

Rethinking Internet Design

What's changed?

- ❑ **operation in untrustworthy world**
 - endpoints can be malicious
 - If endpoint not trustworthy, but want trustworthy network → more mechanism in network core
- ❑ **more demanding applications**
 - end-to-end best effort service not enough
 - new service models in network (Intserv, diffserv)
 - new application-level service architecture built on top of network core (e.g., CDN, p2p)
 - multiway communication (e.g., teleconferencing, audio/video broadcasting, multiplayer game)

Part 1.7 27

Rethinking Internet Design

What's changed?

- ❑ **ISP service differentiation**
 - ISP doing more (than other ISPs) in core is competitive advantage
- ❑ **Rise of third party involvement**
 - interposed between endpoints (even against will)
 - e.g., US recording industry
- ❑ **less sophisticated users**
 - e.g., thin clients (PDAs, set-top boxes, cell-phones)
 - users lack expertise
 - Historical example: Kerberos 5 pulls inside the programming interface for the function of replay prevention

All five changes motivate shift away from end-end!

Part 1.7 28

What's at stake?

"At issue is the conventional understanding of the 'Internet philosophy'

- freedom of action
- user empowerment
- end-user responsibility for actions taken
- lack of control "in" the net that limit or regulate what users can do

The end-to-end argument fostered that philosophy because it enables the freedom to innovate, install new software at will, and run applications of the users choice"

[Blumenthal and Clark, 2001]

Part 1.7 29

Technical response to changes

- **Trust:** emerging distinction between what is "in" network (us, trusted) and what is "on" (i.e. attached to) the network (them, untrusted).
 - E2e argument → services "in" the network is undesirable
 - E2e argument also guides the placement of services "on" the network
- **Modify endpoints**
 - Harden endpoints against attack
 - Endpoints/routers do content filtering: Net-nanny
 - CDN, ASPs: rise of structured, distributed applications in response to inability to send content (e.g., multimedia, high bw) at high quality

Part 1.7 30

Technical response to changes

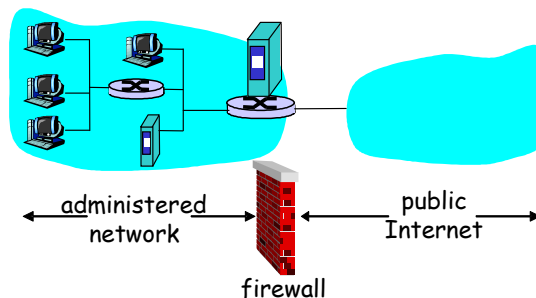
- Add functions to the network core:
 - filtering firewalls
 - application-level firewalls
 - NAT boxes
 - active networking
- ... All operate within network, making use of application-level information
 - which addresses can do what at application level?
 - If addresses have meaning to applications, NAT must "understand" that meaning
 - E.g. arguments for FTP PORT command include an IP address in ASCII! → NAT needs to adjust TCP sequence number!
- ... Need a really compelling case!
 - E.g. multicast never really made it for 20 years

Part 1.7 31

Firewalls

firewall

isolates organization's internal net from larger Internet, allowing some packets to pass, blocking others.



Part 1.7 32

Firewalls: Why

prevent denial of service attacks:

- SYN flooding: attacker establishes many bogus TCP connections, no resources left for "real" connections.

prevent illegal modification/access of internal data.

- e.g., attacker replaces CIA's homepage with something else

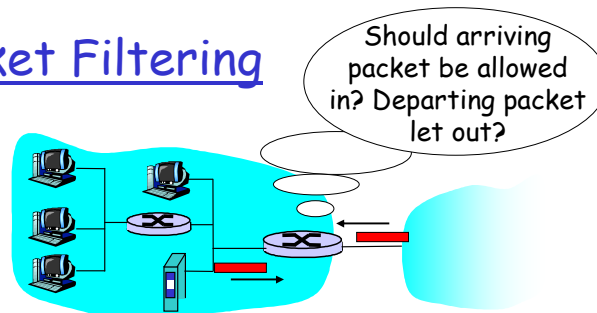
allow only authorized access to inside network (set of authenticated users/hosts)

two types of firewalls:

- packet-filtering
- application-level

Part 1.7 33

Packet Filtering



- internal network connected to Internet via **router firewall**
- router **filters packet-by-packet**, decision to forward/drop packet based on:
 - source IP address, destination IP address
 - TCP/UDP source and destination port numbers
 - ICMP message type
 - TCP SYN and ACK bits

