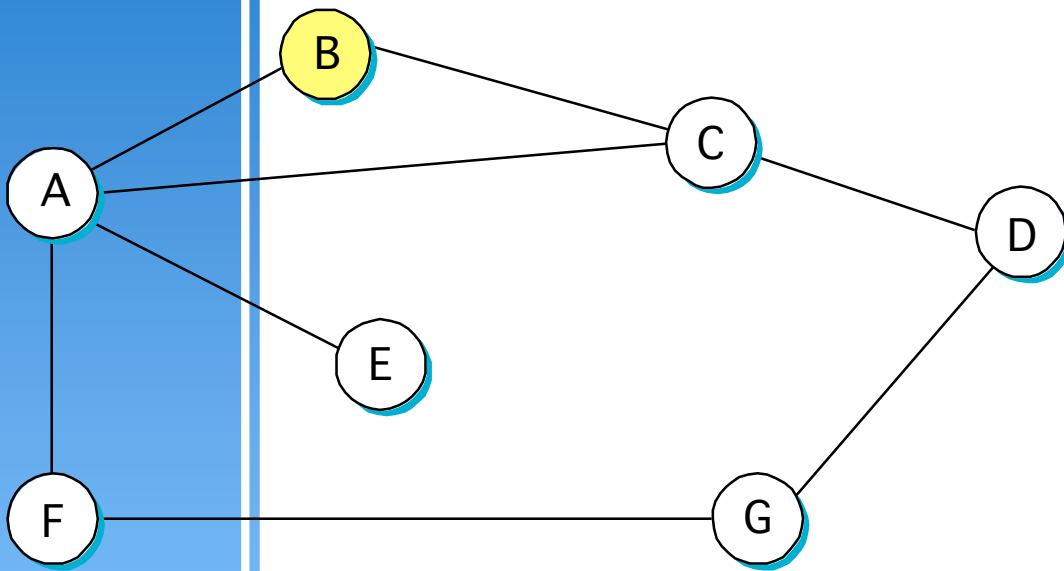


Outline

- ▶ We are focusing on routing algorithms in the last lecture
- ▶ We will look at subnetting and how routing protocols are applied to the larger network hierarchy



Example



Destination	Cost	NextHop
A	1	A
C	1	C
D	2	C
E	2	A
F	2	A
G	3	A



Routing Loops

► Example 1

- F detects that link to G has failed
- F sets distance to G to infinity and sends update to A
- A sets distance to G to infinity since it uses F to reach G
- A receives periodic update from C with 2-hop path to G
- A sets distance to G to 3 and sends update to F
- F decides it can reach G in 4 hops via A

► Example 2: count to infinity problem

- link from A to E fails
- A advertises distance of infinity to E
- B and C advertise a distance of 2 to E
- B decides it can reach E in 3 hops; advertises this to A
- A decides it can reach E in 4 hops; advertises this to C
- C decides that it can reach E in 5 hops...



Loop-Breaking Heuristics

- ▶ Set infinity to 16
- ▶ Split horizon: node does not send routing updates back to the neighbor
- ▶ Split horizon with poison reverse: sends negative information back to the neighbor



Link State (e.g. OSPF)

► Strategy

- send to all nodes (not just neighbors) information about directly connected links (not entire routing table)

► Link State Packet (LSP)

- id of the node that created the LSP
- cost of link to each directly connected neighbor
- sequence number (SEQNO)
- time-to-live (TTL) for this packet



Link State (cont)

► Reliable flooding

- store most recent LSP from each node
- forward LSP to all nodes but one that sent it
- generate new LSP periodically
 - increment SEQNO
- start SEQNO at 0 when reboot
- decrement TTL of each stored LSP
 - discard when TTL=0



Route Calculation

- ▶ Dijkstra's shortest path algorithm
- ▶ Let
 - N denotes set of nodes in the graph
 - $I(i, j)$ denotes non-negative cost (weight) for edge (i, j)
 - s denotes this node
 - M denotes the set of nodes incorporated so far
 - $C(n)$ denotes cost of the path from s to node n

