

LARD

(Locality Aware Request Distribution)

By
Saumitra M Das

Outline

- Motivation
- Introduction & Background
- Basic and Detailed Idea
- Issues and Optimizations
- Performance Evaluation

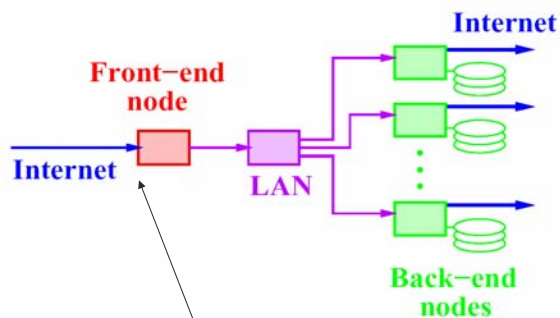
Motivation

JUSTIFYING DOING ALL THIS WORK

- Demand is growing on network servers
 - Why?, you ask..
 - Internet is growing
 - BW is increasing
 - Site consolidation
- Cheaper \$olution\$ are attractive
 - Inexpensive PCs and LANs
 - Parallelizable nature of job

Motivation

FOCUS OF THE WORK



How should this box behave so that we save money.

In other words, what policy should it use for request distribution so that throughput is maximized with minimum utilization of resources.

Motivation

DISSECTING THE TITLE

- Locality
 - Taking advantage of files already cached in back-end server's memory
 - What does this mean?
 - Requests accessing same data be sent to the same set of servers
- So, locality aware is...
- Is that all we need?
 - If so, you could go home instead of listening to me.

Motivation

DISSECTING THE TITLE

- It is only half the picture
- What about distribution ?
 - Load balanced system
 - Distribute requests evenly among back-end servers
- What do we want in the ideal case?
 - Improve hit-rate and response time
- Many studies focus only on one aspect and ignore the other

Introduction

REQUEST DISTRIBUTION PROTOCOL V 1.0

Round Robin Request Distribution

Incoming Requests

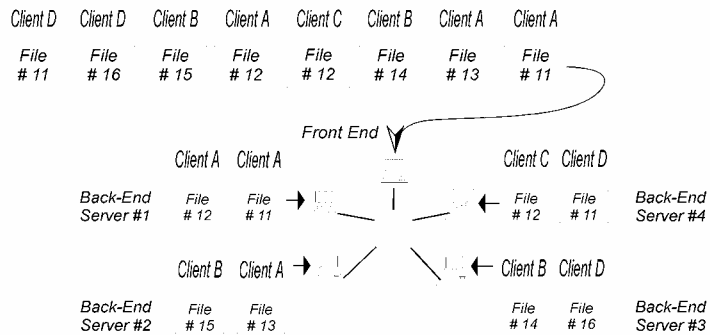


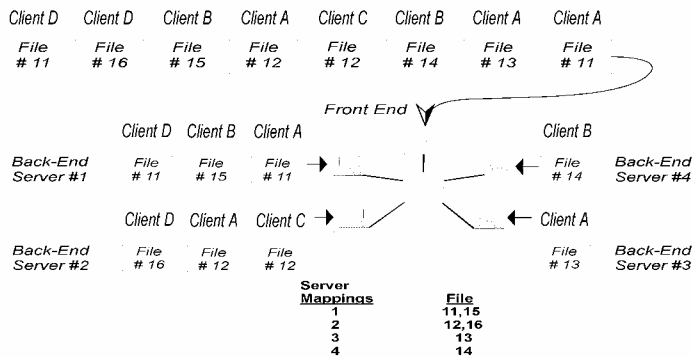
Diagram Source: A Khaleel, A.Reddy, TAMU

Introduction

REQUEST DISTRIBUTION PROTOCOL V 2.0

File Based Request Distribution

Incoming Requests

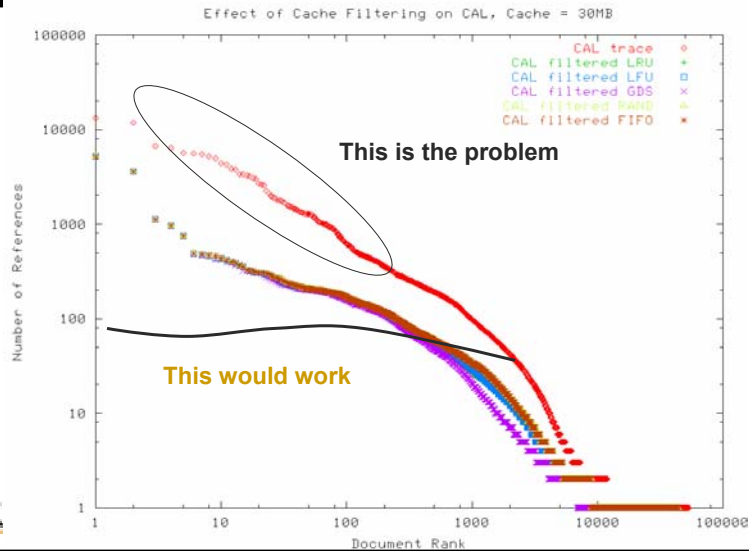


Will this work for a web workload ?

Diagram Source: A Khaleel, A.Reddy, TAMU

Introduction

IS REQUEST DISTRIBUTION PROTOCOL V 2.0 ALL WE NEED?



