# Network Applications: Load Balancing among Homogeneous Servers

Qiao Xiang, Congming Gao, Qiang Su

https://sngroup.org.cn/courses/cnnsxmuf25/index.shtml

10/16/2025

# Outline

- Admin and recap
- □ Multiple network servers

## Recap: Operational Laws

- □ Utilization law: U = XS
- □ Forced flow law: Xi = Vi X
- □ Bottleneck device: largest Di = Vi Si
- Little's Law: Qi = Xi Ri
- Bottleneck bound of interactive response (for the given closed model):

$$X(N) \le \min\{\frac{1}{D_{\max}}, \frac{N}{D+Z}\}$$

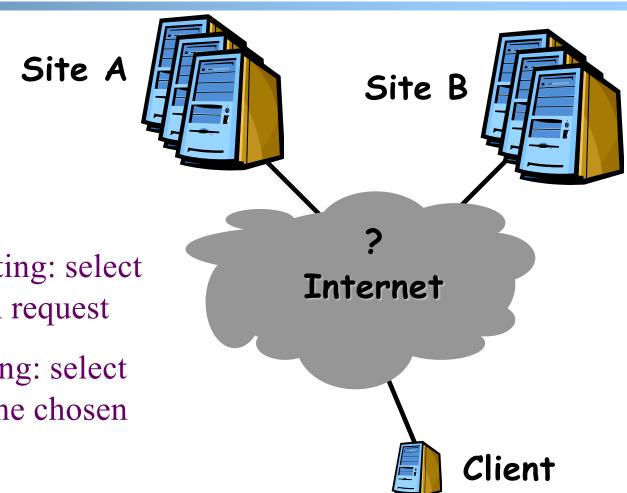
$$R(N) \ge \max\{D, ND_{\max} - Z\}$$

# Recap: Why Multiple Servers?

#### Scalability

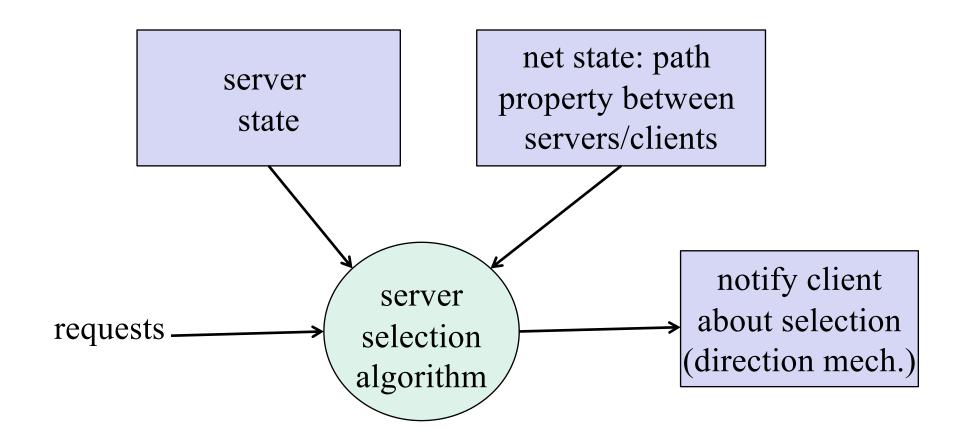
- Scaling beyond single server throughput
  - · There is a fundamental limit on what a single server can
    - process (CPU/bw/disk throughput)
    - store (disk/memory)
- Scaling beyond single geo location latency
  - There is a limit on the speed of light
  - Network detour and delay further increase the delay

## Request Routing: Overview



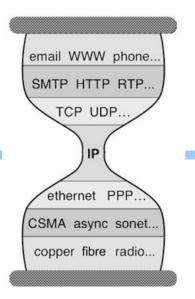
- Global request routing: select a server site for each request
- Local request routing: select a specific server at the chosen site