Part 1.7 34

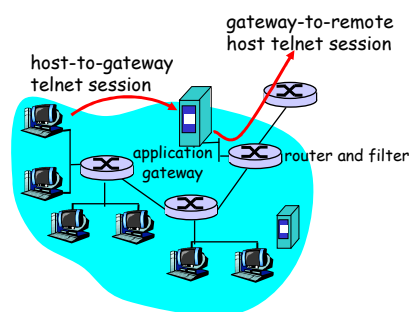
Packet Filtering

- **Example 1: block incoming and outgoing datagrams with IP protocol field = 17 and with either source or dest port = 23.**
 - all incoming and outgoing UDP flows and telnet connections are blocked.
- **Example 2: Block inbound TCP segments with ACK=0.**
 - prevents external clients from making TCP connections with internal clients, but allows internal clients to connect to outside.

Part 1.7 35

Application gateways

- Filters packets on application data as well as on IP/TCP/UDP fields.
- **Example:** allow selected internal users to telnet outside.



1. Require all telnet users to telnet through gateway.
2. For authorized users, gateway sets up telnet connection to dest host. Gateway relays data between 2 connections
3. Router filter blocks all telnet connections not originating from gateway.

Part 1.7 36

NAT: Network Address Translation

- **Motivation:** local network uses just one IP address as far as outside world is concerned:
 - no need to be allocated range of addresses from ISP: just one IP address is used for all devices
 - can change addresses of devices in local network without notifying outside world
 - can change ISP without changing addresses of devices in local network
 - devices inside local net not explicitly addressable, visible by outside world (a security plus).

Part 1.7 37

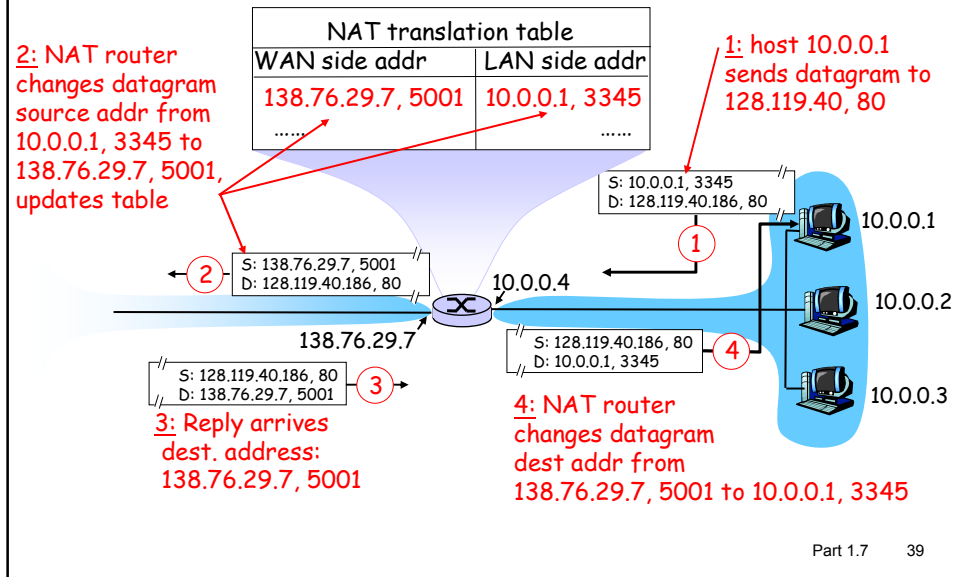
NAT: Network Address Translation

Implementation: NAT router must:

- *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)
... remote clients/servers will respond using (NAT IP address, new port #) as destination addr.
- *remember (in NAT translation table)* every (source IP address, port #) to (NAT IP address, new port #) translation pair
- *incoming datagrams: replace* (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table
- *fix TCP sequence number* if necessary (e.g. for FTP PORT command - "PORT a1,a2,a3,a4,b1,b2")

Part 1.7 38

NAT: Network Address Translation



NAT: Network Address Translation

- ❑ 16-bit port-number field:
 - 60,000 simultaneous connections with a single LAN-side address!
- ❑ NAT is controversial:
 - routers should only process up to layer 3
 - violates end-to-end argument
 - NAT possibility must be taken into account by app designers, eg, P2P applications
 - address shortage should instead be solved by IPv6

