Communication Networks

Prof. Laurent Vanbever































DNS resolution can either be recursive or iterative

		_
Records	Name	Value
A	hostname	IP address
NS	domain	DNS server name
MX	domain	Mail server name
CNAME	alias	canonical name
PTR	IP address	corresponding hostname









What Errors Might A Router See?

- · Dead-end: No route to destination
- · Sign of a loop: TTL expires
- · Can't physically forward: packet too big · And has DF flag set
- · Can't keep up with traffic: buffer overflowing
- · Header corruption or ill-formed packets
-

What should network tell host about?

- No route to destination?
- · Host can't detect or fix routing failure.
- TTL expires?
 - Host can't detect or fix routing loop.
- Packet too big (with DF set)?
- · Host can adjust packet size, but can't tell difference between congestion drops and MTU drops • Buffer overflowing?
- · Transport congestion control can detect/deal with this
- · Header corruption or ill-formed packets?
- Host can't fix corruption, but can fix formatting errors

Router Response to Problems?

- · Router doesn't really need to respond
 - · Best effort means never having to say you're sorry
 - · So, IP could conceivably just silently drop packets
- · Network is already trying its best
 - Routing is already trying to avoid loops/dead-ends
 - Network can't reduce packet size (in DF packets)
 - · Network can't reduce load, nor fix format problems
- · What more can/should it do?

Error Reporting Helps Diagnosis

- · Silent failures are really hard to diagnose
- · IP includes feedback mechanism for network problems, so they don't go undetected
- Internet Control Message Protocol (ICMP)
- · The Internet "print" statement
- · Runs on IP, but viewed as integral part of IP

Internet Control Message Protocol

- · Triggered when IP packet encounters a problem • E.g., Time Exceeded or Destination Unreachable
- · ICMP packet sent back to the source IP address Includes the error information (e.g., type and code)
- IP header plus 8+ byte excerpt from original packet
- · Source host receives the ICMP packet • Inspects excerpt (e.g., protocol/ports) to identify socket
- · Exception: not sent if problem packet is ICMP
 - And just for fragment 0 of a group of fragments

Types of Control Messages

- · IP packet too large for link layer, DF set
- - (who generates Port Unreachable?)
- Source Quench

- Tells source to use a different local router

Using ICMP

- · ICMP intended to tell host about network problems
 - Diagnosis
 - Won't say more about this....
- · Can exploit ICMP to elicit network information
 - Discovery
 - Will focus on this....

Discovering Network Path Properties

- PMTU Discovery: Largest packet that can go through the network w/o needing fragmentation · Most efficient size to use
- (Plus fragmentation can amplify loss)
- Traceroute:
 - What is the series of routers that a packet traverses as it travels through the network?
- Pina:
- · Simple RTT measurements

- Need Fragmentation
 - TTL Expired
 - Decremented at each hop; generated if ⇒ 0
 - Unreachable
 - Subtypes: network / host / port

 - · Old-style signal asking sender to slow down
 - Redirect

Ping: Echo and Reply

- · ICMP includes simple "echo" functionality
- Sending node sends an ICMP Echo Request message
- · Receiving node sends an ICMP Echo Reply
- · Ping tool
- · Tests connectivity with a remote host
- ... by sending regularly spaced Echo Request
- ... and measuring delay until receiving replies

Path MTU Discovery

- MTU = Maximum Transmission Unit · Largest IP packet that a link supports
- Path MTU (PMTU) = minimum end-to-end MTU
- Must keep datagrams no larger to avoid fragmentation · How does the sender know the PMTU is?
- Strategy (RFC 1191):
- · Try a desired value
- Set **DF** to prevent fragmentation
- Upon receiving Need Fragmentation ICMP ...
- ... oops, that didn't work, try a smaller value

Issues with Path MTU Discovery

- · What set of values should the sender try?
- Usual strategy: work through "likely suspects"
- E.g., 4352 (FDDI), 1500 (Ethernet),
- 1480 (IP-in-IP over Ethernet), 296 (some modems)
- What if the PMTU changes? (how could it?)
- · Sender will immediately see reductions in PMTU (how?) · Sender can periodically try larger values
- What if Needs Fragmentation ICMP is lost?
- Retransmission will elicit another one
- How can The Whole Thing Fail?
- "PMTU Black Holes": routers that don't send the ICMP

Discovering Routing via Time Exceeded

- Host sends an IP packet
- · Each router decrements the time-to-live field
- If TTL reaches 0
- · Router sends Time Exceeded ICMP back to the source
- · Message identifies router sending it
- Since ICMP is sent using IP, it's just the IP source address
- And can use PTR record to find name of router







Sharing Single Address Across Hosts

- Network Address Translation (NAT) enables many hosts to share a single address
 - Uses port numbers (fields in transport layer)
- · Was thought to be an architectural abomination when first proposed, but it:
 - · Probably saved us from address exhaustion
 - And reflects a modern design paradigm (indirection)

Special-Purpose Address Blocks

- Limited broadcast
 - Sent to every host attached to the local network Block: 255.255.255.255/32
- Loopback
- Address blocks that refer to the local machine
- Address blocks matrices and Block: 127.0.0.0/8 Usually only 127.0.0.1/32 is used
- Link-local

 - By agreement, not forwarded by any router Used for single-link communication only Intent: autoconfiguration (especially when DHCP fails) Block: 169.254.0.0/16
- Private addresses
- By agreement, not routed in the public Internet For networks not meant for general Internet conne Blocks: 10.0.0.0/8, 172.16.0.0/12, 192.168.0.0/16







What is SDN and how does it help?