#### Request Routing: Basic Architecture

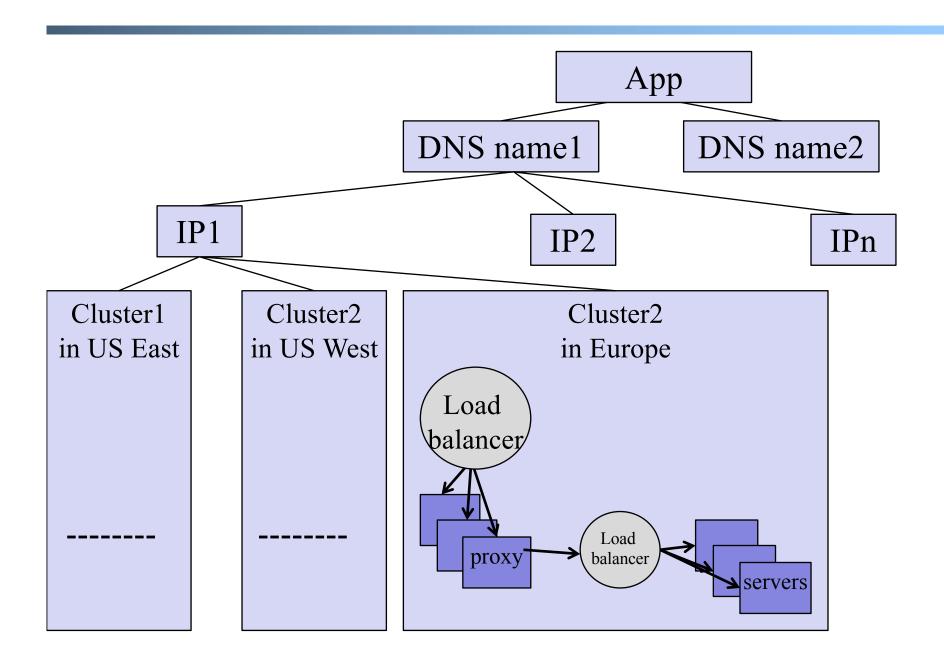


#### Client Direction Mechanisms

- Key difficulty
  - May need to handle a large of clients
- Basic types of mechanisms
  - Application layer, e.g.,
    - App/user is given a list of candidate server names
    - HTTP redirector
  - DNS: name resolution gives a list of server addresses
  - IP layer: Same IP address represents multiple physical servers
    - IP anycast: Same IP address shared by multiple servers and announced at different parts of the Internet.
       Network directs different clients to different servers
    - Smart-switch indirection: a server IP address may be a virtual IP address for a cluster of physical servers



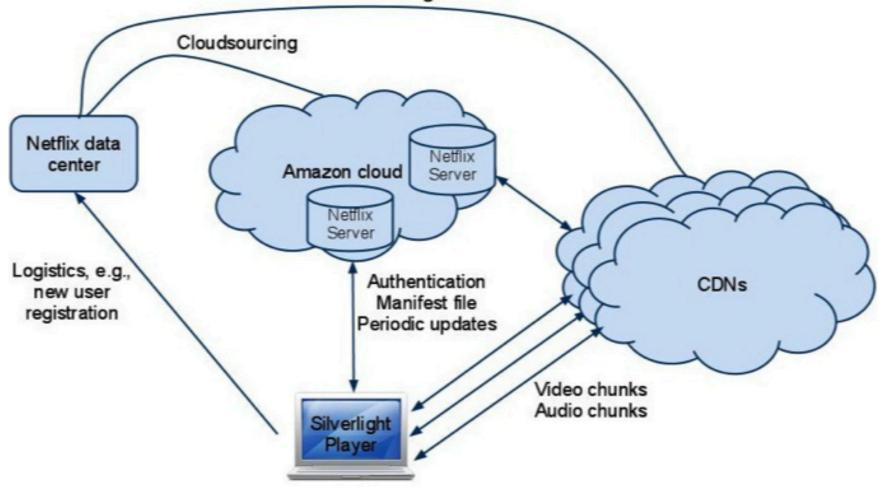
#### Direction Mechanisms are Often Combined



#### Example: Netflix

Hostname	Organization
www.netflix.com	Netflix
signup.netflix.com	Amazon
movies.netflix.com	Amazon
agmoviecontrol.netflix.com	Amazon
nflx.i.87f50a04.x.lcdn.nflximg.com	Level 3
netflix-753.vo.llnwd.net	Limelight
netflix753.as.nflximg.com.edgesuite.net	Akamai

