CS 54001-1: Large-Scale Networked Systems

Professor: Ian Foster TAs: Xuehai Zhang, Yong Zhao

Lecture 3

www.classes.cs.uchicago.edu/classes/archive/2003/winter/54001-1

Reading

- 1. Internet design principles and protocols
 - P&D 1
 - "End-to-End Principles in Systems Design"
- 2. Internetworking, transport, routing

- P&D 4.1, 4.2, 4.3, 5.1, 5.2, 6.3

- 3. Mapping the Internet and other networks
 - Reading to be defined tomorrow.
- 4. Security
 - P&D 8

Week 1:

Internet Design Principles & Protocols

- An introduction to the mail system
- An introduction to the Internet
- Internet design principles and layering
- Brief history of the Internet
- Packet switching and circuit switching
- Protocols
- Addressing and routing
- Performance metrics
- A detailed FTP example

Week 2: Routing and Transport

- Routing techniques
 - Flooding
 - Distributed Bellman Ford Algorithm
 - Dijkstra's Shortest Path First Algorithm
- Routing in the Internet
 - Hierarchy and Autonomous Systems
 - Interior Routing Protocols: RIP, OSPF
 - Exterior Routing Protocol: BGP
- Transport: achieving reliability
- Transport: achieving fair sharing of links

Recap: An Introduction to the Internet



5

Characteristics of the Internet

- Each packet is individually routed
- No time guarantee for delivery
- No guarantee of delivery in sequence
- No guarantee of delivery at all!
 - Things get lost
 - Acknowledgements
 - Retransmission
 - > How to determine when to retransmit? Timeout?
 - > Need local copies of contents of each packet.
 - > How long to keep each copy?
 - > What if an acknowledgement is lost?
- CS 54001-1 Winter Quarter

Characteristics of the Internet (2)

- No guarantee of integrity of data.
- Packets can be fragmented.
- Packets may be duplicated.

Routing in the Internet

- The Internet uses hierarchical routing
- Internet is split into Autonomous Systems (ASs)
 - Examples of ASs: Stanford (32), HP (71), MCI
 Worldcom (17373)
 - Try: whois -h whois.arin.net ASN "MCI Worldcom"
- Within an AS, the administrator chooses an Interior Gateway Protocol (IGP)
 - Examples of IGPs: RIP (rfc 1058), OSPF (rfc 1247).
- Between ASs, the Internet uses an Exterior
 Gateway Protocol
 - ASs today use the Border Gateway Protocol, BGP-4 (rfc 1771)
- CS 54001-1 Winter Quarter

TCP Characteristics

- TCP is connection-oriented
 - 3-way handshake used for connection setup
- TCP provides a stream-of-bytes service
- TCP is reliable:
 - Acknowledgements indicate delivery of data
 - Checksums are used to detect corrupted data
 - Sequence numbers detect missing, or mis-sequenced data
 - Corrupted data is retransmitted after a timeout
 - Mis-sequenced data is re-sequenced
 - (Window-based) Flow control prevents over-run of receiver
- TCP uses congestion control to share network capacity among users
- CS 54001-1 Winter Quarter

Week 3:

Measurement & Characterization

- What does the Internet look like?
- What does Internet traffic look like?
- How do I measure such things?
- How do such characteristics evolve?
- What Internet characteristics are shared with other networks?
- Are all those Faloutsos' related?

This Week's Reading (Phew)

- Growth of the Internet
 - How fast are supply and demand growing?
- On Power-Law Relationships of the Internet Topology
 - What is the structure of the Internet?
- Experimental Study of Internet Stability and Wide-Area Backbone Failures
 - How reliable is the Internet?
- Graph structure in the Web
 - Another area in which interesting structures arise
- Emergence of Scaling in Random Networks
 - Why do power-law structures arise?
- Search in Power-law Networks
 - How can we exploit this structure in useful way?

This Week's Reading (Phew)

- Growth of the Internet
 - How fast are supply and demand growing?
- On Power-Law Relationships of the Internet Topology
 - What is the structure of the Internet?
- Experimental Study of Internet Stability and Wide-Area Backbone Failures
 - How reliable is the Internet?
- Graph structure in the Web
 - Another area in which interesting structures arise
- Emergence of Scaling in Random Networks
 - Why do power-law structures arise?
- Search in Power-law Networks
 - How can we exploit this structure in useful way?