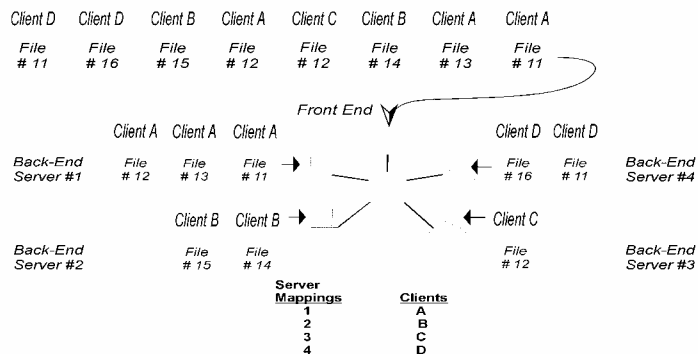
9

Introduction

REQUEST DISTRIBUTION PROTOCOL V 3.0

Client Based Request Distribution

Incoming Requests



10

Introduction

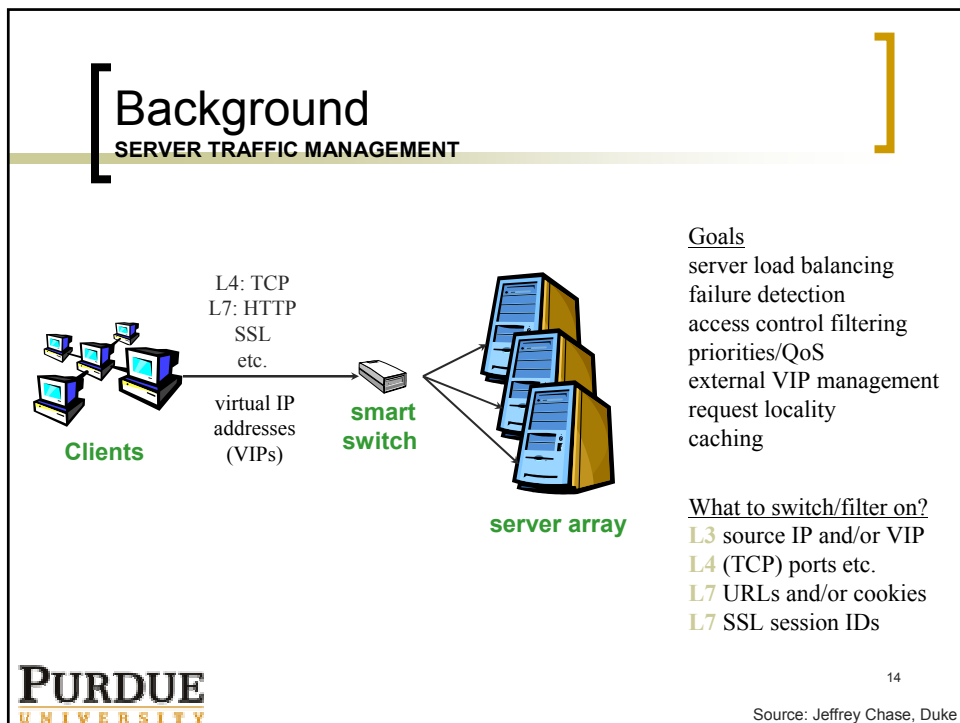
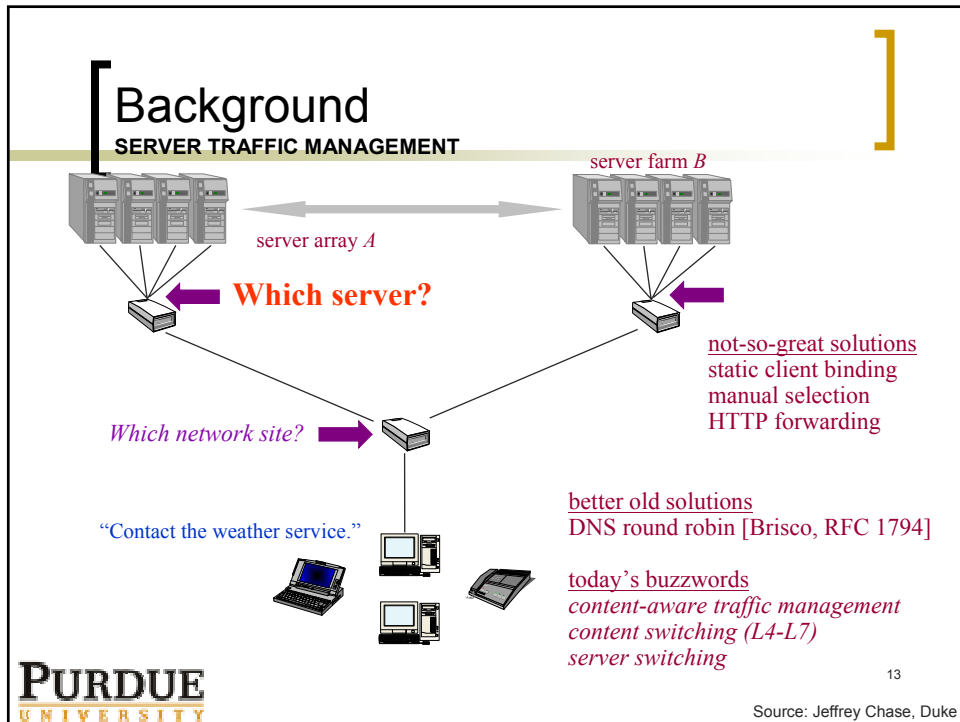
REQUEST DISTRIBUTION PROTOCOL V 4.0

- Locality Aware Request Distribution
 - Broadly based on file-based scheme
 - Addresses the issue of load balancing
 - **Each file assigned a dynamic set of servers instead of just one server**
 - **Each server assigned on average a static set of files instead of a dynamic set**

Introduction

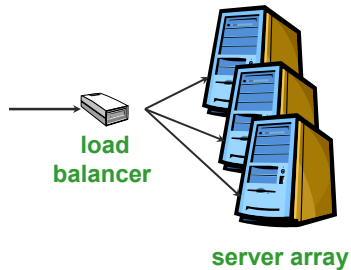
CONTRIBUTION OF THE PAPER

- Scheme designed to achieve both locality and load balancing
- Simulation studies and comparisons
- Prototype Implementation and Evaluation



Background

WHAT CAN YOU DO WITH A L4 SWITCH



Issues

- switch redundancy
- mechanics of L4 switching
- handling return traffic
- server failure detection (health checks)

Policies

- random
- weighted round robin (WRR)
- lightest load
- least connections

Key point: the heavy lifting of server selection happens only on connect request (SYN).