$M = \{s\}$

for each n in $N - \{s\}$

$C(n) = I(s, n)$

while $(N \neq M)$

$M = M \cup \{w\}$ such that $C(w)$ is the minimum for all w in $(N - M)$

for each n in $(N - M)$

$C(n) = \min(C(n), C(w) + I(w, n))$



Route cost metrics

- ▶ Original ARPANET metric
 - measures number of packets queued on each link
 - took neither latency or bandwidth into consideration
- ▶ New ARPANET metric
 - stamp each incoming packet with its arrival time (AT)
 - record departure time (DT)
 - when link-level ACK arrives, compute
 - $\text{Delay} = (\text{DT} - \text{AT}) + \text{Transmit} + \text{Latency}$
 - if timeout, reset DT to departure time for retransmission
 - link cost = average delay over some time period
- ▶ Fine Tuning
 - compressed dynamic range
 - replaced Delay with link utilization



Mobility

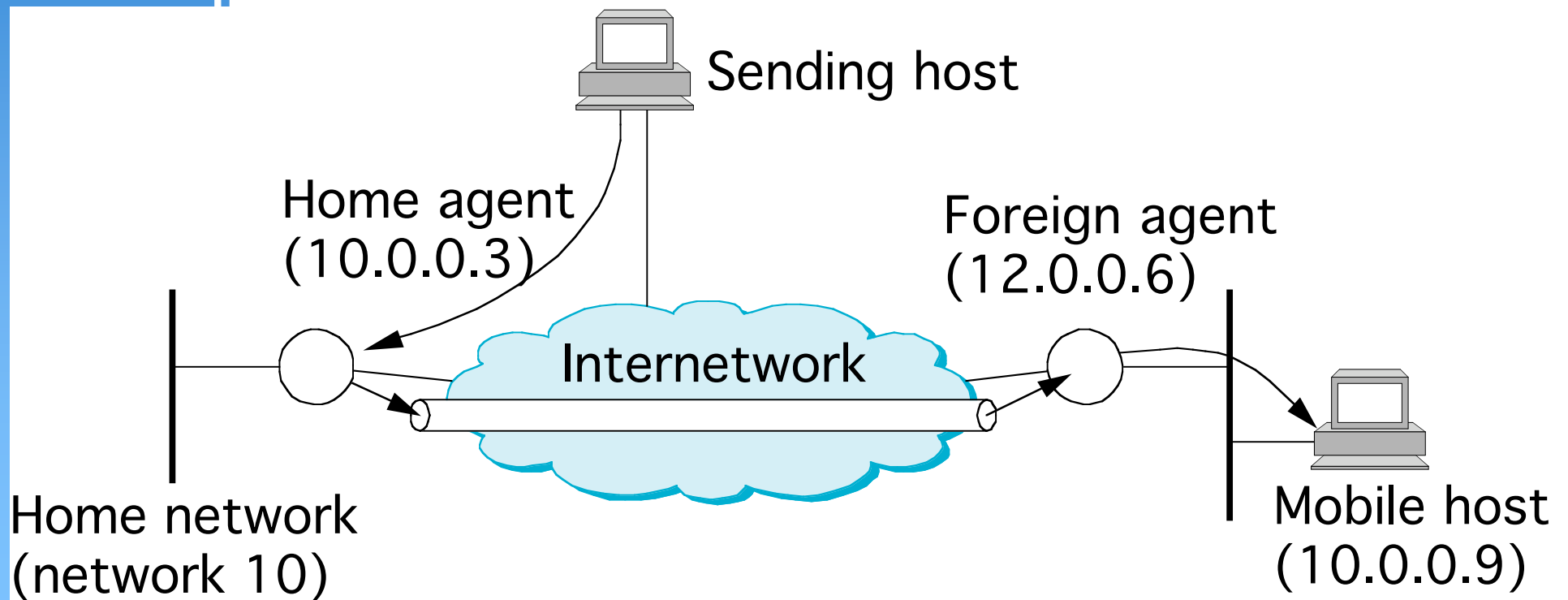
► What if nodes move

- You need a new IP address when you move
- Communications (sockets) have to be reestablished
- One solution is to use Dynamic DNS with DHCP
 - Used at ND
 - When a host moves, DHCP gives it a new address and Dynamic DNS updates the DNS entry with the new DHCP address
 - For example, my laptop is called kural.cse.nd.edu, but may map into different IP addresses depending on where I am
 - Works for new connections, old connections break
 - Can only work within the same domain (because DNS servers are only administered for the domain)

