

# Software Architecture

## Distributed Architecture

# Dispatcher (Load Balancer) Architecture Style

Eunmi Choi  
Kookmin University

### Types of Distributed Architecture Style

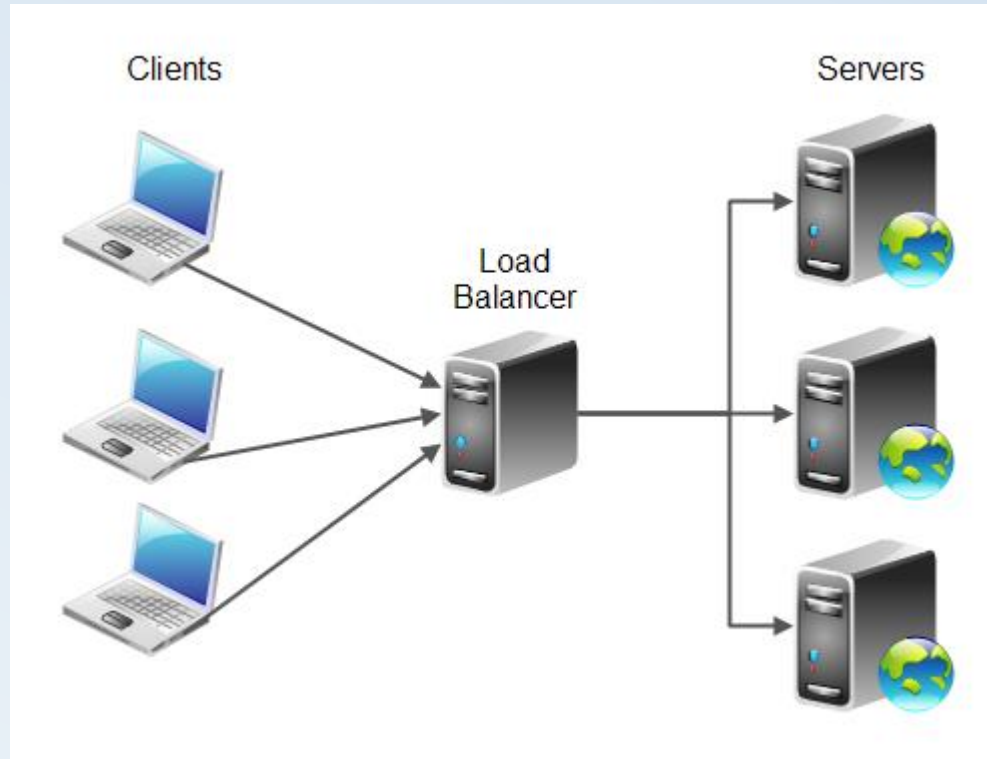
- Client-server
- Multitier
- Proxy
- Dispatcher (Load Balancer)
- P2P
- Broker
- Service-oriented architecture
- Microservice architecture



# Dispatcher(Load Balancer) Architecture Style

- Synopsis

- a middleware architecture used in distributed computing to make a server farm appear to clients as a single server so that the server farms achieve high scalability and high availability through server load balancing



Farm, (Server  
)  
(  
가 )  
가  
!



# Dispatcher(Load Balancer) Architecture Style: Context

- A single server that handles all the incoming requests may not have the capacity to handle high traffic volumes once the business becomes popular.<sup>가</sup>
  - In such a situations companies often start to *scale up*: they upgrade the existing infrastructure by buying a larger box with more processors or add more memory to run the applications.<sup>6</sup>
  - Scaling up, though, is only a short-term solution. And it's a limited approach because the cost of upgrading is disproportionately high relative to the gains in server capability.
- Scaling up, though, is only a short-term solution. And it's a limited approach because the cost of upgrading is disproportionately high relative to the gains in server capability.
  - For these reasons most successful Internet companies follow a *scale out* approach.<sup>d4</sup>
  - Application components are processed as multiple instances on server farms, which are based on low-cost hardware and operating systems. As traffic increases, servers are added.

# Dispatcher(Load Balancer) Architecture Style: Forces

- On the software side, you must design applications so that they can run as multiple instances on different servers.
  - You do this by splitting the application into smaller components that can be deployed independently.

HTTP  
Stateful

?

(Stateless).

!