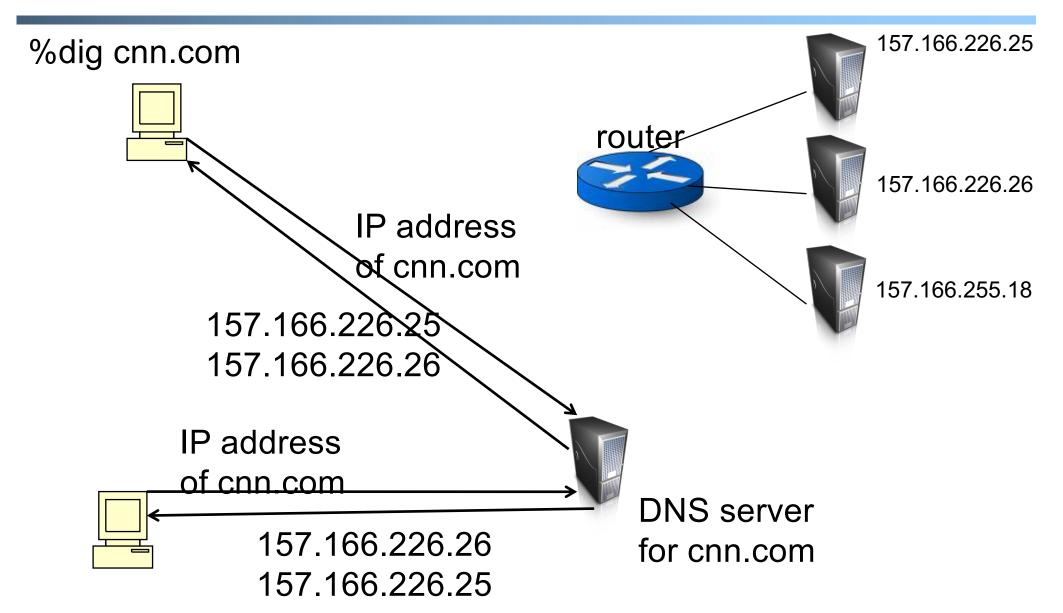
#### Outsourcing



### Outline

- Recap
- □ Single network server
- □ Multiple network servers
  - Why multiple servers
  - Request routing mechanisms
    - Overview
    - Application-layer
    - > DNS

#### Basic DNS Indirection and Rotation



# CDN Using DNS (Akamai Architecture as an Example)

- Content publisher (e.g., cnn)
  - provides base HTML documents
  - runs origin server(s); but delegates heavy-weight content (e.g., images) to CDN
- □ Akamai runs
  - (~240,000) edge servers for hosting content
    - · Deployment into 130 countries and 1600 networks
    - Claims 85% Internet users are within a single "network hop" of an Akamai CDN server.
  - customized DNS redirection servers to select edge servers based on
    - · closeness to client browser, server load

### Linking to Akamai

Originally, URL Akamaization of embedded content: e.g.,

<IMG SRC= http://www.provider.com/image.gif >
 changed to

<IMG SRC = http://a661. g.akamai.net/hash/image.gif>

Note that this DNS redirection unit is per customer, not individual files.

URL Akamaization is becoming obsolete and supported mostly for legacy reasons

### Exercise

https://cdn.nba.com/manage/2021/08/G ettylmages-1234610091-scalede1665274852371-784x441.jpg

- Check any web page of nba.com and find a page with an image
- ☐ Find the URL
- □ Use

%dig [+trace]

to see DNS load direction

https://www.nba.com/news/reports-lakers-extend-gm-rob-pelinka-through-2025-26-season

#### Exercise

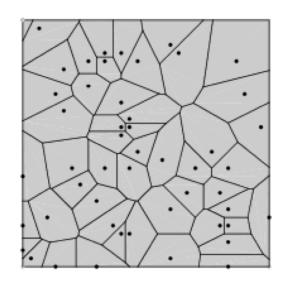
□ https://cdn.nba.com/manage/2025/10/16x 9-43.jpg?w=1568&h=882

#### Recap: CNAME based DNS Name

#### Typical design

- Use cname to create aliases, e.g., https://cdn.nba.com/manage/2021/08/GettyImages-1234610091-scaled-e1665274852371-1568x882.jpg
- o cname: e8017.dscb.akamaiedge.net
  - why two levels in the name?

