

Engineering a Content Delivery Network

Bruce Maggs

Network Deployment



20000+
Servers

1200+
Networks

72+
Countries

● Current Installations



Part I: Services

<http://www.yahoo.com>

<http://www.amazon.com>

<http://windowsupdate.microsoft.com>

<http://www.apple.com/quicktime/whatson>

<http://www.fbi.gov>

Design Themes

- Redundancy
- Self-assessment
- Fail-over at multiple levels
- Robust algorithms

FirstPoint – DNS (e.g., Yahoo!)

- Selects from among several mirror sites operated by content provider



Embedded Image Delivery (e.g., Amazon)

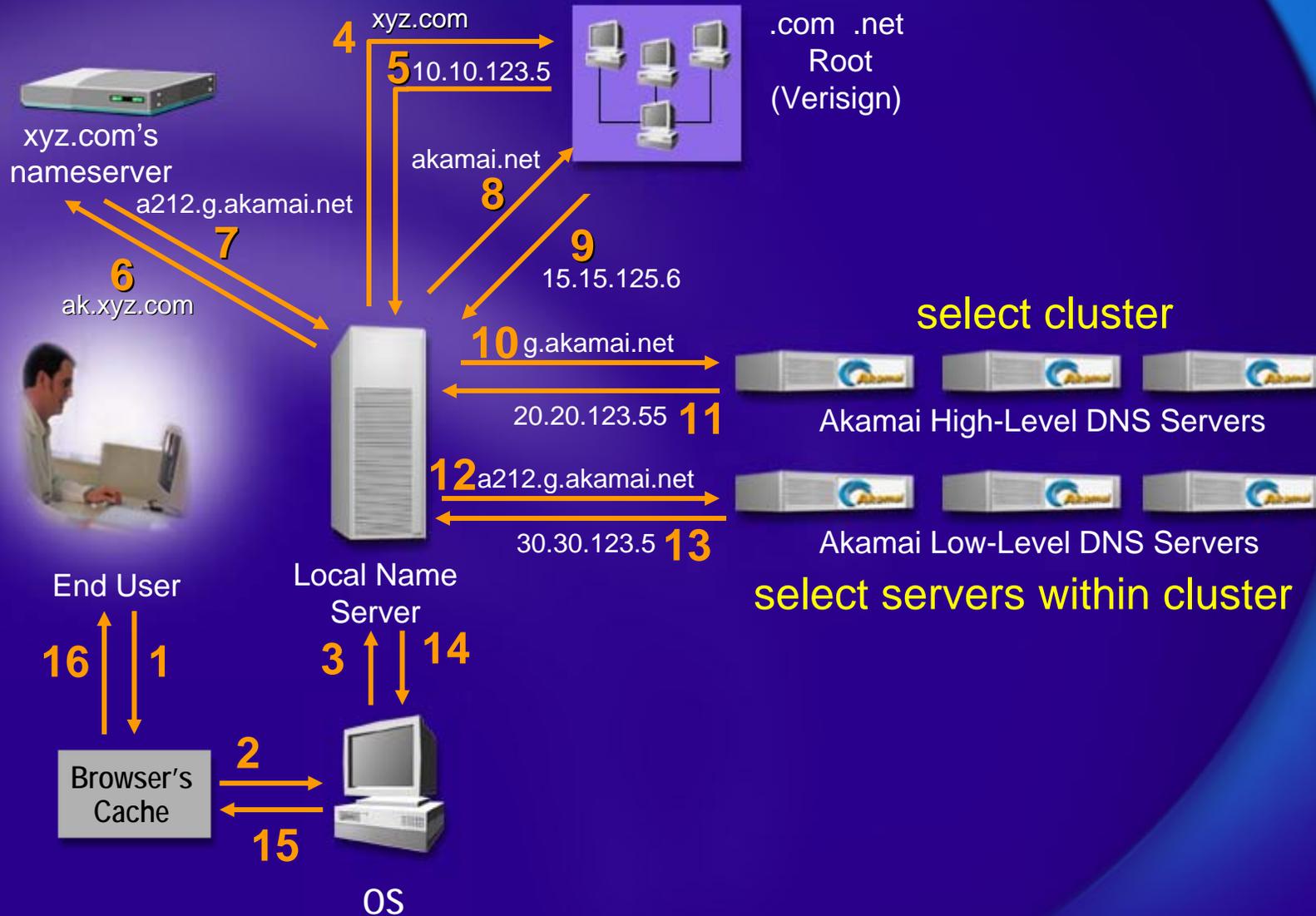
Embedded URLs are Converted to ARLs

```
<html>
<head>
<title>Welcome to xyz.com!</title>
</head>
<body>
  
  
  <h1>Welcome to our Web site!</h1>
  <a href="page2.html">Click here to enter</a>
</body>
</html>
```