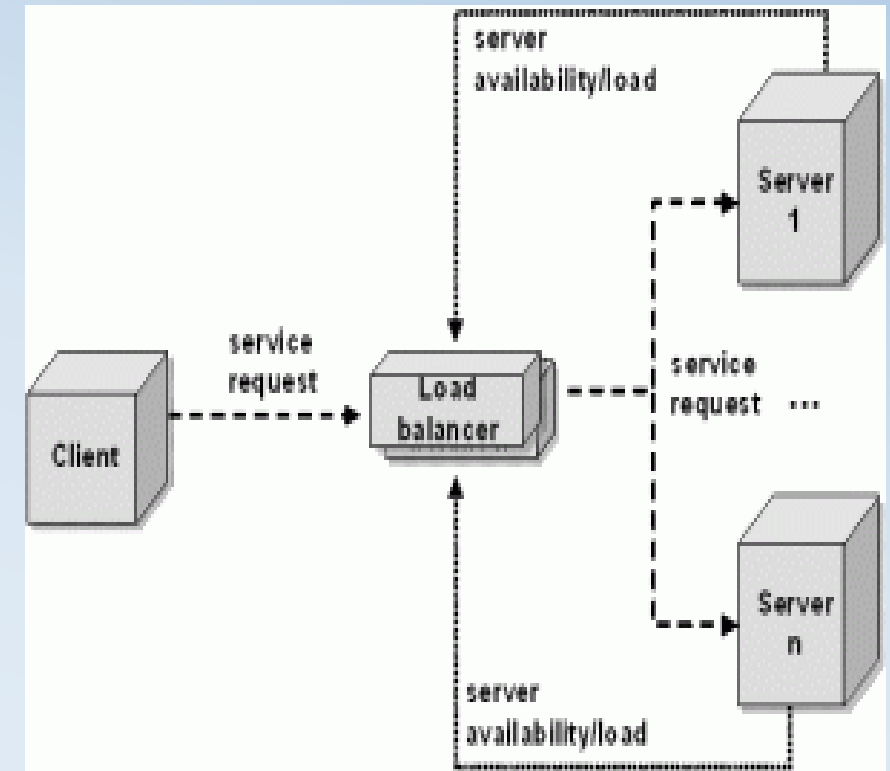
- A more challenging problem arises when the application components are stateful with application-session data.
- On the infrastructure side, the processing load must be distributed among the group of servers (Server Load Balancing).

# Dispatcher(Load Balancer) Architecture Style - Solution

- Structure

- In a widely used server load balancing architecture, the incoming request is directed to a dedicated server **load balancer** that is transparent to the client.

- Based on parameters such as availability or current server load, the load balancer decides which server should handle the request and forwards it to the selected server.
- To provide the load balancing algorithm with the required input data, the load balancer also retrieves information about the servers' health and load to verify that they can respond to traffic.



1. 가 . )
2. 가 ? = . 가 , ... ( )  
- > 가 .

# Availability vs. Scalability

- **Availability** is defined by *uptime* -- the time between failures. (Downtime is the time to detect the failure, repair it, perform required recovery, and restart tasks.)
  - During uptime the system must respond to each request within a predetermined, well-defined time. If this time is exceeded, the client sees this as a server malfunction.
  - High availability, basically, is redundancy in the system: if one server fails, the others take over the failed server's load transparently. The failure of an individual server is invisible to the client.

가 . 가 , ...? . = ( < - > , )

가            가            .            .            가,,            가            가            가

가            가,  
가            가            가            가            가            가

.            ...

- **Scalability** means that the system can serve a single client, as well as thousands of simultaneous clients, by meeting quality-of-service requirements such as response time.
  - Under an increased load, a high scalable system can increase the throughput almost linearly in proportion to the power of added hardware resources.
- High scalability is reached by distributing the incoming request over the servers. If the load increases, additional servers can be added, as long as the load balancer does not become the bottleneck. To reach high availability, the load balancer must monitor the servers to avoid forwarding requests to overloaded or dead servers. Furthermore, the load balancer itself must be redundant too.