#### Two-Level Direction



 high-level DNS determines proximity, directs to low-level DNS;

Input: dscb.akamaiedge.net & and client IP,

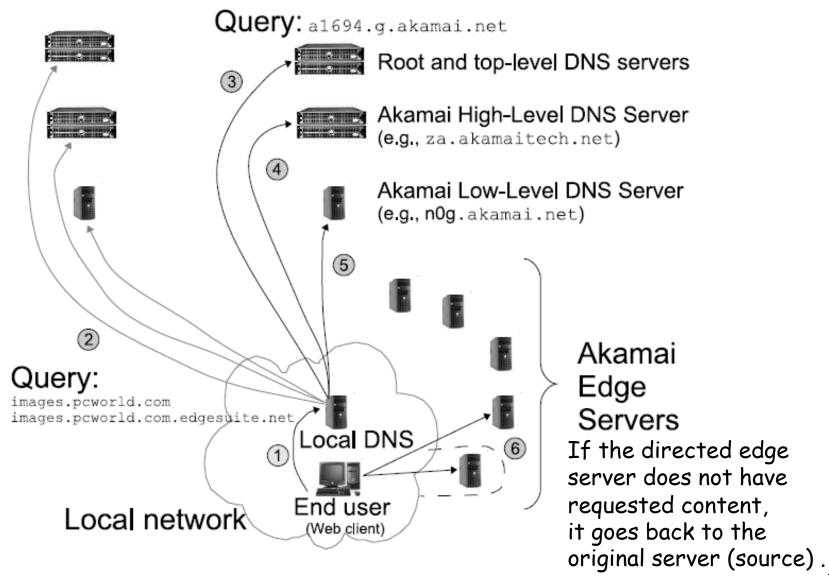
Output: region (low-level) DNS

 low-level DNS: who manages a close-by cluster of servers with different IP addresses

Input: e8017.dscb.akamaiedge.net & and client IP

Output: specific servers

#### Akamai Load Direction



# Two-Level DNS Mapping Alg

- High-level
  - Typically geo-location matching from client to reg
- Low-level
  - □ Typically secret, e.g., details of Akamai algorithms are proprietary
  - Typical goal: load balancing among servers

# Akamai Local DNS LB Alg

- A Bin-Packing algorithm (column 12 of Akamai Patent) every T second
  - Compute the load to each publisher k (called serial number)
  - Sort the publishers from increasing load
  - For each publisher, associate a list of random servers generated by a hash function
    - Hash range may be increasing, e.g., first hash [0,1], second [0, 3]
  - Assign the publisher to the first server that does not overload

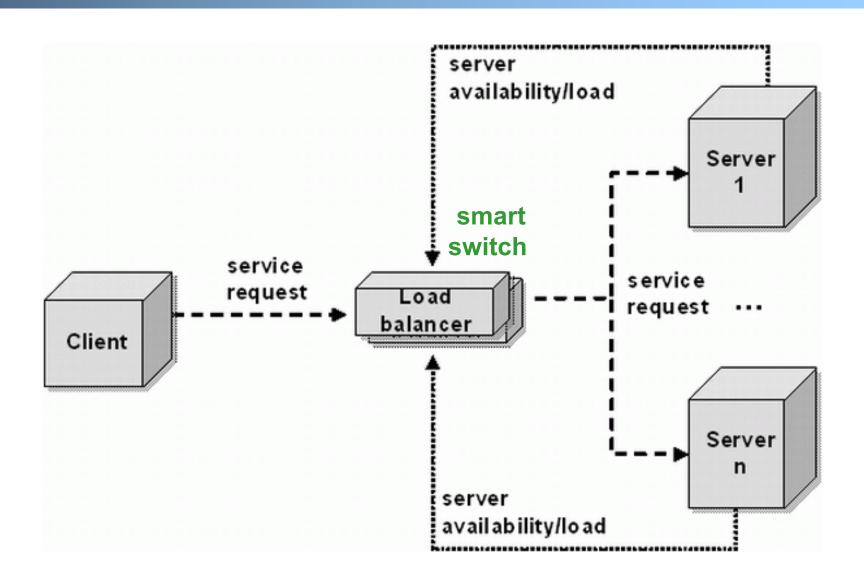
#### Discussion

- Advantages of using DNS for using multiple servers (LB)
  - Leveraging existing DNS features (e.g., cname, hierarchy name for natural hierarchical redirection)
  - Leveraging existing DNS deployment/optimization
- Disadvantages of using DNS
  - Distributed caching may lead to slow response
  - Only in the unit of IP addresses