SDN is a new approach to networking

- Not about "architecture": IP, TCP, etc.
- But about design of network control (routing, TE,...)
- SDN is predicated around two simple concepts

 Separates the control-plane from the data-plane
 Provides open API to directly access the data-plane
- While SDN doesn't do much, it enables a lot



NAT: Early Example of "Middlebox"

- Boxes stuck into network to delivery functionality
 NATs, Firewalls,....
- Don't fit into architecture, violate E2E principle
- But a very handy way to inject functionality that:
 Does not require end host changes or cooperation
 In under constant (a.g. cooperation)
 - Is under operator control (e.g., security)
- An interesting architectural challenge:How to incorporate middleboxes into architecture

Networks are Hard to Manage Operating a network is expensive More than half the cost of a network Yet, operator error causes most outages Buggy software in the equipment Routers with 20+ million lines of code Cascading failures, vulnerabilities, etc. The network is "in the way" Especially a problem in data centers ... and home networks

Rethinking the "Division of Labor"





- Router
- Match: longest destination IP prefix
- Action: forward out a link
- Switch
 - Match: destination MAC address
 - Action: forward or flood

• Firewall

- Match: IP addresses and TCP/UDP port numbers - Action: permit or deny
- NAT
- Match: IP address and
- port
- and port



- Server load balancing
- Network virtualization
- Using multiple wireless access points

SDN/OpenFlow controller

(Un)install rules, v statistics.

Send packets

- Energy-efficient networking
- Adaptive traffic monitoring
- Denial-of-Service attack detection
- le (true): ad event e: ch up: update topology populates switch tabl (Un)install rules, Topology changes, Traffic statistics, Arriving packets Query statistics, Send packets



- Controller and switches
- Rules installed in the switches







- Bring SDN to the Internet
- Enable SDN in existing networks
- Boost the performance of existing networks using SDN
- Verify controller programs and interactions
- Improve network monitoring
- Improve network security and anonymity
- ... and many more!

- SDN is exciting

 Enables innovation
 - Simplifies management
 - Rethinks networking from the ground-up
- Significant momentum
- In both research and industry
 Size of the SDN market already several billion \$\$
- Great research opportunity
 - Practical impact on future networks practices
- Placing network on a strong foundation





...to Google's data-center

Insight

Key concepts and problems in Networking

Naming Layering Routing Reliability Sharing

List any technologies, principles, applications.. used after typing in:

> www.google.ch

and pressing enter in your browser

Skill Build, operate and configure networks



The Internet is organized as layers, providing a set of service provided layer service provided L5 Application network access L4 Transport end-to-end delivery (reliable or not) L3 Network global best-effort delivery L2 Link local best-effort delivery L1 Physical physical transfer of bits



W rc	We started with the fundamentals of routing and reliable transport			
		Application	network access	
	L4	Transport	end-to-end delivery (reliable or not)	
	L3	Network	global best-effort delivery	
		Link	local best-effort delivery	
		Physical	physical transfer of bits	



We saw three ways to compute valid routing state					
	Intuition	Example			
#1	Use tree-like topologies	Spanning-tree			
#2	Rely on a global network view	Link-State SDN			
#3	Rely on distributed computation	Distance-Vector BGP			

We saw how to design a reliable transport protocol		
goals		
correctness	ensure data is delivered, in order, and untouched	
timeliness	minimize time until data is transferred	
efficiency	optimal use of bandwidth	
fairness	play well with other concurrent communications	















Transport Protocols: UDP & TCP The requirements

Data delivering, to the correct application

- · IP just points towards next protocol
- Transport needs to demultiplex incoming data (ports)
- Files or bytestreams abstractions for the applications
- Network deals with packets
- Transport layer needs to translate between them
 Reliable transfer (if needed)
 Not overloading the receiver

Not overloading the network

Transport Protocols: UDP & TCP The implementation

Demultiplexing: identifier for application process

- Going from host-to-host (IP) to process-to-process
- Translating between bytestreams and packets:
- Do segmentation and reassembly
- Reliability: ACKs and all that stuff

Corruption: Checksum

- Not overloading receiver: "Flow Control"
- Limit data in receiver's buffer
- Not overloading network: "Congestion Control"