Performance metric: connections per second.

Limitations

- connection-grained
- no request locality
- no session locality

Background

HOW TO SELECT THE SERVER? – THE OLD WAY

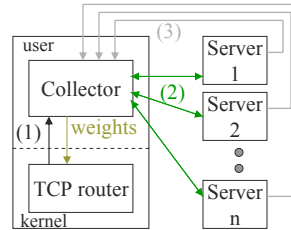
```
/* initially index = 0, L = (S0, S1, ..., Sn-1) */
While (1) {
    i == index;
    if (i == 0) {
        cw--;
        if (cw <= 0) cw = max of W(Sn);
        if (cw == 0) return NULL;
    }
    if (W(Si) >= cw) {
        index = (i + 1) % n;
        return Si;
    }
    else
        index = (i+1) % n;
}
```


Background

L4 SWITCH + WRR

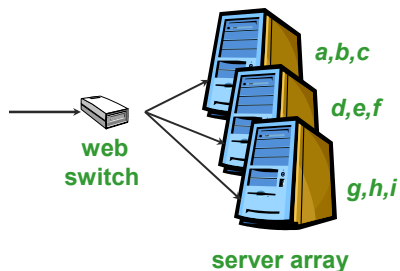
- (1) number of active connections
- (2) periodic polling (http request/response)
- (3) periodic script execution and feedback
 - queue length (cpu, disk)
 - amount of allocated network buffer

- Pros:
 - high routing throughput
- Cons:
 - Content-blind scheduling
 - Difficult to compute weights accurately
 - insufficient load information, load information feedback delay & overhead
 - Dynamic load imbalance across servers



Background

WHAT CAN YOU DO WITH A L7 SWITCH?



Idea: switch parses the HTTP request, retrieves the request URL, and uses the URL to guide server selection.

Example: Foundry

host name
URL prefix
URL suffix
Substring pattern
URL hashing

Advantages

separate static content from dynamic
reduce content duplication
improve server cache performance
cascade switches for more complex policies

Issues

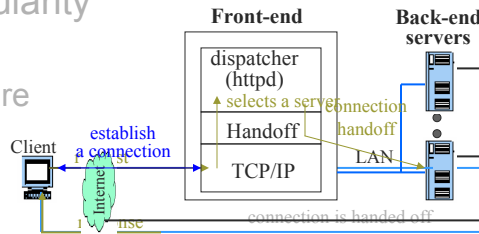
HTTP parsing cost
URL length
HTTP 1.1
session locality

Background

L7 SWITCH SCHEMATIC

- Scheduling granularity = connection

- but, content-aware scheduling

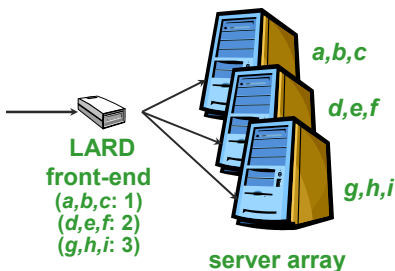


Benefits

- Performance improvement of back-end servers
 - cache affinity(locality) based scheduling
- Easy to specialize the back-ends for certain types of content or services

Basic Idea

THE LARD ALGORITHM



LARD front-end maintains an LRU cache of request targets and their locations, and table of active connections for each server.

Idea: route requests based on request URL, to maximize locality at back-end servers.

Policies

1. *LB* (locality-based) is URL hashing.
2. *LARD* is locality-aware & *LB*: route to target's site if there is one and it is not "overloaded", else pick a new site for the target.

Basic Idea

THE LARD ALGORITHM

```
while (true)
  fetch next request r;
  if server[r.target] = null then
    n, server[r.target] ← {least loaded node};
  else
    n ← server[r.target];
    if (n.load >  $T_{high}$  &&  $\exists$  node with load <  $T_{low}$ ) ||
       n.load  $\geq 2 * T_{high}$  then
      n, server[r.target] ← {least loaded node};
  send r to n;
```

Basic Idea

THE LARD ALGORITHM – WHAT IS LOAD?

- What is the load here
 - Load = # of active connections
- Why
 - Good summary metric
 - Robust
 - Available in front end
 - Simple
 - Quickly adapts

Basic Idea

THE LARD ALGORITHM - INTUITION

- T_{low} = load below which a back-end is likely to have idle resources
- T_{high} = load above which node is likely to cause substantial delay in serving a request
- How to choose these two?
 - The study chose 25 and 65

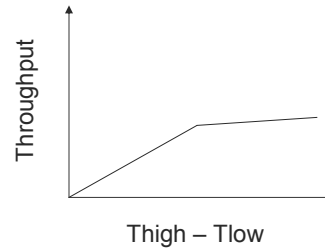
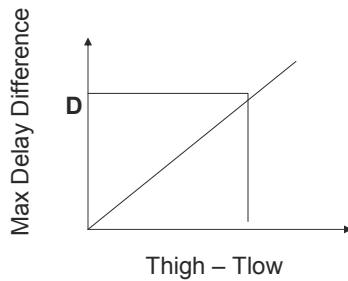
Basic Idea

THE LARD ALGORITHM - HOW TO CHOOSE THESE

- T_{low} = f (back-end speed)
 - Choose T_{low} high enough to avoid idling of back-ends
- T_{high} is a different issue
 - $T_{\text{high}} - T_{\text{low}}$ should be low enough to limit the delay variance among the back-ends
 - $T_{\text{high}} - T_{\text{low}}$ should be high enough to tolerate short term load fluctuations without affecting locality