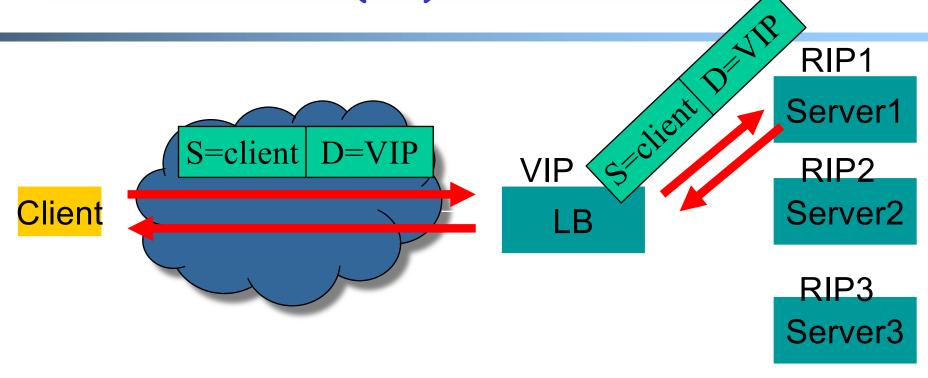
#### Outline

- Recap
- □ Multiple network servers
  - Basic issues
  - Load direction
    - · DNS (IP level)
    - > Load balancer/smart switch (sub-IP level)

# Smart Switch: Big Picture



Load Balancer (LB): Basic Structure



Problem of the basic structure?



#### Problem

- Client to server packet has VIP as destination address, but real servers use RIPs
  - if LB just forwards the packet from client to a real server, the real server drops the packet
  - reply from real server to client has real server IP as source -> client will drop the packet

state: listening address: {RealIP.6789, \*:\*} completed connection queue: C1; C2 sendbuf: recvbuf:

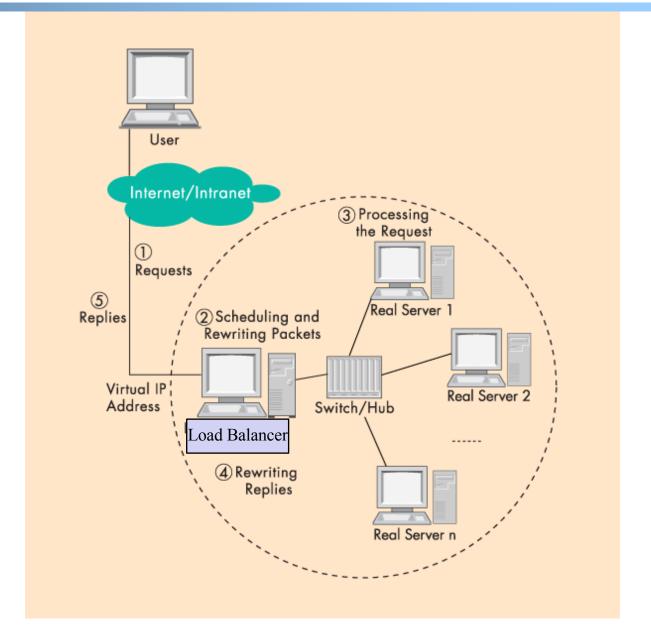
state: address: {RealIP:6789, 198.69.10.10.1500} sendbuf: recvbuf:

state: address: {VIP: <b>6789</b> , 198.69.10.10. <b>1500</b> } sendbuf: recvbuf:	

server client

# Solution 1: Network Address Translation (NAT)

- □ LB does rewriting/ translation
- □ Thus, the LB is similar to a typical NAT gateway with an additional scheduling function