From year-end 1997 to year-end 2001 (U.S. only)

- long distance fiber deployment: fiber miles growth of 5x
- transmission capacity: DWDM advances of 100x
- cumulative fiber capacity growth of around 500x
- actual demand growth: around 4x

Two fundamental mistakes:

- (i) assume astronomical rate of growth for Internet traffic
- (ii) extrapolate that rate to the entire network

Bandwidth and Growth Rate of U.S. Long Distance Networks, year-end 1997



Traffic on Internet backbones in U.S.

For each year, shows estimated traffic in terabytes during December of that year.

<u>Year</u>	TB/month
1990	1.0
1991	2.0
1992	4.4
1993	8.3
1994	16.3
1995	?
1996	1,500
1997	2,500 - 4,000
1998	5,000 - 8,000
1999	10,000 - 16,000
2000	20,000 - 35,000
2001	40,000 - 70,000
2002	80,000 - 140,000

Distribution of Internet costs: almost all at edges

U.S. Internet connectivity market (excluding residential, web hosting, . . .) \approx \$15 billion/year

U.S. backbone traffic: ≈ 100,000 TB/month

Current transit costs (at OC3 bandwidth): ≈ \$150/Mbps

Hence, if utilize purchased transit at 30% of capacity, cost for total U.S. backbone traffic: \approx \$2 billion/year

Backbones are comparatively inexpensive and will stay that way!

Residential broadband costs:

DSL and cable modem users: average data flow around 10Kb/s per user

If provide 20 Kb/s per user, at current costs for backbone transit of \$150 per Mb/s per month, each user will cost around \$3/month for Internet connectivity.

Most of the cost at edges, backbone transport almost negligible

"Moore's Law" for data traffic:

Usual pattern of large, well-connected institutions: approximate doubling of traffic each year

Note: Some large institutions report growth rates of 30-40% per year, the historical pre-Internet data traffic growth rate

SWITCH traffic and capacity across the Atlantic



Traffic between the University of Minnesota and the Internet



20

The dominant and seriously misleading view of data network utilization



Typical enterprise traffic profile: Demolishes myth of insatiable demand for bandwidth and many (implicit) assumptions about nature of traffic



Weekly traffic profile on an AboveNet OC192 link from Washington, DC to New York City:



Traffic Growth Rate: Key Constraint

• Adoption rates of new services.

"Internet time" is a myth.

New technologies still take on the order of a decade to diffuse widely.

Multimedia file transfers a large portion of current traffic, streaming traffic in the noise

Internet traffic at the University of Wisconsin in Madison



Conclusion:

- Internet traffic is growing vigorously
- Internet bubble caused largely by unrealistic expectations, formed in willful ignorance of existing data
- Main function of data networks: low transaction latency
- QoS likely to see limited use
- File transfers, not streaming multimedia traffic, to dominate

This Week's Reading (Phew)

- Growth of the Internet
 - How fast are supply and demand growing?
- On Power-Law Relationships of the Internet Topology
 - What is the structure of the Internet?
- Experimental Study of Internet Stability and Wide-Area Backbone Failures
 - How reliable is the Internet?
- Graph structure in the Web
 - Another area in which interesting structures arise
- Emergence of Scaling in Random Networks
 - Why do power-law structures arise?
- Search in Power-law Networks
 - How can we exploit this structure in useful way?

What Does the Internet Look Like?

- Like the telephone network?
 - Topology: big telephone companies know their telephone networks
 - Traffic: voice phone connections were quickly identified as Poisson/exponential
- But for the Internet
 - Topology: changes are highly decentralized and dynamic. No-one knows the network!
 - Traffic: computers do most of the talking;
 data connections are not Poisson/exponential

Why Is Topology Important?

- Design efficient protocols
- Create accurate model for simulation
- Derive estimates for topological parameters
- Study fault tolerance and anti-attack properties

Two Levels of Internet Topology

Router Level and AS Level



Random Graph Erdös-Rényi model (1960)



Connect with probability p

$$p=1/6$$

N=10
 $\langle k \rangle \sim 1.5$



Pál Erdös (1913-1996) Poisson distribution



- Int-11-97: the inter-domain topology of the Internet in November of 1997 with 3015 nodes, 5156 edges, and 3.42 avg. outdegree.
- Int-04-98: the inter-domain topology of the Internet in April of 1998 with 3530 nodes, 6432 edges, and 3.65 avg. outdegree.
- Int-12-98: the inter-domain topology of the Internet in December of 1998 with 4389 nodes, 8256 edges, and 3.76 avg. outdegree.