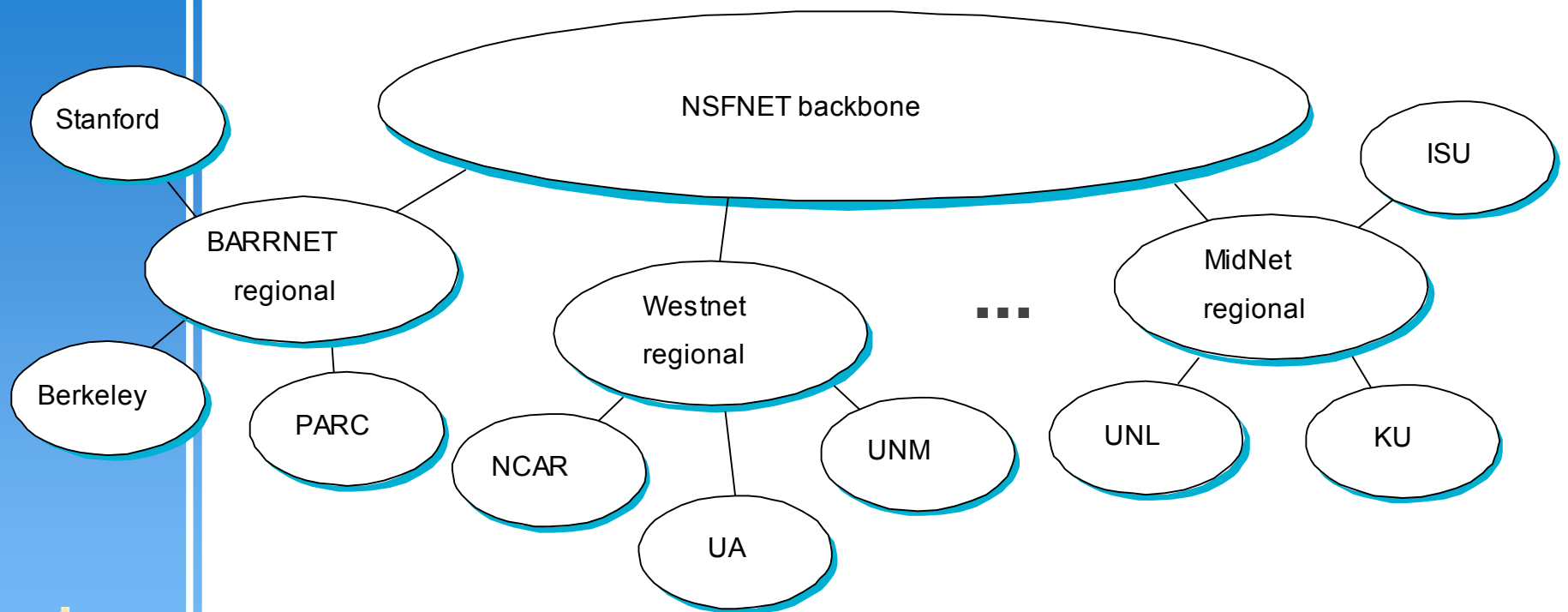


Mobile IP

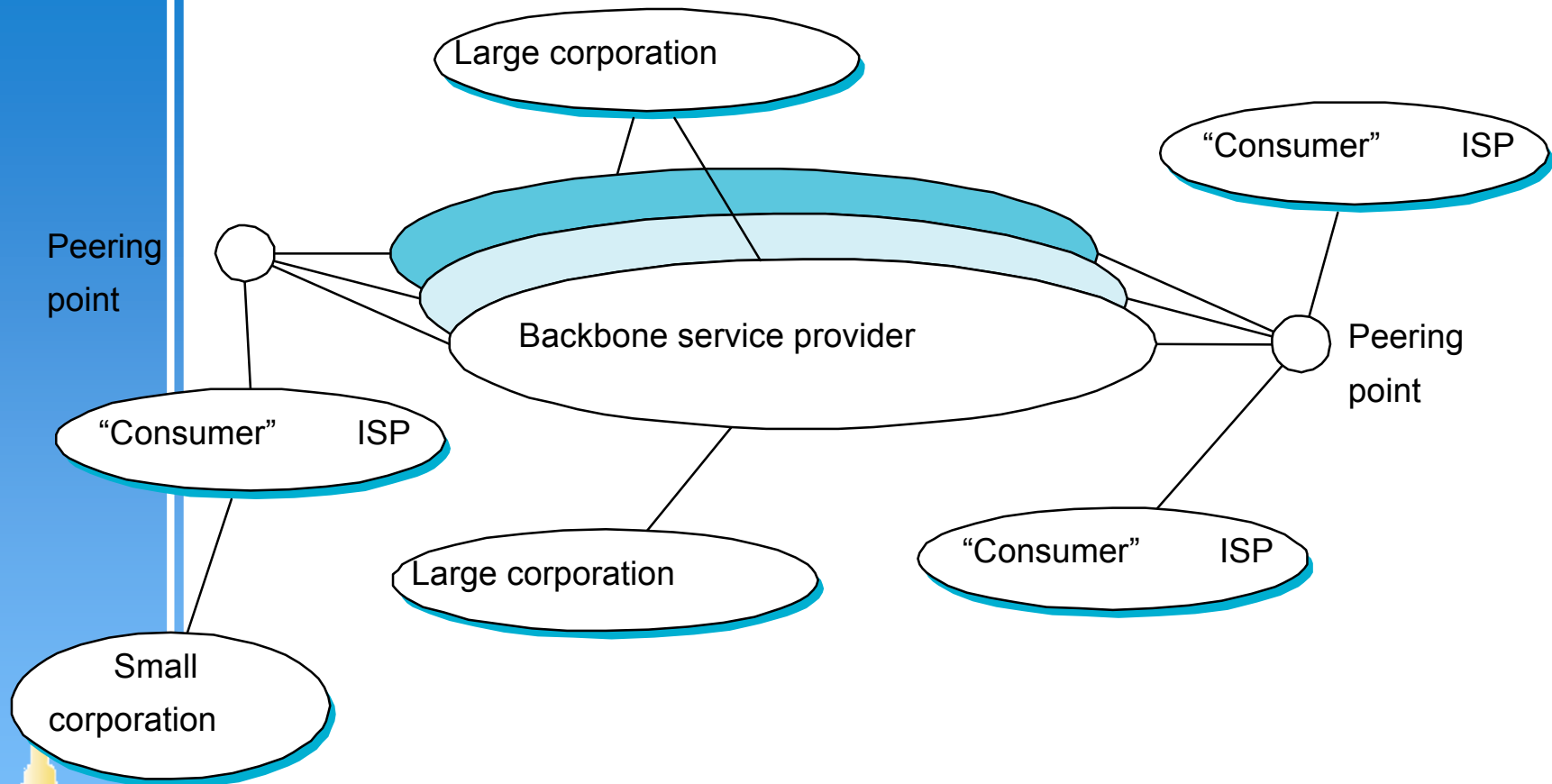
- ▶ Mobile host registers with Foreign Agent. FA informs Home Agent. HA tunnels packets to FA. Communications through the Home address.



Internet Structure - Past



Internet Structure - Today



Subnetting

- ▶ Add another level to address/routing hierarchy: subnet
- ▶ Subnet masks define variable partition of host part
- ▶ Subnets visible only within site

Network number	Host number
----------------	-------------

Class B address

11111111111111111111111111111111	100000000
----------------------------------	-----------

Subnet mask (255.255.255.0)