# LB/NAT Advantages and Disadvantages

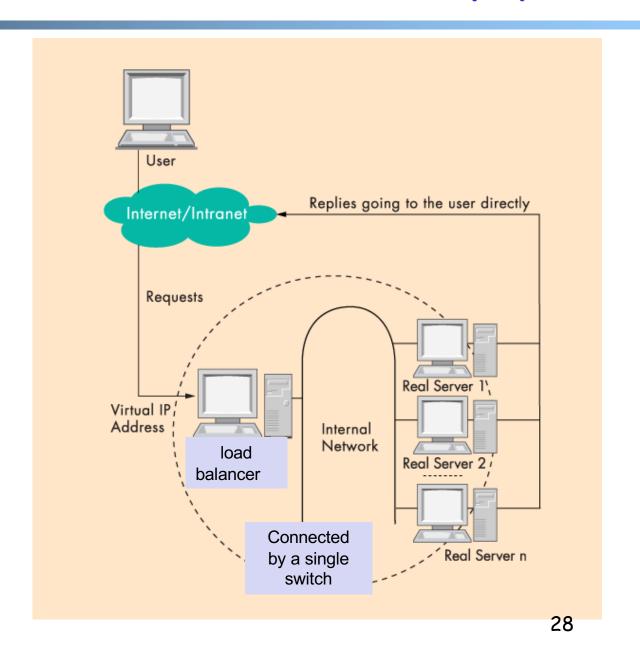
#### Advantages:

- Only one public IP address is needed for the load balancer; real servers can use private IP addresses
- Real servers need no change and are not aware of load balancing

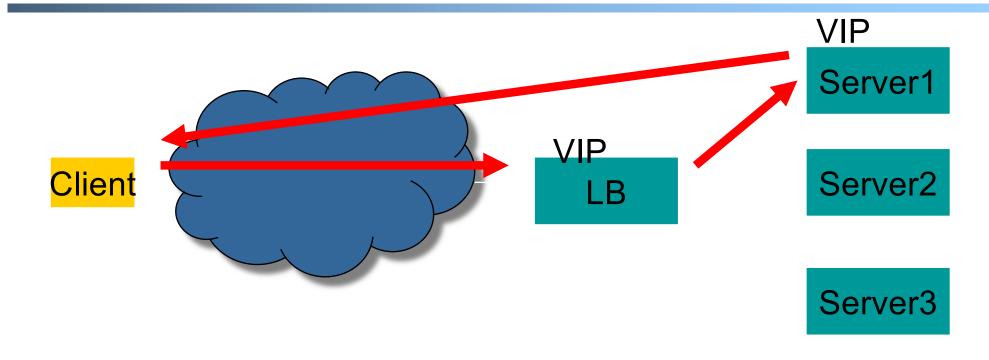
#### Problem

- The load balancer must be on the critical path and hence may become the bottleneck due to load to rewrite request and response packets
  - Typically, rewriting responses has more load because there are more response packets

# Goal: LB w/ Direct Reply



# LB with Direct Reply: Implication

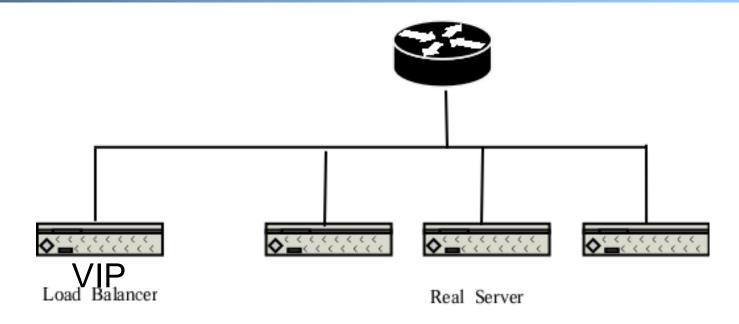


Direct reply Each real server uses VIP as its IP address

Address conflict: multiple

Address conflict: multiple devices w/ the same IP addr

# Why IP Address Matters?

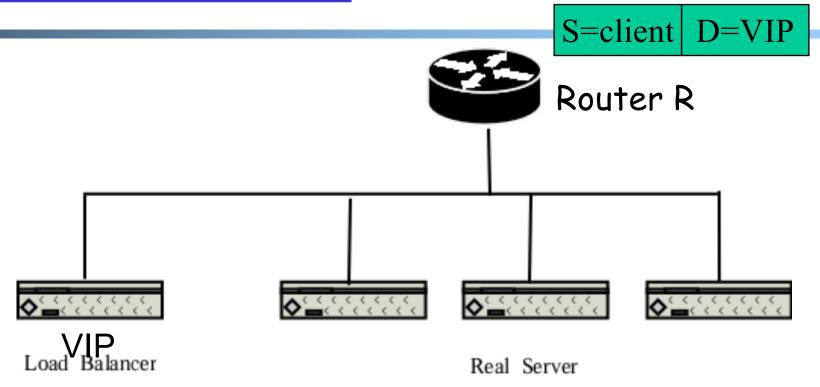


- Each network interface card listens to an assigned MAC address
- □ A router is configured with the range of IP addresses connected to each interface (NIC)
- □ To send to a device with a given IP, the router needs to translate IP to MAC (device) address
- The translation is done by the Address Resolution Protocol (ARP)