Power-Laws 1 (Faloutsos et al.)

Power-Law 1 (rank exponent) The outdegree, d_v , of a node v, is proportional to the rank of the node, r_v , to the power of a constant, \mathcal{R} :

$d_v \propto r_v^{\mathcal{R}}$



Rank Plots



35

Rank Plots


Power-Law 2

Power-Law 2 (outdegree exponent)

The frequency, f_d , of an outdegree, d, is proportional to the outdegree to the power of a constant, O:

$f_d \propto d^{\mathcal{O}}$

Outdegree Plots



Outdegree Plots



Self Similarity

Distributions of packets/unit ook alike in ifferent time cale



CS 54001-1 Winter Quarter

This Week's Reading (Phew)

- Growth of the Internet
 - How fast are supply and demand growing?
- On Power-Law Relationships of the Internet Topology
 - What is the structure of the Internet?
- Experimental Study of Internet Stability and Wide-Area Backbone Failures
 - How reliable is the Internet?
- Graph structure in the Web
 - Another area in which interesting structures arise
- Emergence of Scaling in Random Networks
 - Why do power-law structures arise?
- Search in Power-law Networks
 - How can we exploit this structure in useful way?

Introduction

Earlier study reveals:

99% routing instability consisted of pathological update, not reflect actual network topological or policy changes.

Causes: hardware, software bugs.

Improved a lot in last several years.

This paper study:

"legitimate" faults that reflect actual link or network failures.

Experimental Methodology

Inter-domain BGP data collection (01/98~11/98)

RouteView probe: participate in remote BGP peering session. Collected 9GB complete routing tables of 3 major ISPs in US.

About 55,000 route entries



Figure 2: Map of major U.S. Internet exchange points.

Intra-domain routing data collection (11/97~11/98) **Case study**:

Medium size regional network --- MichNet Backbone.

Contains 33 backbone routers with several hundred customer routers.

Data from:

u A centralized network management station (CNMS) log data

u Ping every router interfaces every 10 minutes.

u Used to study frequency and duration of failures.

- u Network Operations Center (NOC) log data.
 - u CNMS alerts lasting more than several minutes.
 - u Prolonged degradation of QoS to customer sites.

u Used to study network failure category.

Analysis of Inter-domain Path Stability

BGP routing table events classes:

Route Failure:

loss of a previously available routing table path to a given network or a less specific prefix destination.

Question: Why "less specific prefix"?

Router aggregates multiple more specific prefix into a single supernet advertisement.

128.119.85.0/24 ® 128.119.0.0/16

Route Repair:

A previously failed route to a network prefix is announced as reachable.

Route Fail-over:

A route is implicitly withdrawn and replaced by an alternative route with different next-hop or ASpath to a prefix destination.

Policy Fluctuation:

A route is implicitly withdrawn and replaced by an alternative route with different attributes, but the same next-hop and ASpath.

Pathological Routing:

Repeated withdrawals, or duplicate announcements, for the exact same route.

Last two events have been studied before, here we study the first three events in BGP experiments.

Inter-domain Route Availability

Route availability: A path to a network prefix or a less specific prefix is presented in the provider's routing table.



Figure 4: Cumulative distribution of the route availability of 3 ISPs

Observations from Route Availability Data

- Less than 25%~35% of routes had availability higher than 99.99%
- 10% of routes exhibited under 95% availability
- Internet is far less robust than telephony:
 Public Switched Telephone Network (PSTN)
 averaged an availability rate better than
 99.999%
- The ISP1 step curve represents the 11/98
 major internet failure which caused several hours loss of connectivity

Route Failure and Fail-over

Failure: loss of previously available routing table path to a prefix or less specific prefix destination.

Fail-over: change in ASpath or next-hop reachability of a route.



Fig5: Cumulative distribution of mean-time to failure and mean-time to fail-over for routes from 3 ISPs.

Observation from route failure and fail-over

- **u** The majority routes(>50%) exhibit a mean-time to failure of 15 days.
- **u** 75% routes have failed at least once in 30 days.