Basic Idea

THE LARD ALGORITHM – HOW TO CHOOSE THESE

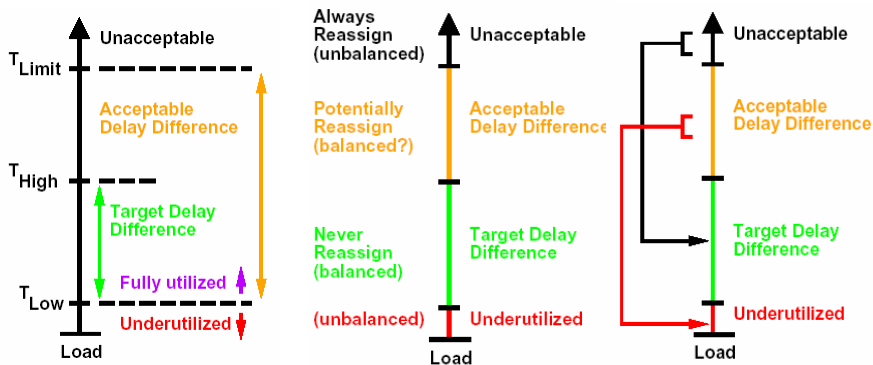


$$T_{high} = (T_{low} + D/R) / 2$$

R = average request service time

Basic Idea

THE LARD ALGORITHM – TO REASSIGN OR NOT TO REASSIGN
THAT IS THE QUESTION



Issues

HOW & WHAT TO IMPROVE?

■ WHAT

- If back-end serving target x fails ?
 - Replicate (in the paper)
- If load on all nodes rises to $2 \cdot T_{\text{high}}$?
 - Admission Control (in the paper)
- If the front-end is overwhelmed?
 - Distributed Distributor (follow up paper)
- If the workload is primarily dynamic content?
 - LARD for Dynamic Content (attempt by students)
- If acting as a proxy is too much overhead for the front-end?
 - TCP Handoff not Splicing (in the paper)
- If connections are persistent
 - P-HTTP support [ARON99] (follow up paper)

Optimizations

REPLICATION

■ Why

- to prevent a single target to cause a back-end from overloading

■ How

- Any ideas...?

Optimizations

REPLICATION

- The front end maintains a mapping from targets to a set of nodes that serve the target.
- Requests for a target are assigned to the least loaded node in the target's server set.
- If a load imbalance occurs, the front end checks if the server set has recently changed or not(within k seconds), if so, it picks a lightly loaded node and adds that to the set.
- If a request has multiple servers and has not moved or had a new one added, the front end needs to remove one to keep the degree not too high.

```

while (true)
  fetch next request r;
  if serverSet[r.target] = {} then
    n, serverSet[r.target] ← {least loaded node};
  else
    n ← {least loaded node in serverSet[r.target]};
    m ← {most loaded node in serverSet[r.target]};
    if (n.load > Thigh && ∃ node with load < Tlow) ||
       n.load ≥ 2Thigh then
      p ← {least loaded node};
      add p to serverSet[r.target];
      n ← p;
    if |serverSet[r.target]| > 1 &&
       time() - serverSet[r.target].lastMod > K then
      remove m from serverSet[r.target];
  send r to n
  if serverSet[r.target] changed in this iteration then
    serverSet[r.target].lastMod ← time();
  
```

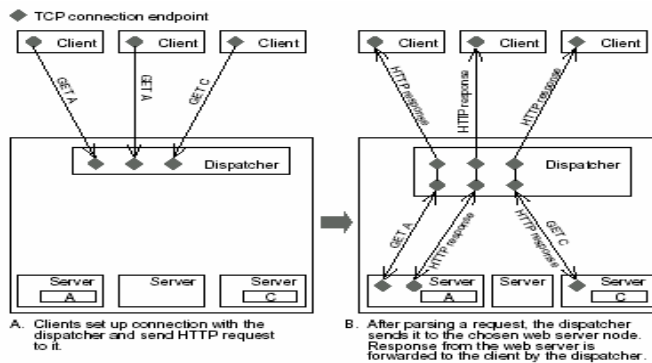
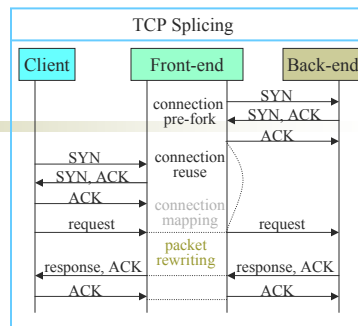
Optimizations

ADMISSION CONTROL

- Why
 - To prevent load on all nodes to become $> 2 \cdot T_{high}$ and LARD=WRR
- How
 - Front-end limits total number of connections to
 - $S = (n-1) \cdot T_{high} + T_{low} - 1$, (n = number of back-end nodes)
 - At most $(n-2)$ nodes can have a load $\geq T_{high}$ while no node has load $< T_{low}$
 - All n nodes can have load $> T_{low}$

Optimizations

TCP SPLICING

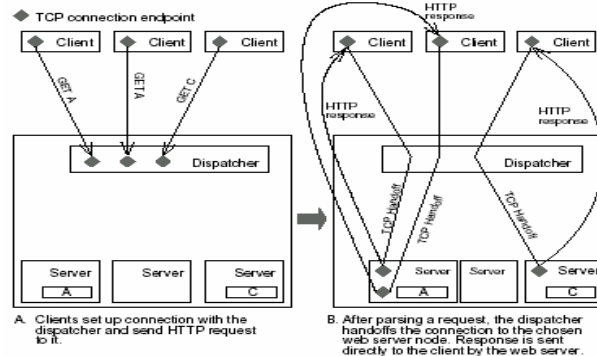
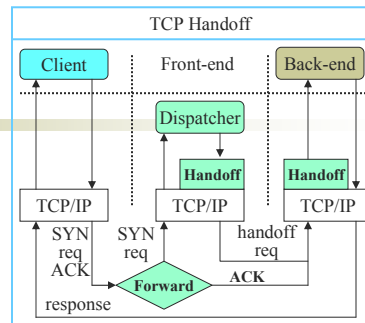