# Dispatcher(Load Balancer) Architecture Style - Solution

- Implementations

- In general, server load balancing solutions are of two main types:
  - *Transport-level* load balancing -- such as the DNS-based approach or TCP/IP-level load balancing -- acts independently of the application payload.
  - *Application-level* load balancing uses the application payload to make load balancing decisions.
- Load balancing solutions can be further classified into software-based load balancers and hardware-based load balancers.
  - Hardware-based load balancers are specialized hardware boxes that include application-specific integrated circuits (ASICs) customized for a particular use. Hardware-based load balancers are often used for transport-level load balancing. In general, hardware-based load balancers are faster than software-based solutions. Their drawback is their cost.
  - In contrast to hardware load balancers, software-based load balancers run on standard operating systems and standard hardware components such as PCs. Software-based solutions runs either within a dedicated load balancer hardware node or directly in the application.



# Dispatcher(Load Balancer) Architecture Style -Examples

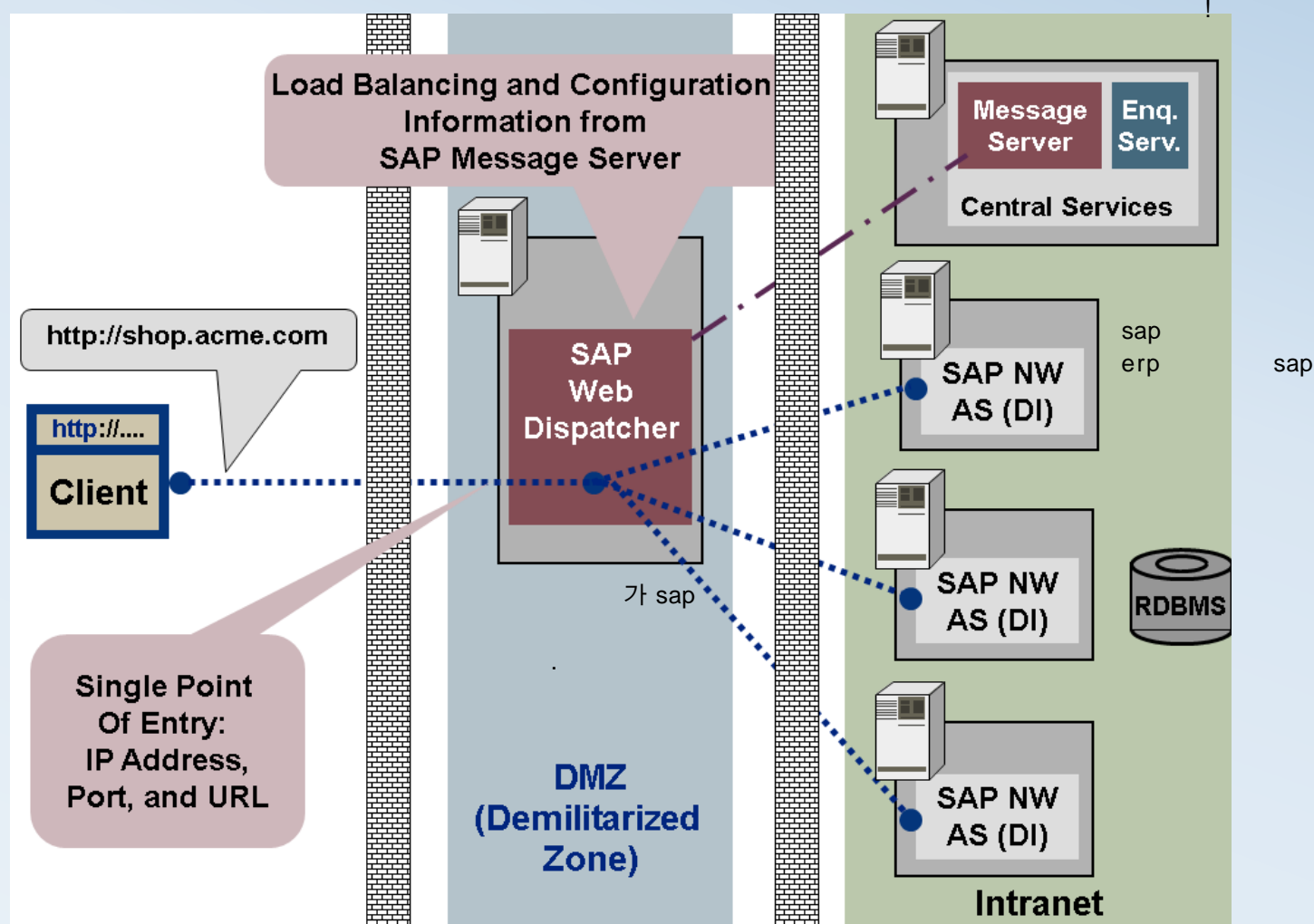
- **DNS-based load balancing**

- The DNS-based approach is based on the fact that DNS allows multiple IP addresses (real servers) to be assigned to one host name.
- DNS is an efficient solution for global server load balancing, where load must be distributed between data centers at different locations.
- Often the DNS-based global server load balancing is combined with other server load balancing solutions to distribute the load within a dedicated data center.

