

Scalable Internet Services and Load Balancing

Kai Shen

11/9/2011

CSC 257/457 - Fall 2011

1

Internet Services

- Internet brings ubiquitous connection
 - Internet-based applications/services accessible to online users through Internet.
 - Examples?
- Trends:
 - Centralization of applications/services
 - Scalability requirements: many simultaneous user accesses; large amount of hosted data, ...

11/9/2011

CSC 257/457 - Fall 2011

2

Search Engine as An Example: Step 1 - Crawling

- Crawling - get all these Web pages out there:
 - First retrieve some root pages;
 - Parse their content and follow hyperlinks to retrieve more pages;
 - Depth-first search or breadth-first search?

11/9/2011

CSC 257/457 - Fall 2011

3

Performance Analysis for Crawling

- What are the resources involved?
 - CPU processing for TCP/HTTP protocol handling and the parsing of page content
 - writing to disk storage
 - network bandwidth to remote web sites
- Assume average page size 10KB
 - raw processing power of a single CPU
 - 1000 pages/sec
 - I/O to a single disk
 - 100 seeks/sec \Rightarrow up to 100 pages/sec
 - network bandwidth from/to the Internet
 - T1 link (1.5Mbit/s) \Rightarrow 12 pages/sec
 - T3 link (45Mbit/s) \Rightarrow 360 pages/sec

11/9/2011

CSC 257/457 - Fall 2011

4

Parallel Crawling

- Challenge
 - Avoid redundant crawling

11/9/2011 CSC 257/457 - Fall 2011 5

Search Engine as An Example: Step 2 - Indexing

- Indexing
 - crawled raw web pages are not easy to search.
 - we index them to formats that are easy to search.
- As part of indexing, we need to give each page an ID
 - using a hash function.

Computer:

Page #123	Page #357
-----------	-----------	--------

Networks:

Page #124	Page #468
-----------	-----------	--------

- Fast intersection?

11/9/2011 CSC 257/457 - Fall 2011 6

Search Engine as An Example: Step 3 - Online Search

Scalability, reliability

11/9/2011 CSC 257/457 - Fall 2011 7

Partitioning and Replication

11/9/2011 CSC 257/457 - Fall 2011 8

Load Balancing over Internet Servers

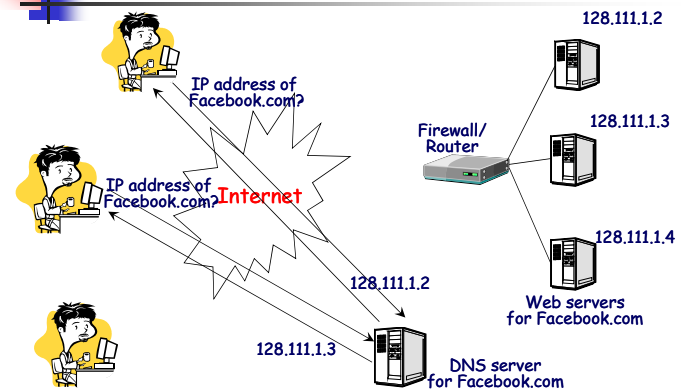
- Popular sites like Google or Facebook receive tens or hundreds of millions of hits per day.
- A large number of replicated servers are used at these sites.
- **Key question:** how to balance client requests over these servers?
- Goals:
 - Performance
 - Ease of deployment

11/9/2011

CSC 257/457 - Fall 2011

9

Load Balancing on Internet Servers Technique 1 - DNS Rotation



11/9/2011

CSC 257/457 - Fall 2011

10

Discussions on DNS Rotation

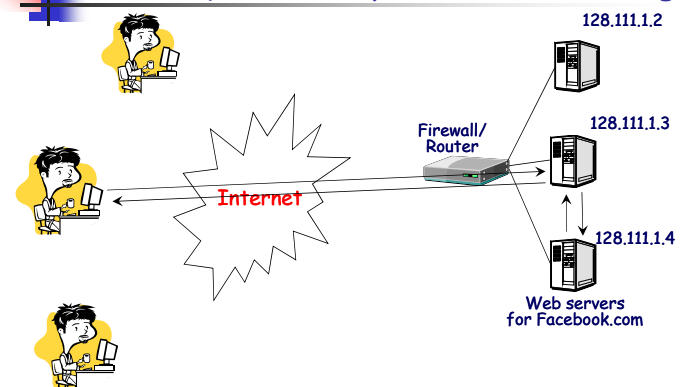
- Advantages
 - Require almost no change on the existing Internet architecture
- Problems
 - DNS Caching
 - Rigid load balancing policy
 - can't balance based on runtime load changes
 - slow or no adjustment in response to failures