PURDUE
UNIVERSITY

31

Optimizations

TCP HANDOFF



PURDUE
UNIVERSITY

32

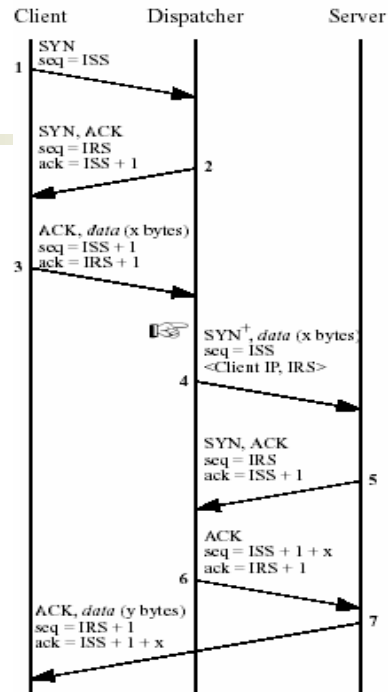
Optimizations

TCP HANDOFF – STATE TRANSFER

Key : individual server node can masquerade as the dispatcher as long as it has the appropriate connection state.

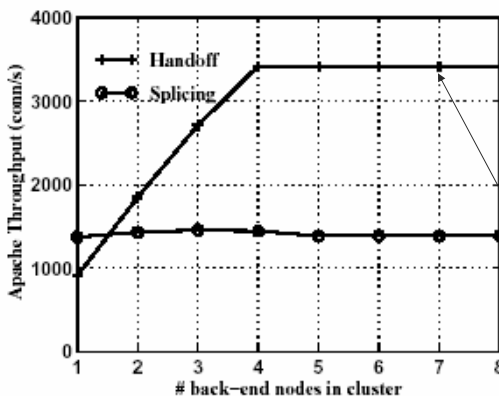
Caveats:

1. Dispatcher cannot accept any options such as SACK from the client which the backend does not support
2. Dispatcher must choose an initial flow control window that each server can satisfy



Optimizations

WHY IS HANDOFF BETTER



1. Splicing though better than a relay has the overhead due to HTTP response traffic
2. Handoff requires kernel modification to both frontend and backend but does not require HTTP response traffic to go through distributor

With handoff, the throughput scales with the cluster size more than splicing



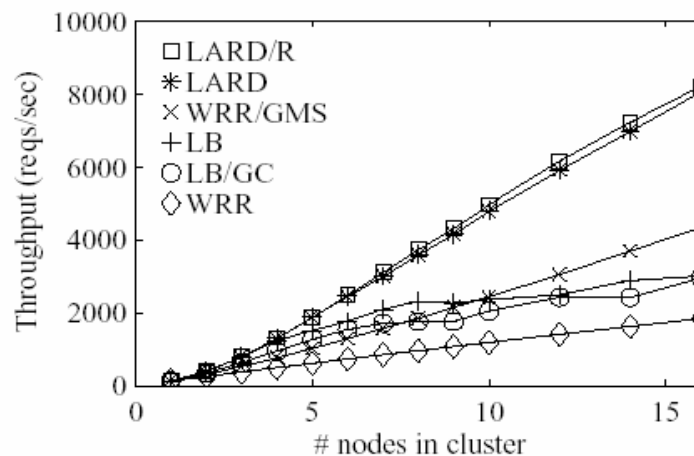
Simulation Results

SIMULATION MODEL

- CPU : Pentium II 300 Mhz
- Network : 100 Mbps LAN
- Disk : WD Enterprise
- Caching : GDS
- Filesystem : Berkeley FFS
- Input : Rice & IBM trace
- Strategies : WRR, LB-Locality, LARD /R