Traversal of UDP through NATs

- Q: Suppose H1, H2 are behind NAT box N1 and N2, respectively, how do H1 and H2 talk to each other?
- A: Using STUN (Simple Traversal of UDP through NATs) to learn the NAT type and addr mapping
 1. H1 send request to a public STUN server S1; S1 replies with current mapping at N1: (GA1,GP1)
 2. Similarly H2 learns current mapping at N2: (GA2,GP2)
 3. H1 and H2 predict ports GP1' and GP2' and exchange mapping through out-of-band communication
 4. H1 sends a packet to (GA2,GP2') (which is dropped by N2 but causes N1 to initiate state)
 5. H2 sends a packet to (GA1,GP1') (causing N2 to initiate state)
- Port prediction
 - port-preserving NAT: dp = 0;
 - symmetric NAT: dp = 1, 2 (there is a window of vulnerability)

Part 1.7 41

Traversal of TCP through NATs

- Q: Suppose H1, H2 are behind NAT box N1 and N2, respectively, how do H1 and H2 talk to each other?
- A: port prediction is similar, but how about seq#?
 - STUNT #1: (STUNT = "STUN and TCP too")
 - H1 and H2 send SYN with small TTL to set up state at N1/N2; communicate this info to STUNT server, which then spoof a SYNACK to both ends
 - STUNT #2
 - Only H1 sends SYN with small TTL. After that H1 reset connection and opens passive TCP socket on same port
 - H2 then sends a SYN to (GA1,GP1)
 - NATBlaster
 - H1 and H2 exchange mapping and sequence number out of band, each sends a SYNACK to the other side
- If you are interested, read IMCO5 paper:
 - Characterization and Measurement of TCP Traversal through NATs and Firewalls

Part 1.7 42

What is an Active Network ?

- Depends on who you ask!
- **active services**: application-level services exploiting position within the network to provide enhanced service
 - CDN
 - streaming media caches
- **capsule approach**: packets carry programs, active node executes program when code-carrying packet arrives to active node
 - code may determine what to do with packet
 - may implement other service: e.g., network management, reliable multicast

Part 1.7 43

The capsule approach to active networks

- Network architecture that allows :
 - *application-customized code to be dynamically deployed* in the network
 - customized-code to be *executed* in controlled framework within network
- Similar to extensible operating systems (SPIN, Synthesize etc)
- the new equation:

Packet = Code + data

sort of like postscript

Part 1.7 44

Capsules



- **Type**
 - Identifier for the forwarding routine to be executed (carries code by reference)
- **Previous address**
 - Where to get the forwarding routine from if it is not available in the present node (*Code Distribution*)
- **Dependent Fields**
 - Parameters for the forwarding code
- **Payload**
 - Header + data of higher layers

Part 1.7 45

Active networking and End-to-End Arguments

- **End-to-end principle:** lower layers should have minimum functionality, but support widest variety of applications possible
 - active networking: support *all* higher-level applications
 - minimum common functionality: ability to execute code: *programmable versus pre-programmed* low layer functionality

Part 1.7 46

Active networking: transparency?

- **Transparency:** use of network by others not very visible (can more or less predict behavior of network)
- **Active networking:** transparency difficult
 - need to constrain interactions among programmable entities in router (who knows *what* they will try to do)
 - like OS trying to constrain interaction among processes!

Part 1.7 47

KISS

- success of LAN protocols, RISC architecture: **KISS!**
- "building complex functions into network optimizes network for small number of services, while substantially increasing cost for uses unknown at design time"
- "end-to-end argument does not oppose active networks per se but instead strongly suggests that enthusiasm for the benefits of optimizing current application needs by making the network more complex may be misplaced"

Part 1.7 48

Epilogue: will IP take over the world?

Reasons for success of IP:

- *reachability*: reach every host, adapts topology when links fail.
 - *heterogeneity*: single service abstraction (best effort) regardless of physical link topology
- Note: our grading for these: A-, A

Epilogue: will IP take over the world?

There are many other claimed (or commonly accepted) reasons for IP's success

- IP already dominates global communications
 - IP is more efficient
 - IP is more robust
 - IP is simpler
 - IP supports real-time apps and telephony
- Are these really true?
... let's take a closer look

1. IP already dominates global communications?

□ business revenues:

- ISPs: 13B
 - 59% US households
- Broadcast TV: 29B
- Cable TV: 29.8B
- Radio broadcast: 10.6B
- Phone industry: 268B

Q: IP equipment cheaper?
Economies of scale?
(lots of routers?)

