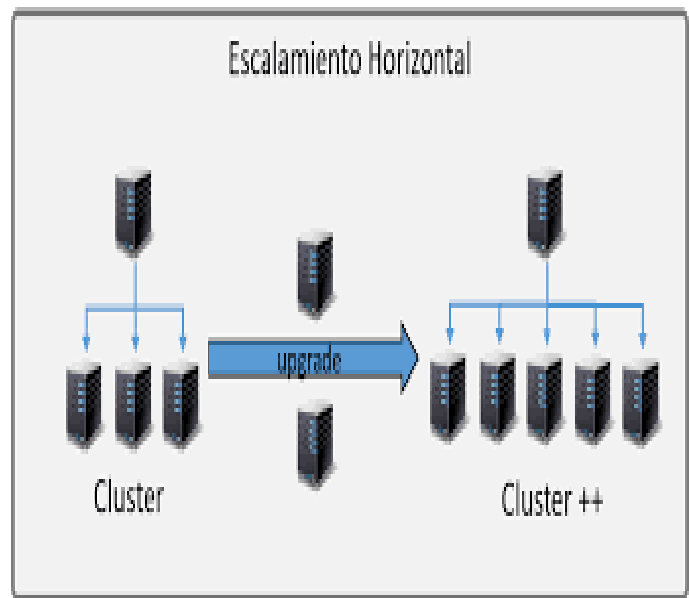


Fig. 1. Horizontal scaling

Another example of scaling, is vertical scaling, that allows us to equally add more servers, but in this case, the servers are better, have better processing capabilities or networking interfaces, etc.

III. LOAD CRITERIA

For being able to determine whether a server has a high load, we must choose a criteria for making that decision. Some examples of server load are the following:

- CPU load - Most used
- Memory usage
- Disk usage
- Memory usage
- Network usage
- GPU load - Rare

The figure 2 shows an example of a linux tool for checking system properties like CPU, memory and other characteristics of a linux machine, that is the most common thing nowadays in servers. For monitoring these properties, more complex