Network number	Subnet ID	Host ID
----------------	-----------	---------

Subnetted address



Subnet Example

Subnet mask: 255.255.255.128

Subnet number: 128.96.34.0

128.96.34.15



H1

128.96.34.1

R1

128.96.34.130

Subnet mask: 255.255.255.1

Subnet number: 128.96.34.1

128.96.34.129

R2

128.96.34.139



H2

H3



128.96.33.14

128.96.33.1

Subnet mask: 255.255.255.0

Subnet number: 128.96.33.0

Forwarding table at router R1

Subnet Number	Subnet Mask	Next Hop
128.96.34.0	255.255.255.128	interface 0
128.96.34.128	255.255.255.128	interface 1
128.96.33.0	255.255.255.0	R2



Forwarding Algorithm

```
D = destination IP address
for each entry (SubnetNum, SubnetMask, NextHop)
    D1 = SubnetMask & D
    if D1 = SubnetNum
        if NextHop is an interface
            deliver datagram directly to D
        else
            deliver datagram to NextHop
```

- ▶ Use a default router if nothing matches
- ▶ Not necessary for all 1s in subnet mask to be contiguous
- ▶ Can put multiple subnets on one physical network
- ▶ Subnets not visible from the rest of the Internet



Supernetting

- ▶ Assign block of contiguous network numbers to nearby networks
- ▶ Called CIDR: Classless Inter-Domain Routing
- ▶ Represent blocks with a single pair
`(first_network_address, count)`
- ▶ Restrict block sizes to powers of 2
- ▶ Use a bit mask (CIDR mask) to identify block size
- ▶ All routers must understand CIDR addressing



Route Propagation

- ▶ Know a smarter router
 - hosts know local router
 - local routers know site routers
 - site routers know core router
 - core routers know everything
- ▶ Autonomous System (AS)
 - corresponds to an administrative domain
 - examples: University, company, backbone network
 - assign each AS a 16-bit number
- ▶ Two-level route propagation hierarchy
 - interior gateway protocol (each AS selects its own)
 - exterior gateway protocol (Internet-wide standard)



Popular Interior Gateway Protocols

- ▶ RIP: Route Information Protocol
 - developed for XNS
 - distributed with Unix
 - distance-vector algorithm
 - based on hop-count
- ▶ OSPF: Open Shortest Path First
 - recent Internet standard
 - uses link-state algorithm
 - supports load balancing
 - supports authentication



EGP: Exterior Gateway Protocol

► Overview

- designed for tree-structured Internet
- concerned with reachability, not optimal routes

► Protocol messages

- neighbor acquisition: one router requests that another be its peer; peers exchange reachability information
- neighbor reachability: one router periodically tests if the another is still reachable; exchange HELLO/ACK messages; uses a k-out-of-n rule
- routing updates: peers periodically exchange their routing tables (distance-vector)



BGP-4: Border Gateway Protocol

▶ AS Types

- stub AS: has a single connection to one other AS
 - carries local traffic only
- multihomed AS: has connections to more than one AS
 - refuses to carry transit traffic
- transit AS: has connections to more than one AS
 - carries both transit and local traffic

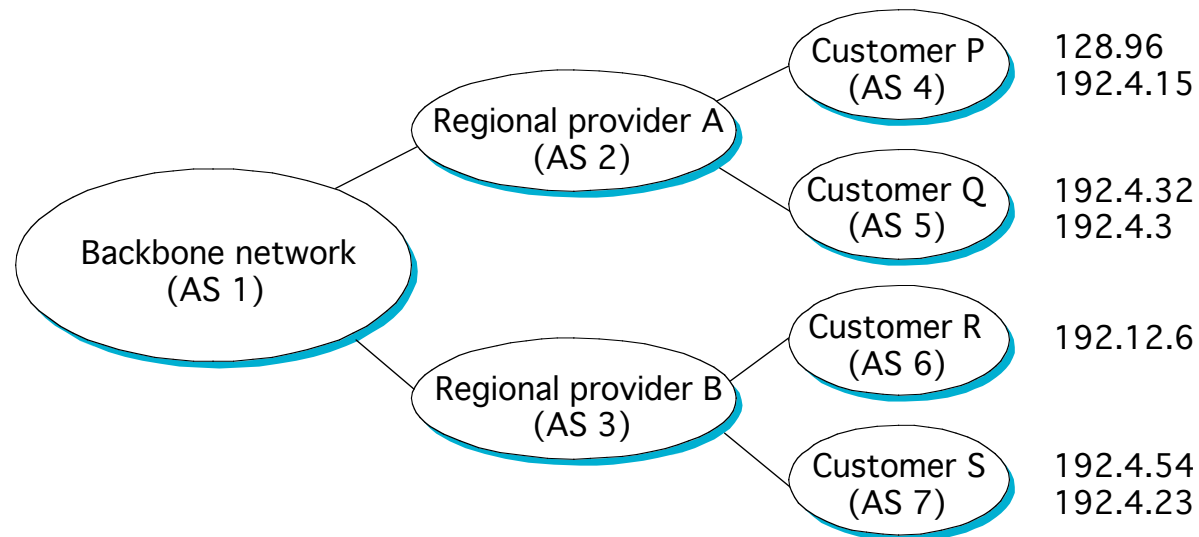
▶ Each AS has:

- one or more border routers
- one BGP speaker that advertises:
 - local networks
 - other reachable networks (transit AS only)
 - gives path information



BGP Example

- ▶ Speaker for AS2 advertises reachability to P and Q
 - network 128.96, 192.4.153, 192.4.32, and 192.4.3, can be reached directly from AS2



- ▶ Speaker for backbone advertises
 - networks 128.96, 192.4.153, 192.4.32, and 192.4.3 can be reached along the path (AS1, AS2).
- ▶ Speaker can cancel previously advertised paths



Peering and Transits

- ▶ Thousands of ISPs. ISPs connect using transit providers and backbone providers to route packets
- ▶ Decisions are made on business goals and \$\$\$
- ▶ Peering does not give access to other peering points, i.e. peering is non-transitive
- ▶ No explicit service level agreement (SLA)
- ▶ Peering can be cheaper
 - For example, Notre Dame can peer with Ameritech and ATT to transfer mutual traffic (from DSL and Cable customers)
 - Lower latency to preferred ISPs



Notre Dame to Saint Marys

► traceroute www.saintmarys.edu

- traceroute to www.saintmarys.edu (147.53.8.10), 30 hops max, 40 byte packets
- 1 eafs-e06.gw.nd.edu (129.74.250.1) 0.664 ms 0.469 ms 0.450 ms
- 2 c245-e01.gw.nd.edu (129.74.245.14) 0.301 ms 0.574 ms 0.345 ms
- 3 monk-fe00.gw.nd.edu (129.74.45.4) 1.046 ms 0.918 ms 0.823 ms
- 4 klimek-i00.gw.nd.edu (129.74.248.102) 4.784 ms 4.569 ms 4.688 ms
- 5 mren-m10-lsd6509.startap.net (206.220.240.86) 4.863 ms 5.884 ms 6.659 ms
- 6 chin-mren-ge.abilene.ucaid.edu (198.32.11.97) 5.234 ms 4.512 ms 4.879 ms
- 7 iplsng-chinng.abilene.ucaid.edu (198.32.8.77) 15.137 ms 22.735 ms 8.524 ms
- 8 ul-abilene.indiana.gigapop.net (192.12.206.250) 8.584 ms 9.009 ms 8.814 ms
- 9 ihets-gw-1-ge15-0.ind.net (157.91.6.37) 8.458 ms 8.581 ms 8.823 ms
- 10 sbn-fa0-0.ind.net (199.8.76.73) 9.256 ms 8.826 ms 8.638 ms
- 11 stmarys-edu-T1.ind.net (199.8.73.110) 30.135 ms 26.131 ms 25.682 ms
- 12 * * smcswitch.saintmarys.edu (147.53.1.1) 31.876 ms !X



Reasons why you don't peer

- ▶ No explicit SLA
- ▶ Use cold-potato algorithm to offset traffic costs
 - Carry traffic in your local network as much as possible rather than use an optimal (possibly more expensive transit route)
 - Transit points use hot potato algorithm, dumping the packets as soon as possible to the back bone (even if it was not optimal)
- ▶ Don't want to help potential competitors
 - Ameritech would want your friends to move to Ameritech so that you all can get faster traffic, not peer with AT&T so that you can enjoy the benefit