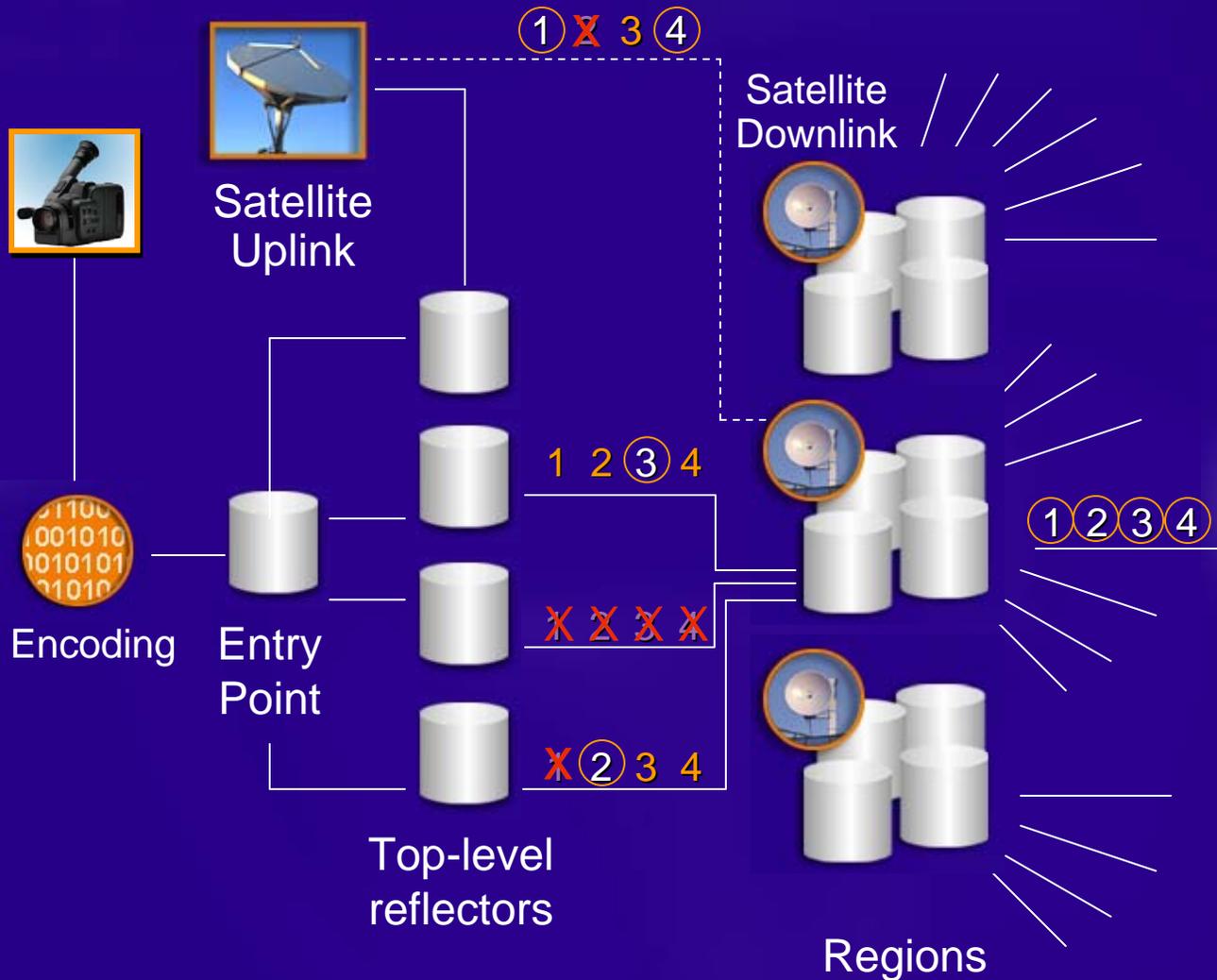
ak



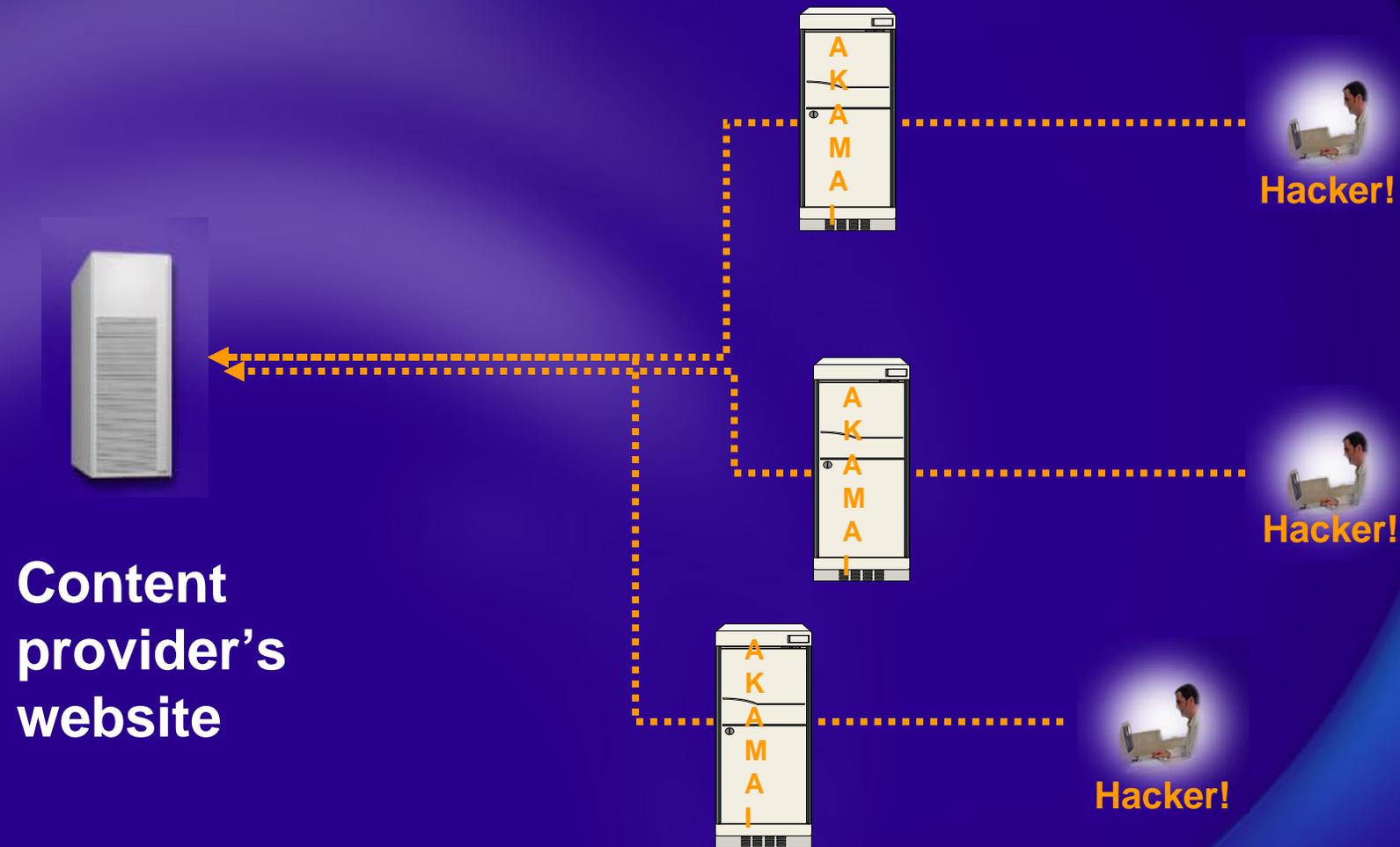
Akamai DNS Resolution



Live Streaming Architecture



SiteShield (www.fbi.gov)



Part II: Failures

- 1. Hardware**
- 2. Network**
- 3. Software**
- 4. Configuration**
- 5. Misperceptions**
- 6. Attacks**

Hardware / Server Failures



Linux boxes with large
RAM and disk capacity,
Windows servers

Sample Failures:

1. Memory SIMMS jumping out of their sockets
2. Network cards screwed down but not in slot
3. Etc.

Akamai Cluster



Servers pool resources

- RAM
- Disk
- Throughput



View of Clusters

Cluster Information

Region 16 : UU-PA 9 %

204.178.110.73	
204.178.110.33	17
204.178.110.34	17
204.178.110.35	0
204.178.110.36	0
204.178.110.46	9
204.178.110.47	9
204.178.110.48	7
204.178.110.49	9
204.178.110.65	7
204.178.110.66	7
204.178.110.75	8
204.178.110.76	13

Region 7 : EX-MA 0 %

209.67.231.142	
209.67.231.134	1
209.67.231.136	0
209.67.231.137	0
209.67.231.166	0
209.67.231.167	0
209.67.231.168	1
209.67.231.169	0
209.67.231.198	1
209.67.231.199	0
209.67.231.200	0
209.67.231.201	0

Region 15 : UU-VAa 13 %

204.178.107.226	
204.178.107.233	17
204.178.107.234	18
204.178.107.235	13
204.178.107.236	16
204.178.110.3	11
204.178.110.4	16
204.178.110.9	0
204.178.110.11	10
204.178.110.12	13
204.178.123.33	0
204.178.123.34	0
204.178.123.35	13
204.178.123.36	16
204.178.123.46	10
204.178.123.47	13
204.178.123.48	13
204.178.123.49	15

Region 44 : UU-SJ 15 %

204.178.118.65	16
204.178.118.66	20
204.178.118.67	21
204.178.118.68	10

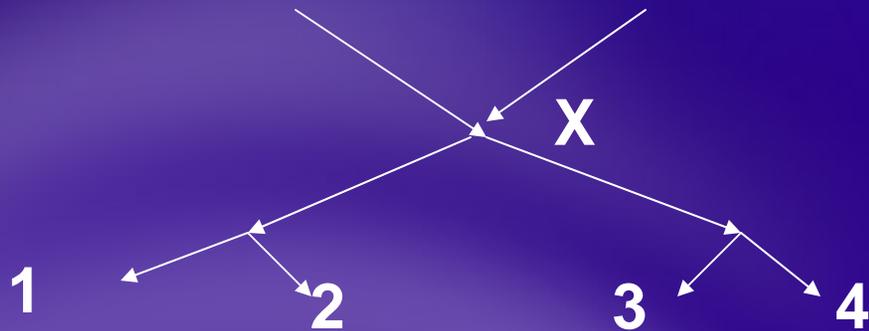
Annotations:

- buddy (points to 204.178.110.33)
- suspended (points to 204.178.110.35)
- hardware failure (points to 204.178.110.36)
- suspended datacenter (points to 209.67.231.136)
- odd man out (points to 204.178.110.9)

Network Failures

E.g., congestion at public and private peering points, misconfigured routers, inaccessible networks, etc., etc., etc.

Core Points



- Core point X is the first router at which all paths to nameservers 1, 2, 3, and 4 intersect.
- X can be viewed as the straddling the core and the edge of the network.

Core Points

500,000 nameservers
reduced to

90,000 core points

7,000 account for 95% end-user load

Engineering Methodology

- C programming language (gcc).
- Reliance on open-source code.
- Large distributed testing systems.
- Burn-in on “invisible” system.
- Staged rollout to production.
- Backwards compatibility.

Perceived Failures

Examples

1. Personal firewalls
2. Reporting tools
3. Customer-side problems
4. Third-party measurements

Cascading Failures

MTU adjustment problem in Linux 2.0.38 kernel

Linux 2.0.38 crashes when TCP connection forces it to reduce MTU to approximately 570 bytes.

Someone in Malaysia configured a router to use this value as its MTU.

Client connecting through the router caused a cascade of Akamai servers to fail.

Attacks

- 8Gb/s attack inflicted on Akamai customer, October 2003
- Attack on Akamai FirstPoint DNS system, July 2004

Lost in Space

The most worrisome “attack” we faced:

One of our servers started receiving properly authenticated control messages from an unknown host.

Fortunately, the messages were not formatted correctly and were discarded by our server.

After two days of investigation, we discovered that the “attacker” was an old server we had lost track of, trying to rejoin the system.

It had been sending these messages for months before we noticed!