Q: per-device, IP is cheaper
(one line into house, multiple devices)

Q: # bits carried in each network?

Q: Internet, more traffic and congestion
is spread among all users (bad?)

□ Router/telco switch markets:

- Core router: 1.7B; edge routers: 2.4B
- SONET/SDH/WDM: 28B,
Telecom Mobile Satellite Service (MSS): 4.5B

Part 1.7 51

2. IP is more efficient?

□ Statistical multiplexing versus circuit switching

□ Link utilization:

- Avg. link utilization in Internet core: 3% to 30% (ISPs: never run above 50%)
- Avg. utilization of Ethernet is currently 1%
- Avg. link utilization of long distance phone lines: 33%

□ low IP link utilization: purposeful!

- predictability, stability, low delay, resilience to failure
- Recall: *overprovisioning is the market's answer to the challenge of Internet QoS*

□ At low utilization, we *forfeit* benefits of statistical multiplexing!

Part 1.7 52

3. IP is more robust?

- ❑ Median IP network availability: downtime: 471 min/yr
 - (99.9% reliability)
- ❑ Avg. phone network downtime: 5 min/yr
 - (99.999% reliability)
- ❑ Convergence time with link failures:
 - BGP: 3 - 15 minutes
 - but ~ 1 sec with intra-domain, e.g., OSPF (with "proper" timer configurations)
 - SONET: 50 ms
- ❑ **Q:** *why is fast convergence challenging in IP networks?*
- ❑ **A:** many reasons ...
 - distributed route computation
 - in-band signaling (signaling and data share same network)
 - more hops in general → backup path computation is more sophisticated
 - note: congestion may occur after reconvergence

Part 1.7 53

4. IP is simpler?

- ❑ Intelligence at edge, simplicity in core
 - Cisco IOS: 8M lines of code
 - Telephone switch: 3M lines of code
 - **Q:** *where does most of the complexity lie?*
- ❑ Line-card complexity:
 - Router: 30M gates in ASICs, 1 CPU, 300M packet buffers
 - Switch: 25% of gates, no CPU, no packet buffers
- ❑ Key: simple abstraction/interface ≠ simple implementation

5. Support of real-time app's & telephony over IP

- ❑ Not yet

Part 1.7 54

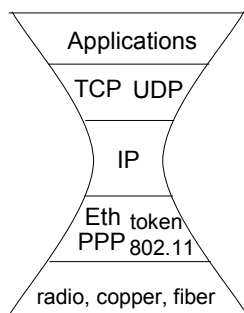
Discussion: benefits of IP?

- IP supports many different types of data applications at a wide range of data rates
 - phone network: 1 of many services (voice, fax, touch-tone service, 800 numbers, teletype, hearing impaired services, lots of enhanced voice services, voicemail...)
- IP traffic, services more diverse (?). IP works at higher bandwidths
 - factually true for end applications, but cores are both high speed (in fact core telephone switches are faster - why?)
- Claim: IP supports short bursty connections "better"
 - Is it really true?
 - implicit: less setup cost, less resources used - not that important given utilization figures
- IP has 1 rtt transaction times, phone network is at least 2 rtt (setup plus transaction)

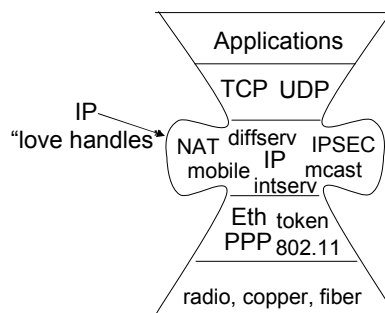
Part 1.7 55

Big picture: supporting new applications - losing the IP hour glass figure?

middle age: a narrowing mind, a widening waist



IP "hourglass"

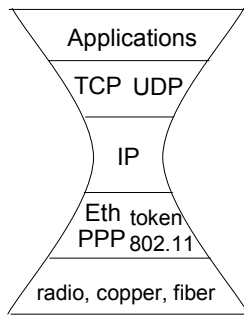


Middle-age IP "hourglass" ?

Part 1.7 56

Big picture: supporting new applications - losing the IP hour glass figure?

middle age: an expanding mind, a slim waist



IP "hourglass"