- **TCP/IP server load balancing**

- TCP/IP server load balancers operate on low-level layer switching. A popular software-based low-level server load balancer is the [Linux Virtual Server \(LVS\)](#).
- The incoming requests on a TCP connection are forwarded to the real servers by the load balancer, which runs a Linux kernel patched to include IP Virtual Server (IPVS) code.

# Dispatcher(Load Balancer) Architecture Style -Examples



# Dispatcher(Load Balancer) Architecture Style - Solution

stateless vs stateful

가

가

- Implementations

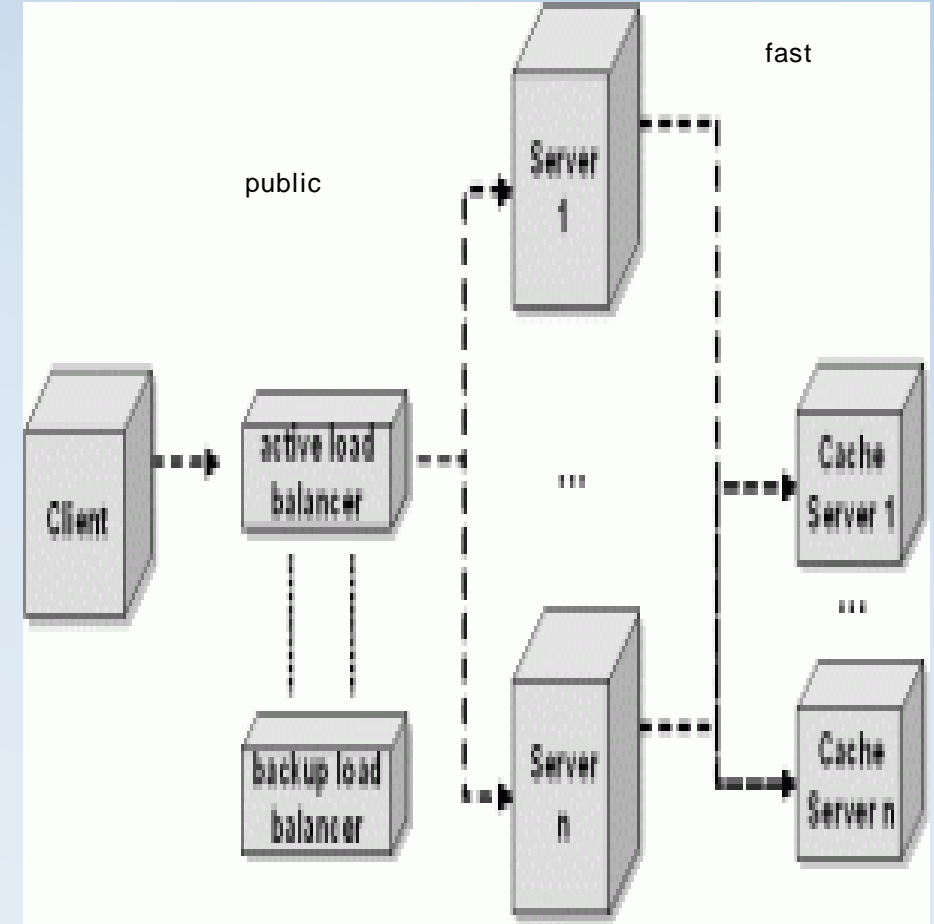
- Caching

- Low-level server load balancer solutions such as LVS reach their limit if application-level caching or application-session support is required. **Caching** is an important scalability principle for avoiding expensive operations that fetch the same data repeatedly. A cache is a temporary store that holds redundant data resulting from a previous data-fetch operation.
    - If a cache is used on the server side, straying requests will become a problem. One approach to handle this is to place the cache in a global space. [memcached](#) is a popular distributed cache solution that provides a large cache across multiple machines.

1. active - active

2. active - passive

가



가

가

가

가

3.

가