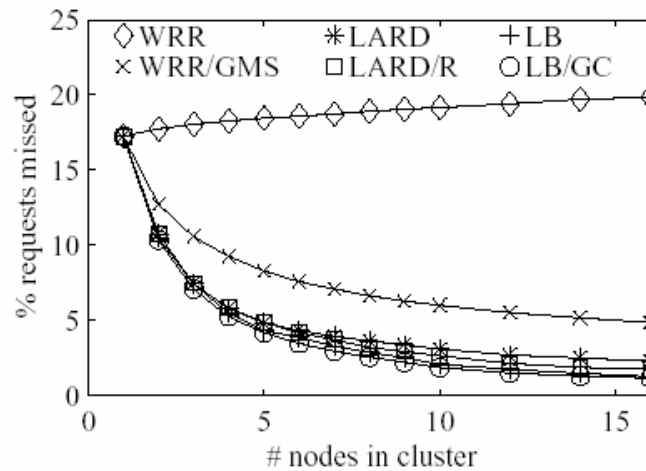
Simulation Results

THROUGHPUT RESULTS



Simulation Results

MISS RATE

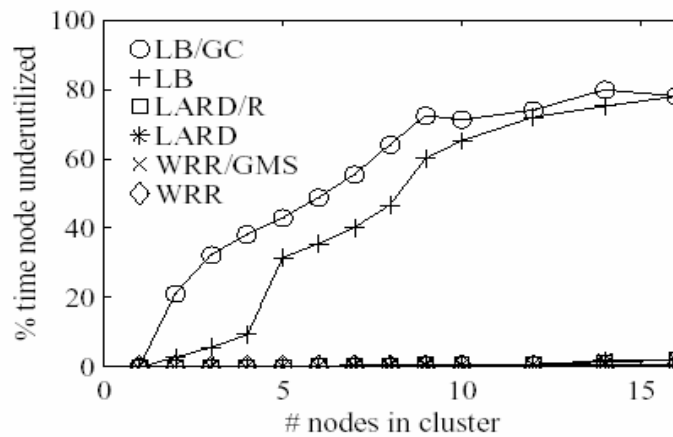


PURDUE
UNIVERSITY

37

Simulation Results

LOAD BALANCING

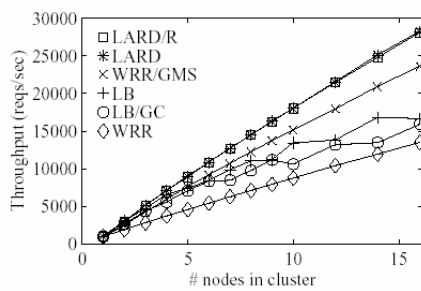


PURDUE
UNIVERSITY

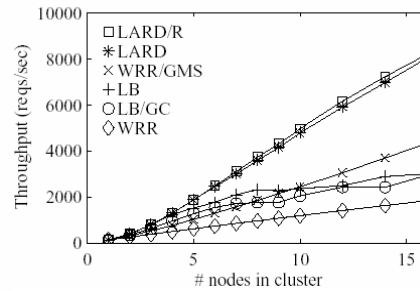
38

Simulation Results

IBM VERSUS RICE WORKLOAD



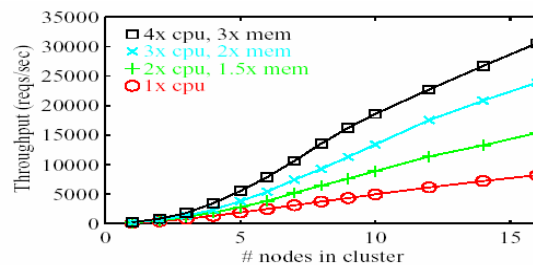
IBM



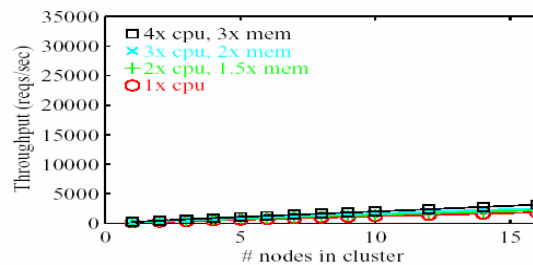
RICE

Simulation Results

SENSITIVITY TO CPU (LARD & WRR)



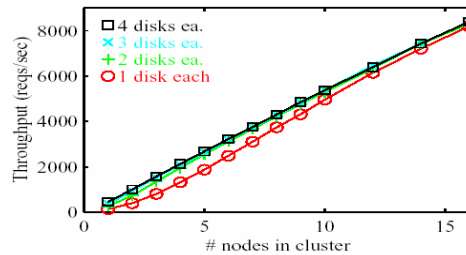
LARD/R



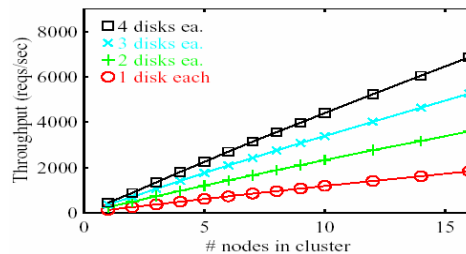
WRR

Simulation Results

SENSITIVITY TO EXTRA DISKS (LARD & WRR)



LARD/R



WRR

PURDUE
UNIVERSITY

41

Experimental Results

ENVIRONMENT

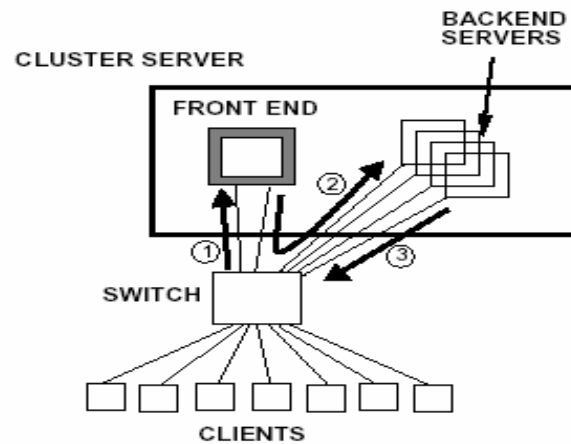
- Pentium II 300 Mhz 128 MB
- 100 Mbps Fast Ethernet
- FreeBSD 2.2.5
- Apache 1.2.4
- Rice Trace

PURDUE
UNIVERSITY

42

Experimental Results

TESTBED

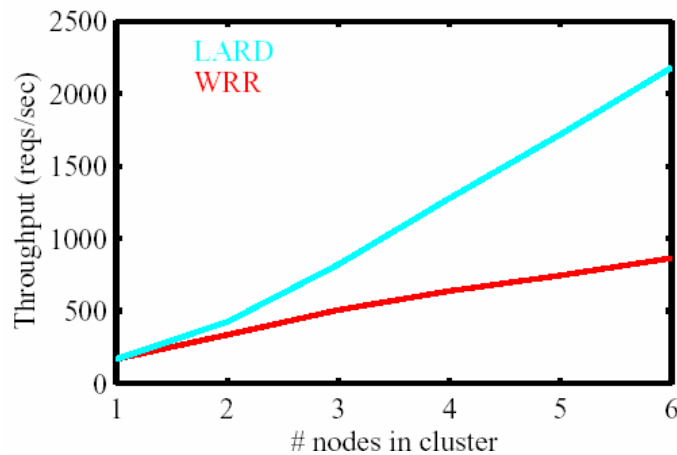


PURDUE
UNIVERSITY

43

Experimental Results

THROUGHPUT



PURDUE
UNIVERSITY

44

Issues

HOW & WHAT TO IMPROVE?