- **u** Majority routes fail-over within 2 days.
- u Only 5%~20% of routes do not fail-over within 5 days.
- **u** A slightly higher incidence of failure today than 1994.

Route Repair Time & Failure Duration

Route Repair: A previously failed route is announced reachable. MTTR: Mean-time to Repair



Fig6: Cumulative distribution of MTTR and failure duration for routes from 3 ISPs.

Observation in MTTR and Failure duration

- u 40% failures are repaired in 10 minutes.
- u Majority (70%) are resolved within 1/2 hour.
- u Heavy-tailed distribution of MTTR: failures not repaired in 1/2 hour are serious outage requiring great effort to deal with.
- u Only 25%~35% outages are repaired within 1 hour.
- u Indication: A small number of routes failed many times, for more than one hour.
- u This agrees with previous results that a small fraction pf routes are responsible for majority of network instability.

Analysis of Intra-domain Network Stability

Backbone router: connect to other backbone router via multiple physical path. Well equipped and maintained.

Customer router: connect to regional backbone via single physical connection. Less ideal maintained.



Observation in MTTR and Failure duration

- Majority interfaces exhibit MTTF 40 days. (while majority inter-domain MTTF occur within 30 days)
- Step discontinuities is because a router has many interfaces.
- 80% of all failures are resolved within 2 hours.
- Heavy-tail distribution of MTTR show that
 longer than 2 hours outages are long-term
 and requires great effort to deal with

Frequency Property Analysis

Frequency analysis of BGP and OSPF update messages.



Fig8: BGP updates measured at Mae-East exchange point(08/96~09/96) ; OSPF updates in MichNet using hourly aggregates.(10/98~11/98) CS 54001-1 Winter Quarter

Observation of update frequency

u BGP shows significant frequencies at 7 days, and 24 hours.

u Low amount instability in weekends.

u Fairly stable of Internet in early morning compared with North American business hours.

u Absence of intra-domain frequency pattern indicates that much of BGP instability stems from Internet congestion.

u BGP is build on TCP. TCP has congestion window. Update or KeepAlive message time out.

u AS Internal congestion make IBGP lost and spread out.

u Some new routers provide a mechanism: BGP traffic has higher priority and KeepAlive message persist under congestion.

Conclusions

- Internet exhibit significantly less availability and reliability than telephony network.
- Major Internet backbone paths exhibit mean-time to failure of 25 days or less, mean-time to repair of 20 minutes or less. Internet backbones are rerouted(either due to failure or policy changes) on average of once every 3 days or less
- The 24 hours, 7 days cycle of BGP traffic and none cycle in OSPF suggest that BGP instability stem from congestion collapse.
- A small number of Internet ASes contribute to a large number of long-term outage and backbone unavailability.
- CS 54001-1 Winter Quarter

This Week's Reading (Phew)

- Growth of the Internet
 - How fast are supply and demand growing?
- On Power-Law Relationships of the Internet Topology
 - What is the structure of the Internet?
- Experimental Study of Internet Stability and Wide-Area Backbone Failures
 - How reliable is the Internet?
- Graph structure in the Web
 - Another area in which interesting structures arise
- Emergence of Scaling in Random Networks
 - Why do power-law structures arise?
- Search in Power-law Networks
 - How can we exploit this structure in useful way?



Erdös-Rényi model (1960)



Connect with probability p p=1/6

N=10 $\langle k \rangle \sim 1.5$



Pál Erdös (1913-1996)

Poisson distribution





- Democratic
- Random

Small Worlds

- Stanley Milgram 's experiment
- Small Worlds by Watts/Strogatz
- $\gamma(v) =$ Clustering coefficient of node v
 - = Percentage of neighbours of v connected to each other
- Clustering coefficient:

 $=\frac{\sum_{v\in V} (v)}{|V|}$

Cluster Coefficient

Clustering: My friends will likely know each other!



Probability to be connected C >> p

 $\frac{\text{\# of links between 1,2,...n neighbors}}{n(n-1)/2}$

Networks are clustered [large C(p)] but have a small characteristic path length [small L(p)].