#### ARP Protocol

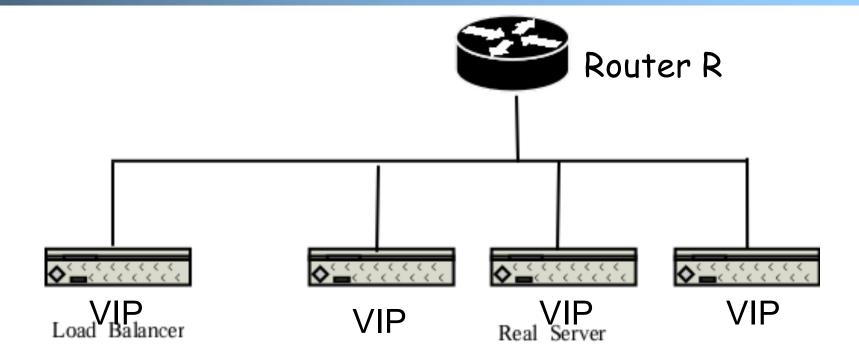
- ARP is "plug-and-play":
  - nodes create their ARP tables without intervention from net administrator
- □ A broadcast protocol:
  - Router broadcasts query frame, containing queried IP address
    - · all machines on LAN receive ARP query
  - Node with queried IP receives ARP frame, replies its MAC address

#### ARP in Action



- Router broadcasts ARP broadcast query: who has VIP?
- ARP reply from LB: I have VIP; my MAC is  $MAC_{LB}$
- Data packet from R to LB: destination  $MAC = MAC_{LB}$

#### LB/DR Problem



#### ARP and race condition:

- When router R gets a packet with dest. address VIP, it broadcasts an Address Resolution Protocol (ARP) request: who has VIP?
- · One of the real servers may reply before load balancer

Solution: configure real servers to not respond to ARP request

# LB via Direct Routing

- □ The virtual IP address is shared by real servers and the load balancer.
- Each real server has a non-ARPing, loopback alias interface configured with the virtual IP address, and the load balancer has an interface configured with the virtual IP address to accept incoming packets.
- □ The workflow of LB/DR is similar to that of LB/NAT:
  - the load balancer directly routes a packet to the selected server
    - the load balancer simply changes the MAC address of the data frame to that
      of the server and retransmits it on the LAN (how to know the real server's
      MAC?)
  - When the server receives the forwarded packet, the server determines that the packet is for the address on its loopback alias interface, processes the request, and finally returns the result directly to the user

# LB/DR Advantages and Disadvantages

#### Advantages:

 Real servers send response packets to clients directly, avoiding LB as bottleneck

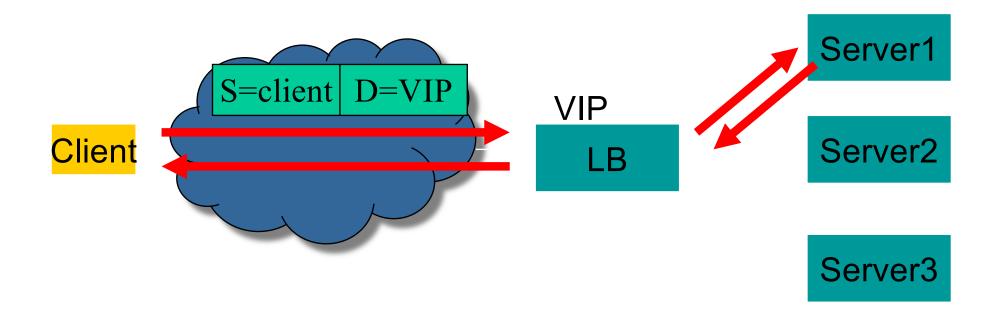
#### ■ Disadvantages:

- Servers must have non-arp alias interface
- The load balancer and server must have one of their interfaces in the same LAN segment
- Considered by some as a hack, not a clean architecture