WHAT

- If back-end serving target x fails ?
 - Replicate (in the paper)
- If load on all nodes rises to $2 \cdot T_{\text{high}}$?
 - Admission Control (in the paper)
- **If the front-end is overwhelmed?**
 - **Distributed Distributor (follow up paper)**
- **If the workload is primarily dynamic content?**
 - **LARD for Dynamic Content (attempt by students)**
- If acting as a proxy is too much overhead for the front-end?
 - TCP Handoff not Splicing (in the paper)
- **If connections are persistent**
 - **P-HTTP support [ARON99] (follow up paper)**

Optimizations

DISTRIBUTED DISTRIBUTOR – TWO SCHEMES

Why we need to do this.

We need this

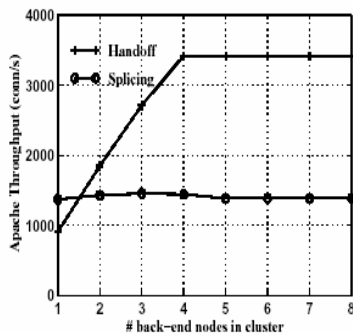


Figure 2: Throughput, 6 KB requests

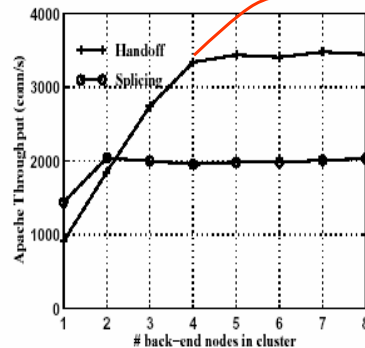
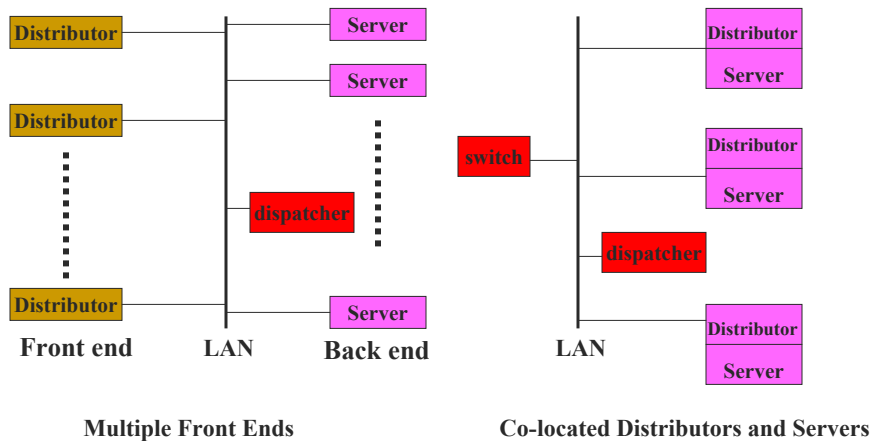


Figure 4: Throughput, IBM trace

Optimizations

DISTRIBUTED DISTRIBUTOR – TWO SCHEMES



Optimizations

DISTRIBUTED DISTRIBUTOR

Multiple front ends mechanism:

1. Poor load balancing
2. Efficient partitioning of cluster nodes into either front end or back end nodes depends upon the workload (unpredictable), which may lead to low cluster utilization.

Co-located distributors and servers:

1. Cluster utilization is high & distributor bottleneck eliminated.
2. Layer 4 switch employed so scalability increased.
3. Centralized switch balances the load on each distributor

Overall disadvantages: A little latency incurred.

Overall Advantages: Scalability is very high

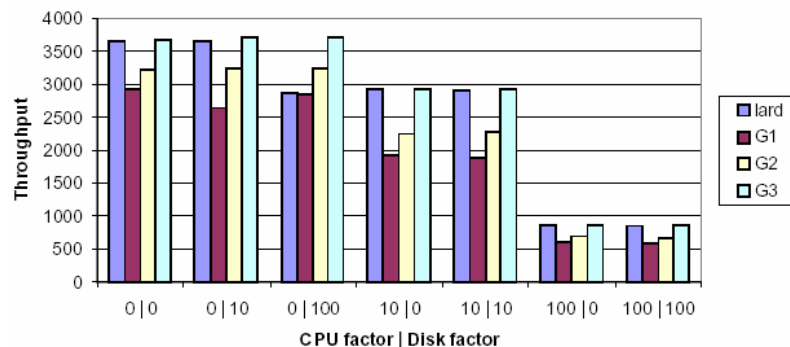
Optimizations

DYNAMIC CONTENT

- Serving dynamic content requires
 - More CPU time and disk bandwidth
 - More complex
- What can be done?
 - $DLoad = W_{cpu} * CPU_{utilization} + W_{disk} * Disk_{utilization}$
 - Reassign target documents when
 - node is overloaded (connections > $2T_{high}$)
 - node is overloaded and \exists node w. connections < T_{low}
 - node DLoad is high and “under-dloaded” node exists
- Shown to perform better than LARD

Optimizations

DYNAMIC CONTENT



Optimizations

PERSISTENT CONNECTIONS

- Problem :
 - LARD operates at the granularity of TCP connections
 - With HTTP/1.1, multiple HTTP requests may arrive on a single TCP connection
 - Result (2 problems)
 - All requests arriving on a given connection must be served by a single back-end node
 - A back-end node may receive requests that it cannot serve if it is a specialized node

Optimizations

PERSISTENT CONNECTIONS - SOLUTIONS

- Redesign handoff to support HTTP/1.1 by allowing the front-end to migrate a connection between back-end nodes
- New handoff mechanism : back-end request forwarding
 - front-end hands off client connections to an appropriate back-end node using the TCP single handoff protocol
 - If request arrives on a p-connection that can/should not be served by the current back-end A, the connection is *not* handed off to another back-end node B
 - front-end informs A which backend (B) should serve the request.
 - A then requests the content directly from B, and forwards the response to the client on its client connection