11/9/2011

CSC 257/457 - Fall 2011

11

Load Balancing on Internet Servers Technique 2 - Cooperative Offloading



11/9/2011

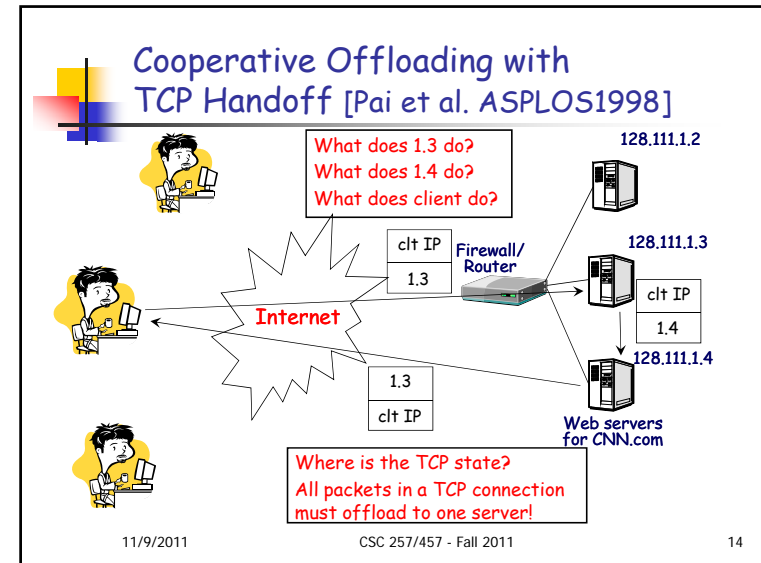
CSC 257/457 - Fall 2011

12

Discussions on Cooperative Offloading

- Can be combined with the DNS rotation.
- Advantages:
 - More flexible policy is possible
 - Be more responsive to runtime workload and server failures (to a certain degree)
- Problems:
 - Need software changes on servers
 - Longer delay

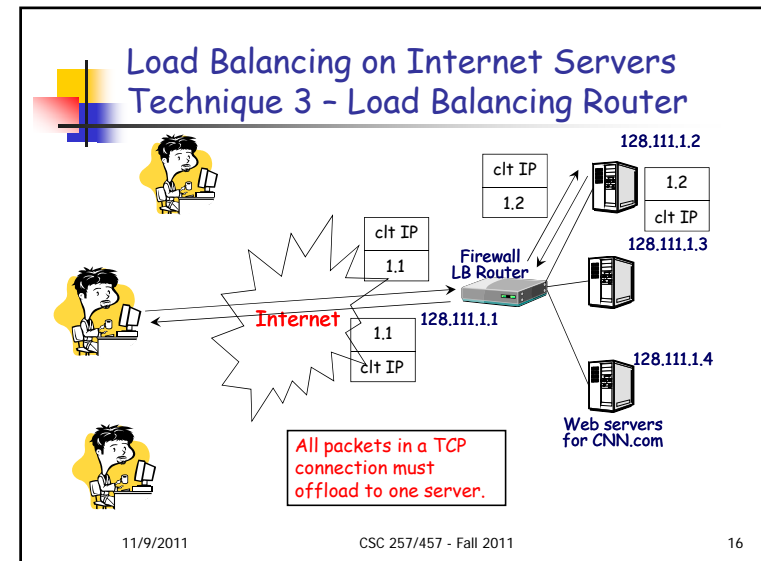
11/9/2011 CSC 257/457 - Fall 2011 13



Cooperative Offloading vs. TCP Handoff

- Software changes on the servers
- Delays

11/9/2011 CSC 257/457 - Fall 2011 15





More About Load Balancing Router

Load balancing policies in LB routers (Goal: transparency, plug-and-play)

- Simple rotation
- Least number of active requests
- Shortest response time

How deep do we look into the network protocol stack?

- Network layer (IP)?
- Transport layer (TCP/UDP)?
- Application layer?

11/9/2011

CSC 257/457 - Fall 2011

17



Summary

- Scalable Internet servers
 - partitioning
 - replication
- Load balancing for Internet servers
 - DNS rotation
 - Cooperative offloading (w. TCP handoff)
 - Load balancing router
- Changes required on the components:
 - DNS server??
 - Web server??
 - Client??
 - Router??

11/9/2011

CSC 257/457 - Fall 2011

18