Network	С	C _{rand}	L	N
WWW	0.1078	0.00023	3.1	153127
Internet	0.18-0.3	0.001	3.7-3.76	3015- 6209
Actor	0.79	0.00027	3.65	225226
Coauthorship	0.43	0.00018	5.9	52909
Metabolic	0.32	0.026	2.9	282
Foodweb	0.22	0.06	2.43	134
C. elegance	0.28	0.05	2.65	282

Watts-Strogatz Model



What did we expect?



CS 54001-1 Winter Quarter

19 degrees of separation



• Finite size scaling: create a network with N nodes with $P_{in}(k)$ and $P_{out}(k)$

< l > = 0.35 + 2.06 log(N)



Power-law Distributions

Gnutella: Node connectivity follows a powerlaw*, i.e. P(k neighbours) ~ $k^{-\gamma}$



* Mapping the Gnutella network: Properties of largescale peer-to-peer systems and implications for system design. M. Ripeanu, A. Iamnitchi, and I. Foster. IEEE Internet Computing Journal 6, 1 (2002), 50-57.

What does it mean?



INTERNET BACKBONE

Nodes: computers, routers Links: physical lines



(Faloutsos, Faloutsos and Faloutsos, 1999)



ACTOR CONNECTIVITIES

Nodes: actors Links: cast jointly





Days of Thunder (1990) Far and Away (1992) Eyes Wide Shut (1999)



 10^{3}



SCIENCE CITATION INDEX

1,000 Most Cited Physicists, 1981-June 1997

Out of over 500,000 Examined

(see http://www.sst.nrel.gov)



* citation total may be skewed because of multiple authors with the same name

SCIENCE COAUTHORSHIP

Nodes: scientist (authors) Links: write paper together


Food Web

Nodes: trophic species Links: trophic interactions





Most real world networks have the same internal structure:

Scale-free networks

Why?

What does it mean?

SCALE-FREE NETWORKS

(1) The number of nodes (N) is NOT fixed.

Networks continuously expand by the addition of new nodes

Examples: WWW : addition of new documents Citation : publication of new papers

(2) The attachment is NOT uniform.

A node is linked with higher probability to a node that already has a large number of links.

Examples : WWW : new documents link to well known sites (CNN, YAHOO, NewYork Times, etc) Citation : well cited papers are more likely to be cited again

Scale-free model

 $\Pi(k_i) = \frac{\kappa_i}{\sum_i k_i}$

(1) **GROWTH** :

At every timestep we add a new node with *m* edges (connected to the nodes already present in the system).

The probability that a new node will be connected to

(2) PREFERENTIAL ATTACHMENT :

node *i* depends on the connectivity k_i of that node



CS 54001-1 Winter Quarter A.-L.Barabási, R. Albert, Science 286, 509 (1999)7

Achilles' Heel of complex network



R. Albert, H. Jeong, A.L. Barabasi, Nature 406 378 (2000)

CS 54001-1 Winter Quarter

78

This Week's Reading (Phew)

- Growth of the Internet
 - How fast are supply and demand growing?
- On Power-Law Relationships of the Internet Topology
 - What is the structure of the Internet?
- Experimental Study of Internet Stability and Wide-Area Backbone Failures
 - How reliable is the Internet?
- Graph structure in the Web
 - Another area in which interesting structures arise
- Emergence of Scaling in Random Networks
 - Why do power-law structures arise?
- Search in Power-law Networks
 - How can we exploit this structure in useful way?

What Does the Web Really Look Like?

- Graph Structure in the Web, Broder et al.
- Analysis of 2 Altavista crawls, each with over 200M pages and 1.5 billion links

Confirm Power Law Structure



But Things Are More Complex Than One Might Think ...



CS 54001-1 winter Quarter

82



Course Outline (Subject to Change)

- 1. (January 9th) Internet design principles and protocols
- 2. (January 16th) Internetworking, transport, routing
- 3. (January 23rd) Mapping the Internet and other networks
- 4. (January 30th) Security
- 5. (February 6th) P2P technologies & applications (Matei Ripeanu) (plus midterm)
- 6. (February 13th) Optical networks (Charlie Catlett)
- 7. *(February 20th) Web and Grid Services (Steve Tuecke)
- 8. (February 27th) Network operations (Greg Jackson)
- *(March 6th) Advanced applications (with guest lecturers: Terry Disz, Mike Wilde)
- 10. (March 13th) Final exam
 - * Ian Foster is out of town.