#### Example Implementation of LB

- □ An example open source implementation is Linux virtual server (linux-vs.org)
  - Used by
    - www.linux.com
    - sourceforge.net
    - wikipedia.org
  - More details on ARP problem: http://www.austintek.com/LVS/LVS-HOWTO/HOWTO/LVS-HOWTO.arp\_problem.html
  - Many commercial LB servers from F5, Cisco, ...
- More details please read chapter 2 of Load Balancing Servers, Firewalls, and Caches

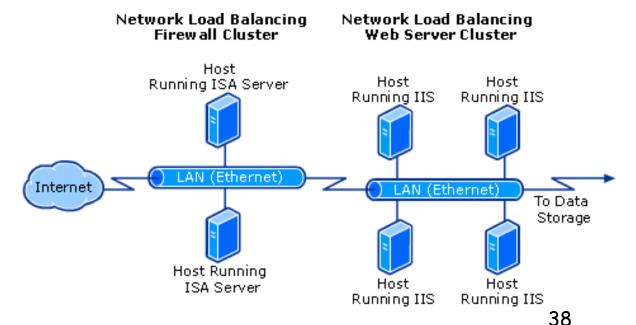
# Problem of the Load Balancer Architecture



One major problem is that the LB becomes a single point of failure (SPOF).

#### Solutions

- Redundant load balancers
  - E.g., two load balancers (a good question to think offline)
- Fully distributed load balancing
  - e.g., Microsoft Network Load Balancing (NLB)



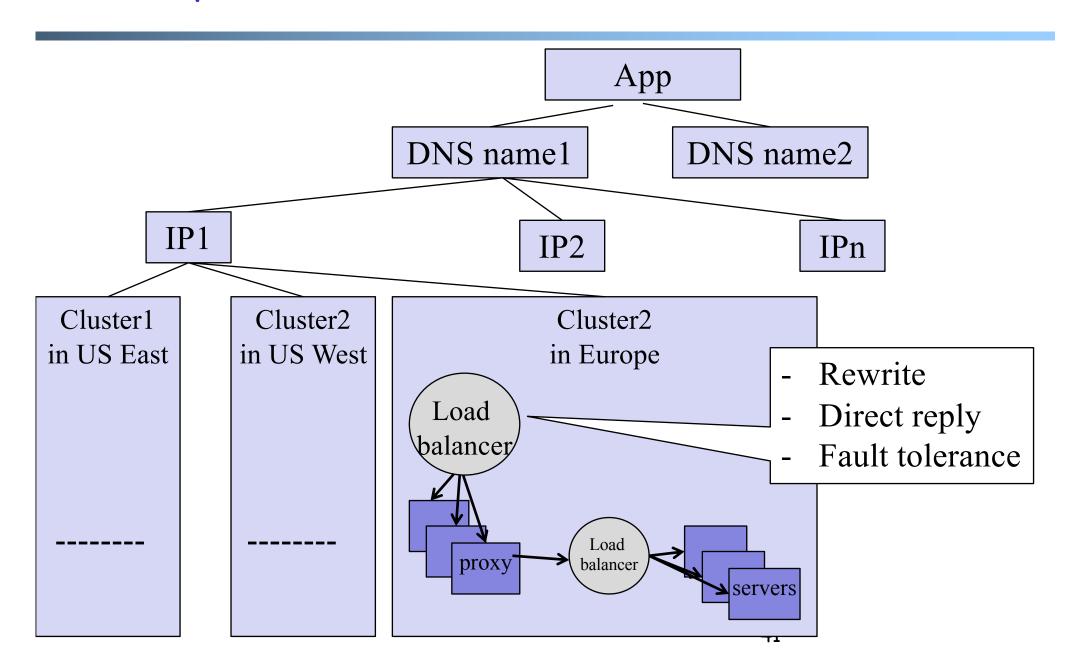
### Microsoft NLB

- □ No dedicated load balancer
- All servers in the cluster receive all packets
- Key issue: one and only one server processes each packet
  - All servers within the cluster simultaneously run a mapping algorithm to determine which server should handle the packet. Those servers not required to service the packet simply discard it.
    - Mapping (ranking) algorithm: computing the "winning" server according to host priorities, multicast or unicast mode, port rules, affinity, load percentage distribution, client IP address, client port number, other internal load information

#### Discussion

□ Compare the design of using Load Balancer vs Microsoft NLB

#### Recap: Direction Mechanisms



#### Scalability of Server-Only Approaches