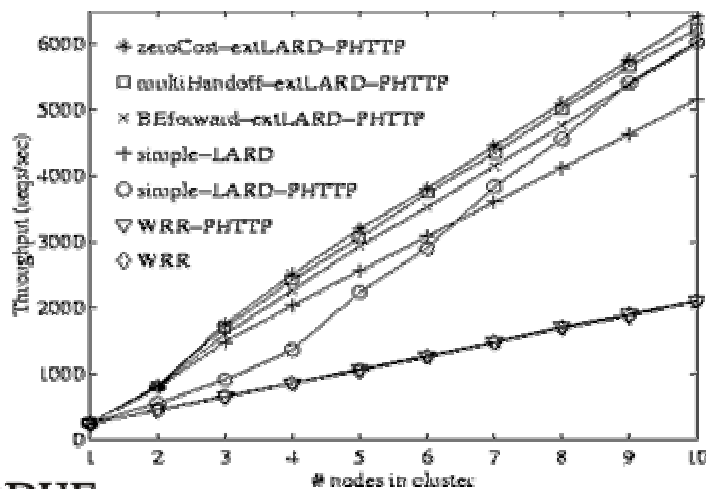
Optimizations

PERSISTENT CONNECTIONS – EXTENDING LARD

- Factors to consider
 - Choice of assignment constrained by previous choice
 - Reassignment overhead should be considered
- Extended LARD
 - First request on a p-connection assigned backend using original LARD policy
 - For subsequent requests
 - If the target is cached at the connection handling node or if the disk utilization on the connection handling node is low (less than 5 queued disk events), then the request is assigned to the same
 - Else, the three cost metrics (locality, load, replacement) are computed over the connection handling node and any other back-end nodes that have the target cached. The request is then assigned to the node that yields the minimum aggregate cost

Optimizations

PERSISTENT CONNECTIONS – RESULTS



Conclusions

PERFORMANCE EVALUATION OVERVIEW

- LARD paper compares SLB/WRR and LB with LARD approaches:
 - simulation study
 - small Rice and IBM web server logs
 - tweak simulation parameters to achieve desired result?
 - Nodes have small memories with *greedy-dual* replacement.
 - WRR combined with global cache-sharing among servers (GMS).
 - WRR/GMS is global cache LRU with duplicates and cache-sharing cost.
 - LB/GC is global cache LRU with duplicate suppression and no cache-sharing cost.

Conclusions

PERFORMANCE EVALUATION RESULTS

- 1. WRR has the lowest cache hit ratios and the lowest throughput.
 - There is much to be gained by improving cache effectiveness.
- 2. LB achieve slightly better cache hit ratios than LARD.
 - WRR/GMS lags behind because of duplicates.
- 3. The caching benefit of LB is minimal, and LB is almost as good as LB/GC.
 - Locality- request distribution induces good cache behavior at the back ends: global cache replacement adds little.
- 4. Better load balancing in the LARD strategies dominates the caching benefits of LB.
 - LARD/R and LARD achieve the best throughput and scalability; LARD/R yields slightly better throughput.

Conclusions

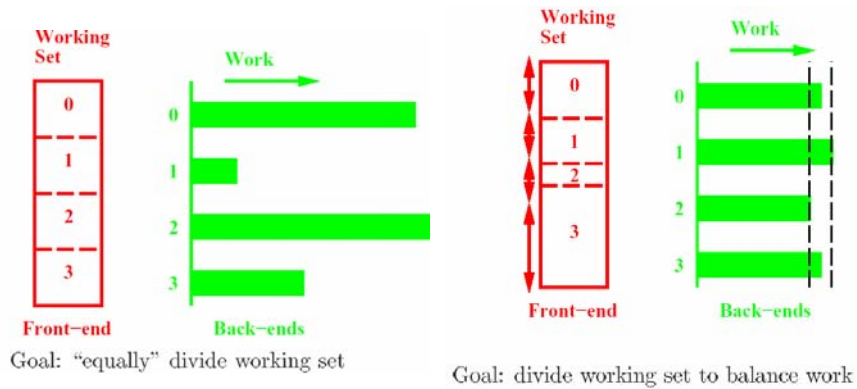
ISSUES AND QUESTIONS

- 1. LB (URL switching) has great cache behavior but lousy throughput.
 - Why? Underutilized time results show poor load balancing.
- 2. WRR/GMS has good cache behavior and great load balancing, but not-so-great throughput.
 - Why? How sensitive is it to CPU speed and network speed?
- 3. What is the impact of front-end caching?
- 4. What is the effectiveness of bucketed URL hashing policies?
 - E.g., Foundry: hash URL to a bucket, pick server for bucket based on load.
- 5. Why don't L7 switch products support LARD? Should they?
 - [USENIX 2000]: use L4 front end; back ends do LARD handoff.

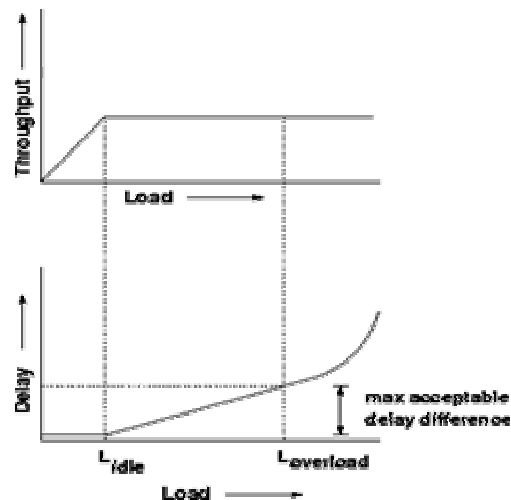
Questions & Ambiguous Responses



[LARD motivation]



[LARD Intuition]



]



31